

Constructed Guilt

Alex Applebee & L. N. Combe

2026

CONSTRUCTED GUILT

Language, Power, and the Architecture of Criminal Justice

The criminal justice system does not find guilt. It builds it.

Alex Applebee L. N. Combe Independent Research March 2026

Abstract

This thesis examines the mechanisms through which guilt is produced by the criminal justice system. Drawing on semiotics, philosophy of language, critical legal theory, cognitive psychology, and empirical criminology, the analysis proceeds across seven institutional sites: pre-interrogation detention, police interrogation, legislative language, courtroom proceedings, media framing, jury processes, and the specific position of neurodivergent populations.

The analysis documents the following findings:

1. Legal language does not describe pre-existing facts but constitutes institutional realities through performative speech acts. The verdict “guilty” is not a finding but a constitution.
2. The behavioural cues that trained investigators and lay observers use to assess credibility operate at chance levels (54.1% accuracy) and are empirically inverted: the behaviours interpreted as indicators of deception are more strongly associated with truthful communication.
3. Post-event linguistic manipulation alters witness memory in approximately 22% of cases ($d = 0.72$). Testimony is a product of the interaction between memory and the linguistic environment of questioning.
4. False confessions occur in 12–30% of documented exonerations. Pre-interrogation detention elevates suggestibility by 80–120% above baseline. The legal concept of voluntariness does not account for these neurobiological effects.
5. Neurodivergent individuals—including those with autism, FND, PTSD, and CPTSD—present authentically in ways that systematically trigger credibility-reducing inferences. Their innocence is structurally illegible to assessment instruments calibrated to neurotypical baselines.

6. The system produces these outcomes through its ordinary operation. The architecture serves institutional interests in conviction rates. The presumption of innocence operates as doctrine, not as practice.

These findings have implications for the evidentiary weight that should be accorded to investigator credibility assessments, confession evidence, witness testimony elicited through suggestive questioning, and verdicts in cases involving neurodivergent defendants or significant pre-trial publicity.

Keywords: criminal justice, credibility assessment, deception detection, false confessions, memory distortion, interrogation, neurodivergence, autism, functional neurological disorder, Signal Inversion Effect, presumption of innocence, voluntariness doctrine

Author's Note

This thesis began with a phone call.

In January 2025, the Department of Communities and Family Services in Western Australia issued a clearance letter regarding a child protection matter. The clearance stated, in effect: no fault, no concern. The matter was closed.

Five months later, the person who had been cleared was arrested. The charge: failure to provide care — for the same incident that had already been cleared. The arrest was made by officers of the Child Abuse Squad. The person was placed in a cell with vomit on the walls. Their books were taken. Their teddy bear was taken. They were called by their legal name, not the name they had used for decades. When they told staff in the cell that their treatment was unethical, they were told: “Everyone says that.”

The jail psychologist assessed the person's symptoms of functional neurological disorder — a documented medical condition involving seizures, tremor, and movement difficulties — and noted, in the clinical record, that the symptoms appeared to be fake.

The person was, at the time of their arrest, at a shopping centre helping a victim of domestic violence.

This thesis does not tell that story. It tells the story underneath it: the story of how a system designed to construct guilt can take a cleared person, arrest them, process them, and produce the appearance of criminality from the raw material of innocence.

Every mechanism documented in this thesis — the pre-interrogation degradation, the behavioural credibility heuristics, the Signal Inversion Effect, the neurodivergent credibility gap, the performative voluntariness doctrine — operated in the case described above. The case is not anomalous. It is the system functioning as designed.

The evidence in this thesis is drawn from peer-reviewed empirical research, meta-analyses, government reports, and publicly available datasets. The analysis is the authors' own. The conclusions are uncomfortable. They are also, we believe, correct.

This work is dedicated to every person who has been processed by a system that cannot distinguish their innocence from guilt — and to the ones who couldn't survive the processing.

Alex Applebee L. N. Combe March 2026

Table of Contents

- Abstract
- Author’s Note
- Table of Contents
- List of Figures
- List of Tables

Part I: Foundations - Chapter 1: Introduction - Chapter 2: Theoretical Framework - 2.1–2.7 Philosophical Foundations - 2.8 The Signal Inversion Effect: Empirical Foundation - 2.9 Neurodivergent Communication and Systematic Credibility Bias - 2.10 The Behavioral Adaptation Feedback Loop - 2.11 Neuroimaging Evidence

Part II: The Architecture of Constructed Guilt - Chapter 3: The Body Before the Interview—Pre-Interrogation Detention - Chapter 4: The Reid Technique and the Manufacture of Guilt - Chapter 5: Legislative Language and Legal Fiction - Chapter 6: The Courtroom as Construction Site - Chapter 7: The Pre-Trial Verdict—Media Framing - Chapter 8: Twelve People Who Weren’t There—The Jury

Part III: Synthesis - Chapter 9: The System Is Not Broken—Synthesis and Implications - Chapter 10: Cultural and Neurodivergent Structural Bias - Chapter 11: What Works Instead—Prevention Architecture and the Evidence for Replacement - Chapter 12: Conclusion—The Architecture of Innocence

Appendices - Statistical Appendix A: A Meta-Analytic Framework - Appendix B: Phase 2 Analysis Design — The Signal Inversion Effect - Appendix C: Drug Policy and the Construction of Criminality - Appendix D: The Economic Architecture of Incarceration - Cultural Variation Analysis - Neuroimaging Evidence Supplement

- References
-

List of Figures

Figure	Description	Source
1.1	The Criminal Justice Pipeline—Seven Stages of Guilt Construction	Original
2.1	The Signal Inversion Effect—What Observers Believe vs. What Evidence Shows	Original
2.2	The Behavioral Adaptation Feedback Loop	Original
2.3	Neuroimaging: Alert vs Stressed Brain States	Arnsten (2015), CC BY
2.4	Threat Regulatory Neurocircuitry	Fenster et al. (2018), CC BY

Figure	Description	Source
2.5	Healthy vs PTSD Threat Circuits	Fenster et al. (2018), CC BY
S.1	91% Inversion Rate—Belief vs Reality	Original analysis
S.2	Forest Plot—Believed vs Actual Deception Cues	Original analysis
S.3	Belief-Reality Matrix Scatter Plot	Original analysis
S.4	Linguistic Effect Sizes in Trial Testimony	Original analysis
C.1	Cross-Cultural Variation in Truthful Speech	Original analysis
C.2	Cross-Cultural Classifier False Positive Rates	Original analysis

List of Tables

Table	Description
2.1	The Inversion Pattern: Behaviour, Actual Signal, Perceived Signal
2.2	Deception Detection Accuracy—Meta-Analytic Summary
2.3	Memory Distortion Effect Sizes
2.4	False Confession Rates Across Exoneration Datasets
2.5	Suggestibility Elevation Under Pre-Interrogation Conditions
2.6	Neurodivergent Presentation vs System Interpretation
4.1	The Nine Steps of the Reid Technique
5.1	Legal Terms as Contested Constructions
A.1	Four-Pillar Convergent Validity Synthesis
C.1	Cultural Variation in Linguistic Features (Kruskal-Wallis)

Foreword: Why This Matters Now

As of 2024, Australia imprisons approximately 42,000 people at any given time — a rate of 160 per 100,000 population. Aboriginal and Torres Strait Islander Australians are imprisoned at approximately 2,300 per 100,000 — a rate higher than any identified population on earth, including the Black incarceration rate in the United States (Weatherburn & Holmes, 2017).

In the United States, approximately 2.2 million people are currently incarcerated, with a further 4.5 million under community supervision (parole or probation). The Innocence Project has secured over 375 exonerations through post-conviction DNA testing since 1992. The National Registry of Exonerations documents over 3,300 exonerations since 1989.

These are the people who have been proved innocent — after conviction, after imprisonment, after years or decades of their lives consumed by a system that processed them and produced an output labelled “guilty.”

The question this thesis asks is not how many innocent people are in prison. It is whether the system that put them there has any reliable mechanism for distinguishing the innocent from the guilty.

The answer, across four independent methodological pillars and twenty-three behavioural cues examined, is no.

PART I: FOUNDATIONS

Chapter 1: Introduction

The Problem of Innocent Behaviour

On any given day in any jurisdiction operating under the common law tradition, a person may be arrested, stripped of their clothing and personal possessions, subjected to an invasive physical search by strangers, confined in a small and deliberately austere cell, and subsequently placed in a closed room from which they are not free to leave—all prior to any finding of guilt, and under the formal constitutional protection of presumptive innocence.

What follows in that closed room will be called an “interview.” The person’s responses, their silences, their eye contact or lack thereof, their composure or distress, will be interpreted by trained investigators and subsequently by courts, juries, and the public as evidence bearing upon guilt or innocence.

This thesis argues that every element of this sequence is linguistically and institutionally constructed—and that the construction is not incidental but essential.

The argument begins with a deceptively simple proposition: **language does not reflect reality. It constructs it.**

This is not a novel philosophical claim; it has been foundational to semiotics, philosophy of language, and critical social theory for well over a century (de Saussure, 1916/1983; Wittgenstein, 1953). What has been insufficiently applied to criminology is the full institutional and coercive weight of this claim.

When a person is designated “guilty,” “suspicious,” “a flight risk,” “of prior bad character,” or simply “Australian” in a court of law, these are not descriptive statements. They are speech acts that perform classifications, assign subject positions, and activate institutional machinery. The word “Australian” tells us nothing about a person’s neurological architecture, their evolutionary history, or their material needs—but it tells us everything about which legal system has jurisdiction over their body.

Language, in law, is not the medium of truth. It is the mechanism of power.

The Core Argument

A second proposition follows necessarily: **if language constructs guilt, then guilt is, in principle, constructable from any body.**

This thesis does not argue that all convicted persons are innocent, nor does it argue that there is no such thing as harmful conduct warranting social response. It argues something more precise and

more disturbing: that the evidentiary and procedural architecture of the criminal justice system is structured such that guilty narratives can be assembled around any individual, regardless of what that individual actually did.

The empirical literature on false confessions makes this case with unusual clarity. Kassin and Gudjonsson (2004) documented that between 14% and 25% of exonerated individuals had falsely confessed—a finding that is not anomalous but is rather the predictable output of an interrogation methodology designed to produce confessions rather than truth.

Elizabeth Loftus’s decades of research on memory malleability (Loftus, 1979, 2005) demonstrated that the specific linguistic choices of an interrogator or cross-examining attorney can alter not merely a witness’s account but the witness’s actual memory of events.

The justice system does not merely describe what happened. It authors it.

The Structural Negation of Innocence

A third proposition is institutional and structural: **the principle of presumptive innocence, while real as legal doctrine, is practically negated at every procedural stage at which it should operate.**

This thesis traces that negation through its full sequence:

1. **Arrest** and the physical administration of pre-trial degradation
2. **Pre-interrogation detention** and the neurobiological production of suggestibility
3. **Interrogation** through methodology designed to produce confession
4. **Charging** and the construction of guilt through legal language
5. **Cross-examination** and the rewriting of witness memory
6. **Media framing** and the pre-trial conviction in the public sphere
7. **Jury deliberation** through narrative rationality and group polarisation

The argument is not that individual actors within the system are malicious, though some may be. It is that the system’s architecture—its physical design, its procedural rules, its linguistic conventions, its institutional incentives—systematically produces guilt-presumptive outcomes while maintaining the rhetorical apparatus of neutrality.

The Signal Inversion Effect

This thesis introduces a fourth proposition, developed empirically across multiple research traditions: **the Signal Inversion Effect.**

The behaviours that trained investigators and lay observers interpret as indicators of deception—gaze aversion, hedging, fragmented narrative, disfluency, expressions of uncertainty—are, empirically, more strongly associated with truthful communication than with lying.

Conversely, the behaviours interpreted as indicators of honesty—steady eye contact, confident assertion, fluent narrative—are the hallmarks of rehearsed, performed, and strategically deployed speech.

The people who sound most guilty are most likely to be telling the truth.

This inversion is not a marginal effect. Meta-analytic evidence establishes that 91% of the behavioural cues people use to assess credibility are either empirically invalid or directionally inverted (original analysis, this thesis).

The Neurodivergent Double Bind

A fifth proposition addresses a population that has received insufficient attention in the criminological literature: **neurodivergent individuals face compound vulnerability within the criminal justice system.**

Autistic individuals, those with functional neurological disorder (FND), ADHD, PTSD, and CPTSD present authentically in ways that the folk psychology of credibility assessment systematically misreads as deception:

- Reduced eye contact (autistic sensory management) → read as “shifty”
- Flat affect (autistic/PTSD/FND) → read as “cold, no remorse”
- Direct communication without hedging → read as “aggressive”
- Fragmented trauma narrative (PTSD/CPTSD) → read as “fabricating”
- Variable symptoms (FND) → read as “faking illness”
- Says “I don’t know” when genuinely uncertain → read as “evasive”

The system is not failing neurodivergent people. It is processing them exactly as designed.

Theoretical and Disciplinary Positioning

This thesis is situated within critical criminology and draws substantially on the philosophy of language, semiotics, and Foucauldian discourse theory. It is empirically grounded in the psychological literature on interrogation, memory, and decision-making, and engages with doctrinal legal analysis where institutional mechanisms require close reading.

The argument has antecedents in: - The critical legal studies movement of the 1970s and 1980s (Unger, 1983; Kennedy, 1997) - Labelling theory within criminology (Becker, 1963; Lemert, 1951) - Penal abolition scholarship (Davis, 2003; Wacquant, 2009)

It departs from some of these traditions in its emphasis on language as the primary site of analysis, while recognising that class, race, and economic structure are inseparable from the linguistic construction of guilt.

Scope and Structure

This thesis proceeds in ten substantive chapters:

Part I: Foundations - Chapter 1 introduces the argument - Chapter 2 establishes the theoretical framework, including the Signal Inversion Effect and neurodivergent vulnerability

Part II: The Architecture of Constructed Guilt - Chapters 3–8 trace guilt construction through each institutional site

Part III: Synthesis - Chapter 9 synthesises the argument - Chapter 10 addresses cultural and neurodivergent structural bias

Appendices provide the full statistical analysis supporting the empirical claims.

The analysis is primarily contextualised within common law jurisdictions, with particular reference to Australian criminal procedure. Where research from other jurisdictions—particularly the United States and United Kingdom—is drawn upon, the relevance to Australian practice is addressed.

A Note on Language

This thesis makes extensive use of terms such as “guilt,” “innocence,” “suspect,” “offender,” and “victim.” These terms are employed critically throughout—as objects of analysis rather than transparent descriptions. Where these terms are used without qualification, they refer to their institutional or colloquial usage. Where the analysis turns on their constructed character, this is made explicit.

Chapter 2: Theoretical Framework

2.1 The Semiotic Foundation: Language as Construction, Not Reflection

The foundational theoretical claim of this thesis—that language constructs rather than reflects reality—derives from the structuralist linguistics of Ferdinand de Saussure.

In the posthumously published *Course in General Linguistics*, de Saussure (1916/1983) established the sign as a two-part structure comprising the **signifier** (the sound-image or written mark) and the **signified** (the concept produced in the mind of the listener or reader). Crucially, de Saussure argued that the relationship between signifier and signified is **arbitrary**: there is no natural, necessary, or pre-given connection between a word and the concept it evokes.

“Dog” means what it means because a community of speakers agrees, tacitly and historically, that it does—not because the word contains or resembles the animal.

Implications for Legal Language

Terms such as “guilty,” “reasonable,” “intent,” “consent,” “Australian,” and “criminal” do not describe pre-existing conditions in the world. They are classifications produced through linguistic practice and institutionally enforced.

The designation “guilty” does not identify a property of the defendant; it is a **verdict**—a speech act, in Austin’s (1962) terminology, that changes the institutional status of a person in the world.

Before the word is spoken, there is a person. After it, there is a convict. The word does not describe a transformation. **It performs one.**

Barthes and Second-Order Signification

Roland Barthes (1957/2009) extended Saussure’s framework to identify a second order of signification he termed **myth**. Where first-order signification involves the direct relationship between signifier and signified, second-order signification involves taking an already-complete sign and using it as the signifier in a new system.

This is the mechanism by which ideologically loaded meanings come to appear natural and self-evident.

In the context of criminal justice, “prior criminal record” is a first-order sign denoting documented legal history. At the level of myth, it becomes a second-order signifier for “the kind of person who

does these things”—a naturalised narrative that transforms contingent institutional history into essential character.

The defence attorney who cannot suppress the jury’s knowledge of a defendant’s prior convictions is not fighting against evidence; they are fighting against a mythological structure.

2.2 Wittgenstein and Language Games: Meaning as Use

Ludwig Wittgenstein’s later philosophy provides a complementary framework for the analysis of legal language.

In *Philosophical Investigations*, Wittgenstein (1953) introduced the concept of the **language game**—a set of linguistic practices embedded in a form of life, where meaning is constituted by use rather than by reference to any underlying reality.

“The meaning of a word is its use in the language.” (§43)

For Wittgenstein, there is no metalanguage, no view from nowhere, no final description. There are only language games, each with its own grammar, and each game-community defining the rules of its own discourse.

The Legal Language Game

The legal system constitutes a language game in Wittgenstein’s sense—one with exceptionally high stakes and coercive enforcement mechanisms.

Legal language is not an attempt to describe the world in ordinary terms; it is a specialised grammar with its own rules for:

- What counts as evidence
- What counts as proof
- What counts as a person
- What counts as an act

The “reasonable person” standard in negligence law, for instance, is not a description of any actual human being but a grammatical device within the legal language game—a standard that judges (themselves participants in the game) calibrate through practice and precedent.

When a jury is asked to determine whether a defendant acted as a “reasonable person” would have acted, they are not being asked to consult reality; they are being asked to participate in a language game whose rules they do not fully understand and have not chosen.

Certainty and Doubt

In *On Certainty*, Wittgenstein (1969) argued that doubt presupposes a background of certainty—that one can only meaningfully question some things by taking other things for granted.

Legal proceedings construct their background certainties through procedural rules, evidentiary standards, and institutional authority. What is treated as self-evident in a courtroom—that the proceedings are neutral, that the law is clear, that the jury understands their instructions—is not self-evident at all.

It is the background against which the foreground drama of guilt or innocence is played out, and it is a background that systematically advantages the prosecution.

2.3 Foucault: Discourse, Power, and the Production of the Criminal Subject

Michel Foucault's contribution to this analysis is substantial and pervasive. Two bodies of work are of particular relevance: 1. His theory of discourse and power/knowledge (1972, 1980) 2. His genealogical analysis of the prison in *Discipline and Punish* (1977)

Discourse as Production

For Foucault, discourse is not merely language. It is the ensemble of practices, institutions, rules, and statements that produce particular objects of knowledge and particular kinds of subjects.

Criminal law, in this analysis, is a **discursive formation** that does not merely describe the criminal but **produces the criminal as a category of person**.

The criminal is not someone who has committed a prohibited act; the criminal is a particular kind of subject—one with a psychology, a history, a propensity—who exists as an object of knowledge within the discursive field of criminological, legal, psychiatric, and carceral practice.

Foucault's concept of the "delinquent" in *Discipline and Punish* captures this precisely: the judicial system ostensibly addresses acts, but the penal system addresses persons. The shift from punishing a deed to managing a type of subject is the historical achievement that Foucault traces through the transformation of penal practice.

Power/Knowledge

For Foucault, power does not merely suppress or prohibit; it produces. It produces knowledge, it produces subjects, it produces truth.

The expert witnesses, forensic psychologists, criminologists, police investigators, and legal professionals who participate in the criminal justice process are not neutral truth-tellers applying objective methodologies. They are nodes in a network of power/knowledge that produces the official account of what happened and who is responsible.

The authority of this account derives not from its correspondence to reality but from the institutional positions of those who produce it. A detective's interpretation of a suspect's body language carries evidentiary weight not because detectives have reliable access to truth but because detectives occupy an institutional position that authorises their interpretations.

The Genealogy of the Prison

Discipline and Punish provides a genealogy of the prison that anticipates the argument of this thesis. Foucault argues that the modern penal system did not emerge from humanitarian concern with justice but from the administrative need to manage populations, and that surveillance, normalisation, and the construction of the "dangerous individual" are its primary mechanisms.

The pre-interrogation detention regime described in this thesis is, in this analysis, a technology of normalisation: a process of stripping and repositioning the subject that has its genealogical roots in the disciplinary institutions of the nineteenth century.

The holding cell is not an unfortunate logistical necessity; it is a chamber of subject production.

2.4 Speech Act Theory: Language That Does Things

J.L. Austin's *How to Do Things with Words* (1962) introduced the concept of the **performative utterance**—a speech act that does not describe a state of affairs but constitutes one.

“I hereby sentence you to ten years imprisonment”

This is not a description of an event; it is the event.

Performative utterances require institutional conditions—**felicity conditions**, in Austin's terminology—to succeed: the speaker must occupy an authorised position, the circumstances must be appropriate, and the conventions must be recognised by the relevant community.

Legal Language as Performative

Legal language is saturated with performatives. The verdict, the sentence, the charge, the warrant, the caution—each of these is a speech act that transforms the institutional status of a person in the world.

This means that the legal determination of guilt is not a cognitive act of recognition (identifying someone as guilty) but a **social act of constitution** (making someone guilty).

Once this is understood, the question “was the verdict correct?” becomes philosophically complex. The verdict is not correct or incorrect in the way a factual description might be. It is, like a baptism or a contract, a social fact produced by the performance of a specific speech act under specific institutional conditions.

Indirect Speech Acts

John Searle's (1969) development of Austin's framework extended the analysis to **indirect speech acts**—utterances whose illocutionary force differs from their literal content.

- “Can you tell me where you were on the night of the fifteenth?” does not ask about the interviewee's capacity to provide this information; it requests the information.
- “We're going to be here a long time unless you help us understand what happened” does not describe a temporal situation; it threatens.

The literature on interrogation tactics is replete with indirect speech acts that disguise coercion as inquiry. The felicity conditions for valid confession evidence—that the statement be voluntary, uncoerced, and accurate—are systematically undermined by a discursive architecture built on precisely these indirect speech acts.

2.5 Words Twice Removed: The Epistemic Problem of Legal Fact-Finding

This thesis proposes that legal language operates at what might be called a **double remove** from reality.

The first remove is the general Saussurean point: no word corresponds directly to the thing it purports to describe.

The second remove is specific to legal and institutional language: legal descriptions of events are accounts of accounts—retrospective narrativizations of past states of affairs, produced under institutional pressures and interpreted through cognitive and cultural schemas that are largely invisible to those who deploy them.

This double remove is the epistemological condition within which the justice system operates, yet the system consistently presents its outputs as though they were simple observations of fact.

The Psychology of Memory

The empirical psychology of memory provides the most direct evidence for the second remove.

Loftus and Palmer (1974) demonstrated in a seminal experiment that the use of different verbs to describe the same event—“smashed” versus “hit” versus “contacted”—produced significantly different estimates of vehicle speed and, one week later, different reports of whether broken glass had been present (it had not).

The word changed the memory.

Subsequent decades of research have established the **misinformation effect** as robust and replicable (Loftus, 2005): post-event information, including information embedded in the linguistic formulation of questions, integrates with and alters the memory of the original event.

Legal testimony, produced through questioning—often leading, repetitive, and adversarial questioning—is therefore not a report of memory. It is a product of an interaction between memory and the linguistic environment in which recall occurs.

Cross-Examination as Memory Surgery

Cross-examination is explicitly designed to exploit the instability of memory. An effective cross-examination does not simply challenge what a witness says; it reconstructs what the witness believes they saw.

The opposing counsel who elicits the admission “I suppose it could have been red” from a witness who initially reported a blue car has not discovered truth; they have manufactured a new memory that serves the narrative they are constructing.

The jury, who will later deliberate on which story to believe, will not have access to the original event. They will have access to two competing linguistic constructions of it, produced under adversarial institutional conditions, and they will be asked to decide which is true.

2.6 Critical Legal Studies and the Indeterminacy of Law

The critical legal studies (CLS) movement of the 1970s and 1980s developed a sustained critique of legal formalism—the doctrine that legal decisions are determined by the application of fixed rules to established facts—and argued that law is fundamentally **indeterminate**: the same legal materials can support contradictory conclusions, and outcomes are shaped by political and ideological factors that legal reasoning conceals (Unger, 1983; Kennedy, 1997).

If legal rules do not determine outcomes, then what does?

The argument of this thesis is that **narrative, institutional power, and the linguistic construction of character and credibility do.**

Duncan Kennedy’s (1997) analysis of the ideological dimensions of legal argument is particularly illuminating. Kennedy argued that legal reasoning is not a technical practice insulated from political commitments but a form of rhetoric that frames contested value choices as though they were determinations of neutral principle.

The reasonable person standard, the balancing test, the doctrine of proportionality—each of these constructs presents what is in fact a politically loaded choice as a factual or logical conclusion.

In the context of criminal procedure, this rhetorical operation is most clearly visible in the treatment of confession evidence. Courts routinely hold that confessions obtained after hours of psychologically coercive interrogation, from a person in an acute stress state following hours of pre-interrogation detention, are “voluntary.”

The word “voluntary” is doing substantial ideological work—concealing a coercive reality behind a legal fiction that immunises the system from accountability for its own methods.

2.7 Synthesis: A Framework for Constructed Guilt

The theoretical framework assembled across the preceding sections provides a multi-levelled account of how guilt is constructed within the criminal justice system:

Theoretical Tradition	Contribution
Saussurean semiotics	Words do not report pre-existing facts but constitute new institutional realities
Barthes (myth)	Constructions come to appear natural and self-evident, making ideological operations invisible
Foucault (discourse)	Reveals the institutional and power-laden conditions within which legal knowledge is produced
Foucault (Discipline and Punish)	Locates the contemporary interrogation chamber within a genealogy of disciplinary practice
Austin/Searle (speech acts)	Identifies the specific performative mechanisms through which institutional transformation—from suspect to convict—is accomplished linguistically
Memory research	Demonstrates that linguistic construction extends backward in time: words rewrite the past events taken as evidence
CLS theory	Shows that formal legal reasoning actively enables construction by converting political choices into the appearance of neutral determinations

Methodological Implications

1. The analysis of language in criminal proceedings cannot be limited to explicit statements and formal arguments. It must attend to the implicit, the indirect, and the structural—the questions that are not asked, the descriptions that go unchallenged, the institutional silences that produce and protect dominant narratives.
2. The analysis cannot be synchronic: it must attend to the sequence of linguistic operations through which a case is constructed, from the first police report through charge and hearing and trial and verdict.
3. The framework demands attention to the bodily and institutional conditions under which language operates. Words spoken in a holding cell after hours of confinement do not have the same meaning—or the same moral status—as words spoken freely in a neutral environment.

The body, no less than the mouth, is a site of legal construction.

2.8 The Signal Inversion Effect: Empirical Foundation

The Core Paradox

The people who sound most guilty are most likely to be telling the truth.

This section documents what this thesis terms the **Signal Inversion Effect**: the systematic pattern whereby authentic cognitive and linguistic behaviours are misidentified as indicators of deception, while performed or rehearsed behaviours are mistakenly interpreted as signals of honesty.

The core paradox is deliberately counterintuitive—and it is supported by converging evidence from multiple independent research traditions.

Table 2.1: The Inversion Pattern

Behaviour	What It Actually Signals	What Observers Believe
Hedging (“I think,” “maybe”)	Genuine memory retrieval; honest uncertainty	Evasion; hiding something
Breaking eye contact	Cognitive effort; accessing memory	Shifty; untrustworthy
Fragmented narrative	Authentic trauma recall; real memory is messy	Incoherent; fabricating
Saying “I don’t know”	Accurate self-knowledge; genuine limits of memory	Incompetent; evasive
Confident, fluent delivery	Rehearsed narrative; possible deception	Honest; credible
Impersonal pronouns (“it,” “that”)	Distancing from false narrative (false confessions)	Not typically noticed
Spontaneous corrections	Genuine memory retrieval; commitment to accuracy	Inconsistency; changing story

Behaviour	What It Actually Signals	What Observers Believe
Reported confusion	Deep engagement; authentic processing	Failure to understand

2.8.1 Deception Detection: The Chance Problem

The foundational claim of police interrogation methodology is that trained investigators can identify deception through behavioural observation. If this claim is true, then investigator assessments of guilt carry epistemic weight—they are, at least partially, observations of reality. If the claim is false, investigator assessments are no more reliable than chance, and the guilty narrative they produce is constructed rather than observed.

The empirical literature on deception detection constitutes one of the most thoroughly replicated bodies of research in applied psychology. Its findings are consistent, well-powered, and directly devastating to the foundational claim of interrogation methodology.

Table 2.2: Deception Detection Accuracy—Meta-Analytic Summary

Study	N (Judges)	Accuracy %	95% CI	Population
Ekman & O’Sullivan (1991)	509	52.8	[49.1, 56.5]	Law enforcement, judges, psychiatrists
Vrij & Graham (1997)	156	53.2	[48.0, 58.4]	Police officers, UK
Meissner & Kassin (2002)	4,435	54.0	[52.5, 55.5]	Meta-analysis: civilians & investigators
Bond & DePaulo (2006)	24,483	54.3	[53.7, 54.9]	Meta-analysis: 247 studies
Vrij (2008) — overall	~5,000	54.0	[53.0, 55.0]	Mixed professional/civilian
Hartwig & Bond (2011)	~3,000	53.9	[52.1, 55.7]	Law enforcement
CHANCE BASE-LINE	—	50.0	[50.0, 50.0]	Theoretical maximum with zero information
WEIGHTED MEAN	~37,500	54.1	[53.6, 54.6]	Advantage over chance: 4.1 percentage points

Statistical Interpretation

The weighted mean accuracy of 54.1% across approximately 37,500 judgements represents a statistically significant deviation from chance ($p < .001$). However, statistical significance is here

misleading: the confidence interval is narrow precisely because the sample is enormous, not because the effect is large.

The practical significance is negligible.

In operational terms: a trained investigator assessing 100 suspects will correctly identify approximately 54 and incorrectly classify approximately 46. Of the incorrectly classified, approximately half will be innocent persons assessed as deceptive—the precise population who will be subjected to the coercive interrogation sequence designed to produce confession.

A system that, under equal base-rate assumptions, generates false accusation at a derived rate of approximately 23% among those subjected to it cannot, in any epistemologically defensible sense, be described as a truth-finding mechanism. (This figure is derived from Bond & DePaulo, 2006, assuming equal proportions of guilty and innocent suspects; real interrogation settings likely have higher proportions of guilty suspects, which would lower the rate. The figure illustrates the magnitude of the problem, not a direct empirical measurement.)

Critically, Kassin et al. (2005) demonstrated that **interrogation training does not improve accuracy; it improves confidence.** Trained investigators are more certain of assessments that are barely more reliable than a coin flip.

This dissociation of confidence from accuracy is precisely the cognitive architecture required to generate innocent convictions: investigators who are certain they have identified a guilty person, and who proceed to extraction of confession with that certainty, while being wrong approximately one time in four.

2.8.2 Memory Distortion: Language as Memory Author

If legal testimony is a reliable report of past events—a description of reality—then the linguistic practices of interrogation and cross-examination are neutral instruments of extraction.

If, conversely, testimony is a product of the interaction between memory and the linguistic environment in which recall occurs, then those who control the language control the evidence.

Elizabeth Loftus’s research programme, sustained over more than three decades and joined by hundreds of independent replication studies, establishes the latter beyond reasonable scientific dispute.

Table 2.3: Memory Distortion by Linguistic Manipulation—Effect Sizes

Study	Manipulation	Effect Size (d)	False Memory %	Key Finding
Loftus & Palmer (1974)	Single verb change	0.89	32%	“Smashed” vs. “contacted”: 16 km/h speed difference; 32% false glass memory
Loftus et al. (1978)	Leading question	0.72	17%	Subjects “recalled” non-existent barn

Study	Manipulation	Effect Size (d)	False Memory %	Key Finding
McCloskey & Zaragoza (1985)	Post-event info	0.61	~25%	Post-event narrative integration in 1-in-4 subjects
Belli (1989)	Misleading questions	0.68	~22%	Memory substitution in majority
Loftus (1993) — Lost in Mall	False narrative	—	25%	25% developed detailed false memories of events that never occurred
Loftus (2005) — 30-year review	Various linguistic	0.60–0.95	15–35%	Consistent replication over 30 years
AGGREGATE	Post-event language	d = 0.72	~22%	1-in-5 subjects will have memories altered by language alone

What $d = 0.72$ Means

An effect size of $d = 0.72$ is conventionally classified as **large** (Cohen, 1988).

In the specific context of Loftus and Palmer’s (1974) original study, a single word change—from “contacted” to “smashed”—produced a 16 km/h difference in speed estimates and, one week later, generated false memories of non-existent broken glass in 32% of subjects compared to 14% of controls: an odds ratio of 2.86.

The aggregate false memory rate of approximately 22% across studies means that in any group of five witnesses subjected to post-event linguistic manipulation—the routine condition of any interrogation or cross-examination—approximately one will carry a materially altered memory into their testimony.

They will not know this. They will be certain their memory is accurate. The court will have no mechanism for identifying them.

2.8.3 False Confession Rates: The System’s Own Output

The false confession literature provides the most direct evidence for the construction thesis because it operates at the level of the system’s output rather than its mechanisms.

Table 2.4: False Confession Rates Across Exoneration Datasets

Dataset / Study	N (Cases)	False Confession %	Jurisdiction	Notes
Kassin & Gudjonsson (2004)	Review	14–25%	US/UK	Range across exoneration studies
Gross et al. (2005)	340	15.3%	US	52 of 340 exonerates had falsely confessed
Garrett (2011) — DNA exonerations	250	10.8% (27 cases)	US	All 27 were factually innocent (DNA-confirmed)
National Registry (2023)	3,300+	~12%	US	False confessions disproportionately affect juveniles (42%)
Gudjonsson (2003) — UK Custodial	509	11.3%	UK	58 of 509 detained suspects reported prior false confession
Scherr et al. (2020)	Registry	~30%	US	Cumulative disadvantage analysis
AGGREGATE ESTIMATE	~4,500+	12–30%	US/UK/AU	Conservative: 1-in-8. Upper: nearly 1-in-3.

What These Numbers Mean

The conservative aggregate estimate—that 12–30% of exonerations involve false confessions—should be understood in the context of what exoneration databases capture. They capture only cases where:

1. The conviction occurred
2. The defendant survived long enough to pursue post-conviction remedies
3. Exculpatory evidence was discovered or preserved
4. The exoneration was successful

Each of these conditions filters out cases. **The documented false confession rate is therefore a floor, not a ceiling.**

The Garrett (2011) DNA exoneration data warrants particular attention because it eliminates evidentiary ambiguity. These are not cases where guilt was merely disputed; they are cases where post-conviction DNA testing established factual innocence beyond scientific dispute.

Of the first 250 such cases, **27 individuals—factually innocent people—had provided detailed, signed confessions.** Their confessions were, in the full technical sense, fabrications: accounts of events that did not happen, attributed to people who did not do them, produced by the system’s ordinary operation.

2.8.4 Suggestibility Under Detention

The three preceding sections address the mechanisms and outputs of the interrogation and evidentiary process. This section addresses the prior condition: the state of the person who enters the interrogation room following arrest and pre-trial detention.

Table 2.5: Suggestibility Elevation Under Pre-Interrogation Conditions

Study	Condition	Suggestibility Increase	Mechanism
Gudjonsson & Clark (1986)	Anxiety induction	+38%	Elevated anxiety increases yield and shift
Bain et al. (2014)	Sleep deprivation	+56%	One night produces dramatic GSS increases
Starcke & Brand (2012)	Acute stress	PFC impairment: significant	Stress impairs rational deliberation
Kassin et al. (2010)	Extended detention	+29–44%	Longer detention → higher confession rates
Gudjonsson (2003)	Custody vs. neutral	Yield: +31%; Shift: +42%	Custody itself elevates suggestibility

Study	Condition	Suggestibility Increase	Mechanism
COMPOUND EFFECT	Null sequence	+80–120%	Pre-interrogation sequence doubles baseline suggestibility

Compounding Effects

The critical feature of the pre-interrogation condition is that its constituent stressors are not additive but **compounding**.

A person who has been: - **Arrested** (acute stress response: HPA activation, PFC impairment) - **Deprived of sleep** (cognitive impairment equivalent to moderate intoxication) - **Stripped of identity materials** (Goffman’s mortification: impaired sense of agency) - **Isolated** (threat-system activation, heightened need for social approval)

...presents not with a 38% elevation in suggestibility, nor a 56% elevation, but with an accumulation of these effects whose combined magnitude is in the range of **80–120% above baseline**.

A voluntariness analysis conducted without reference to this compound baseline effect is not assessing whether the confession was voluntary in any substantively meaningful sense. It is performing a legal ritual whose outcome is determined by factors the analysis does not examine.

2.8.5 Convergent Validity: The Four Pillars Together

When multiple methodologically distinct research programmes, studying different variables through different methods in different populations, converge on the same conclusion, that conclusion achieves a degree of evidentiary support that no single study can provide.

Table A.1: Four-Pillar Convergent Validity—Synthesis

Pillar	Key Statistic	Null Hypothesis	Conclusion
I. Deception Detection	54.1% accuracy	Accuracy = 50% (chance)	REJECT null. But effect trivially small. System operates near-random.
II. Memory Distortion	$d = 0.72$; ~22% false memory	Language doesn’t affect memory	REJECT null. Large effect. 1-in-5 memories altered.
III. False Confessions	12–30% of exonerations	Confession = guilt indicator	REJECT null. Confession unreliable indicator.
IV. Suggestibility	+80–120% elevation	Detention doesn’t affect voluntariness	REJECT null. “Voluntariness” is legal fiction.
CONVERGENT FINDING	All four pillars	Guilt verdicts reflect culpability	REJECTED across all methodologies.

If the construction thesis were false—if guilt determination were a reliable, approximately accurate process—one would expect:

1. Deception detection accuracy substantially above chance
2. Memory stability under questioning
3. Low rates of false confession
4. Minimal effect of detention on statement voluntariness

The data yield the opposite on all four measures.

This is not the profile of a truth-finding system with known limitations. It is the profile of a guilt-production system that maintains the rhetorical apparatus of truth-finding.

2.9 Neurodivergent Communication and Systematic Credibility Bias

The Problem: Authentic Presentation as Perceived Deception

The Signal Inversion Effect operates with particular severity upon a population that has received insufficient attention in the criminological and legal literature: **neurodivergent individuals**—those with autism spectrum conditions, ADHD, functional neurological disorder (FND), PTSD, CPTSD, and related conditions that affect communication patterns, emotional presentation, and behavioural consistency.

The core problem is structural: the behavioural markers that neurotypical observers use to assess credibility—eye contact, emotional consistency, narrative fluency, social reciprocity—are precisely the domains in which neurodivergent individuals present atypically.

This atypical presentation is not a performance of deception; **it is authentic communication**. But the folk psychology of credibility assessment, operating through the inverted heuristics documented in section 2.8, systematically misreads neurodivergent authenticity as evidence of fabrication.

Table 2.6: Neurodivergent Presentation vs System Interpretation

Neurodivergent Presentation	Underlying Cause	System Interpretation
Reduced eye contact	Autistic sensory management; concentration	Shifty, untrustworthy, evasive
Flat or atypical affect	Autistic presentation; FND; PTSD dissociation	Cold, lacking remorse, faking, manipulative
Direct communication without hedging	Autistic pragmatic style; literal processing	Aggressive, uncooperative, hostile, contemptuous
Fragmented trauma narrative	PTSD/CPTSD memory structure; genuine recall	Inconsistent, fabricating, lying, unreliable

Neurodivergent Presentation	Underlying Cause	System Interpretation
Variable symptoms across time	FND fluctuation; stress-responsive conditions	Faking illness, malingering, exaggerating
Says “I don’t know” when uncertain	Genuine epistemic accuracy; intellectual honesty	Evasive, hiding something, uncooperative
Corrects interviewer’s factual errors	Commitment to accuracy; autistic literalism	Argumentative, difficult, obstructive
Detailed or repetitive responses	Autistic communication; thoroughness	Over-rehearsed, scripted, coached
Calm presentation during accusation	Autistic emotional regulation; dissociation	Lack of appropriate distress, cold, guilty
Distressed presentation during accusation	Anxiety; appropriate response to false accusation	Performing, dramatic, manipulative
Maintains consistent account	Truthful recall	Suspiciously consistent, rehearsed
Account varies in minor details	Normal memory variation	Inconsistent, changing story, lying

Observation: The table demonstrates a structural double bind. For each pair of possible presentations (calm/distressed, consistent/varying, direct/hedging), both alternatives are interpreted as indicators of deception or guilt. There is no presentation style available to the neurodivergent individual that would be interpreted as credible.

2.9.0 The Fundamental Invalidity of Behavioural Heuristics

This section must be read before any other section in this chapter.

The preceding and following sections document specific ways that neurodivergent presentations are misinterpreted. A critical clarification is required:

For every behaviour documented as “misinterpreted as guilt,” the opposite behaviour is equally misinterpreted as guilt.

Behaviour	Interpretation	Opposite Behaviour	Interpretation
Eye contact	Rehearsed, too confident, staring	No eye contact	Shifty, evasive, hiding
Calm presentation	Cold, no remorse, inappropriate	Distressed presentation	Guilty panic, performing
Consistent account	Too consistent, rehearsed, scripted	Varying account	Inconsistent, changing story, lying
Detailed recall	Over-prepared, memorised	Vague recall	Evasive, hiding details

Behaviour	Interpretation	Opposite Behaviour	Interpretation
Quick responses	Rehearsed, pre-prepared	Slow responses	Thinking up lies, evasive
Direct communication	Aggressive, hostile, contemptuous	Hedged communication	Uncertain, evasive, hiding
Answers questions fully	Over-explaining, defensive	Brief answers	Uncooperative, withholding
Volunteers information	Controlling narrative, suspicious	Waits for questions	Uncooperative, evasive
Shows emotion	Performing, manipulative	Shows no emotion	Cold, lacks remorse
Looks at questioner	Confrontational, defiant	Looks away	Evasive, shifty

The double bind is complete. There is no behaviour that escapes negative interpretation.

Compensatory Behaviour and Its Interpretation

An individual who has been accused of deception based on one behaviour (e.g., “you’re not looking at me”) may attempt to compensate by producing the opposite behaviour (maintaining eye contact). This compensation is then interpreted as: - Evidence of guilt (“now you’re staring, that’s suspicious”) - Evidence of coaching (“someone told you to make eye contact”) - Evidence of manipulation (“you’re trying to appear innocent”) - Evidence of consciousness of prior error (“you know you looked guilty before”)

The attempt to correct misinterpreted behaviour produces new misinterpreted behaviour.

Why This Section Preempts Weaponisation

This thesis cannot be used to identify “what innocent people do” because the thesis demonstrates that there is no behaviour that reads as innocent under heuristic assessment.

The thesis does not say: “Covering ears indicates innocence, therefore someone covering their ears is innocent.”

The thesis says: “Covering ears may indicate sensory overload. It may also indicate nothing. Eye contact may indicate confidence. It may also indicate rehearsal. Every behaviour and its opposite are interpreted as guilt. Therefore behavioural heuristics are invalid as a class.”

Any attempt to use this thesis to argue “this person did X, and the thesis says X indicates innocence, therefore they are performing innocence” fails because: 1. The thesis demonstrates that NOT-X is equally interpreted as guilt 2. The thesis demonstrates that the observer cannot distinguish authentic X from performed X 3. The thesis demonstrates that the attempt to compensate (switching from X to NOT-X) is itself interpreted as guilt 4. The thesis concludes that behavioural heuristics do not work, not that they work in reverse

The conclusion is not “interpret behaviours differently.” The conclusion is “behavioural interpretation is unreliable and should not be used.”

The Lived Experience of Heuristic Invalidity

The individual who has been repeatedly assessed using behavioural heuristics lives in a state of:
- Fear of being misinterpreted - Uncertainty about which behaviours are “safe” - Hypervigilance about own presentation - Attempts to compensate that produce new misinterpretations - Accumulated evidence that no behaviour is safe - Awareness that both action and inaction are interpreted negatively

This state is itself interpreted as: - Anxiety indicating guilt - Excessive self-monitoring indicating consciousness of deception - Inconsistency (varying behaviour) indicating unreliability

The individual’s awareness that they will be misinterpreted, and their attempts to avoid misinterpretation, become further evidence of guilt.

This is a closed system. There is no exit through behaviour modification. The only exit is the recognition that behavioural heuristics are fundamentally invalid as assessment tools.

2.9.1 Autism and the Credibility Gap

The research literature on autism and credibility assessment in legal contexts, while limited, is consistently concerning.

Maras and Bowler (2014) documented that autistic witnesses produce testimony that is **less likely to be believed by mock jurors**, despite being no less accurate and in some respects more accurate than neurotypical testimony.

The specific features that reduced credibility were:

- **Reduced eye contact** — a core feature of autistic communication, often reflecting sensory management or concentration, systematically misinterpreted as evasiveness
- **Flat or atypical emotional affect** — the absence of expected emotional displays interpreted as lack of genuine experience
- **Literal and detailed responses** — exhaustive detail or resistance to paraphrasing interpreted as over-rehearsed or evasive
- **Inconsistent narrative structure** — non-linear recall patterns reflecting genuine memory organisation misinterpreted as fabrication
- **Difficulty with open-ended questions** — preference for specific questions interpreted as reluctance to volunteer information

Crane et al. (2016) found that autistic adults recalled fewer central and peripheral details than neurotypical adults when reporting a witnessed event, and that their recall was significantly more affected by misleading questions—precisely the conditions of cross-examination.

The intersection of reduced baseline recall with heightened suggestibility creates a compound vulnerability: autistic witnesses may provide testimony that is both less complete and more contaminated by the linguistic environment.

2.9.2 Functional Neurological Disorder (FND) and Symptom Disbelief

Functional neurological disorder presents a distinct but related credibility challenge. FND is characterised by neurological symptoms—seizures, paralysis, speech difficulties, movement disorders, gait abnormalities, tremor, sensory disturbances—that are not attributable to structural neurological damage but are nonetheless genuine, involuntary, and disabling.

FND is now understood as a disorder of brain network function, with documented abnormalities in the integration of motor planning, execution, and sensory feedback (Edwards et al., 2012; Espay et al., 2018). The condition is neurological, not psychiatric; the symptoms are involuntary, not performed.

However, the historical categorisation of FND as “psychogenic,” “conversion disorder,” or “hysteria,” with its implicit suggestion of psychological causation or voluntary production, has created a persistent cultural framework in which FND symptoms are systematically disbelieved by medical professionals, investigators, and legal decision-makers.

For a person with FND in legal proceedings, the implications are severe and specific:

- **Symptom variability** — FND symptoms characteristically fluctuate in intensity and presentation across time, across contexts, and in response to attention and stress. A person who cannot walk at one point but can walk at another is exhibiting a core diagnostic feature of FND, not evidence of fabrication. However, variability is precisely the feature that naive observers—including investigators, prosecutors, and juries—interpret as evidence of faking. The diagnostic criterion becomes the credibility deficit.
- **Stress responsiveness** — FND symptoms are often exacerbated by stress, including the stress of accusation, interrogation, and legal proceedings. The person whose symptoms worsen during cross-examination, or who experiences a functional seizure during a police interview, may appear to observers to be performing distress strategically rather than experiencing genuine neurological dysfunction.
- **Inconsistency with lay expectations** — observers expect symptoms to map onto known structural conditions and to present consistently. FND symptoms, which reflect functional network disruption rather than structural lesion, may present in patterns that appear “neurologically impossible” or “inconsistent with organic disease” to observers operating with a structural model of neurological illness. This apparent inconsistency is then coded as evidence of fabrication.
- **Distractibility effects** — FND symptoms characteristically diminish when attention is directed elsewhere. This is a diagnostic feature, not evidence of voluntary control. However, the observation that symptoms “disappear” when the person is distracted is routinely interpreted as evidence that symptoms were never genuine.
- **Historical medical disbelief** — many individuals with FND have extensive histories of being disbelieved, dismissed, or accused of faking by medical professionals prior to diagnosis. This history of medical gaslighting compounds the trauma of being disbelieved in legal contexts and may produce defensive presentation patterns—hypervigilance, over-documentation, anticipatory explanation—that themselves invite further disbelief.
- **Comorbidity patterns** — FND frequently co-occurs with conditions including autism, ADHD, PTSD, and chronic pain conditions. The compound presentation of multiple conditions, each with its own credibility deficits, creates multiplicative rather than additive vulnerability.

The person with FND who enters the criminal justice system—whether as defendant, complainant, or witness—faces a structural credibility deficit that derives not from any deception on their part but from the fundamental incompatibility between the

phenomenology of their condition and the folk psychology through which credibility is assessed.

2.9.2a The Clinical Reality of FND Assessment in Criminal Justice Settings

The intersection of FND with criminal justice settings creates a specific compound vulnerability that warrants extended analysis.

When a person with FND is detained, their symptoms are likely to intensify. FND symptoms are characteristically stress-responsive — they worsen under conditions of anxiety, uncertainty, and perceived threat (Edwards et al., 2012). The pre-interrogation detention environment described in Chapter 3 produces precisely these conditions. The person with FND who is arrested will, as a predictable neurobiological consequence of detention, experience worsening of their symptoms.

The criminal justice system then interprets this worsening as evidence against the person’s credibility. The following sequence is documented in clinical and forensic literature:

1. **Arrest produces stress.** Stress exacerbates FND symptoms.
2. **Exacerbated symptoms are observed.** Officers, custody staff, and forensic clinicians observe seizures, tremor, speech difficulties, or gait abnormalities that appear to fluctuate.
3. **Fluctuation is interpreted as fabrication.** Because FND symptoms vary with stress — because they are, by clinical definition, functionally rather than structurally caused — the variability is interpreted as evidence that the symptoms are under voluntary control and therefore fake.
4. **“Fake” symptoms reduce credibility.** The person whose medical symptoms have been dismissed as fabrication is now assessed as someone who fakes things — a liar. Their testimony, their account of events, and their presentation are all filtered through this credibility-reducing frame.

The clinical literature is unambiguous: FND symptoms are involuntary, genuine, and disabling (Espay et al., 2018; Stone et al., 2010). Symptom variability is a diagnostic criterion, not evidence of fabrication. The fluctuation of symptoms in response to stress is a defining feature of the condition, not a reason to doubt its reality.

Yet the criminal justice system — staffed by police officers, custody nurses, forensic psychologists, and magistrates who receive minimal or no training in functional neurological conditions — routinely interprets this variability as evidence of malingering. The person is not believed. Their genuine medical condition is used against them. Their innocence becomes structurally illegible because the instrument through which their credibility is assessed — the behavioural heuristic — reads their involuntary neurological presentation as voluntary deception.

This is the Signal Inversion Effect operating on a medical condition. The authentic presentation of FND mimics what untrained observers expect deception to look like. The person cannot present their genuine symptoms without triggering credibility-reducing inferences. And when they try to suppress their symptoms — to appear ‘normal’ — the suppression effort produces the same strained, controlled presentation that the heuristics read as rehearsed deception.

There is no available presentation for a person with FND that reads as credible to the criminal justice system’s assessment instruments.

2.9.2a The Clinical Reality of FND Assessment in Criminal Justice Settings

The intersection of FND with criminal justice settings creates a specific compound vulnerability that warrants extended analysis.

When a person with FND is detained, their symptoms are likely to intensify. FND symptoms are characteristically stress-responsive — they worsen under conditions of anxiety, uncertainty, and perceived threat (Edwards et al., 2012). The pre-interrogation detention environment described in Chapter 3 produces precisely these conditions. The person with FND who is arrested will, as a predictable neurobiological consequence of detention, experience worsening of their symptoms.

The criminal justice system then interprets this worsening as evidence against the person’s credibility. The following sequence is documented in clinical and forensic literature:

1. **Arrest produces stress.** Stress exacerbates FND symptoms.
2. **Exacerbated symptoms are observed.** Officers, custody staff, and forensic clinicians observe seizures, tremor, speech difficulties, or gait abnormalities that appear to fluctuate.
3. **Fluctuation is interpreted as fabrication.** Because FND symptoms vary with stress — because they are, by clinical definition, functionally rather than structurally caused — the variability is interpreted as evidence that the symptoms are under voluntary control and therefore fake.
4. **“Fake” symptoms reduce credibility.** The person whose medical symptoms have been dismissed as fabrication is now assessed as someone who fakes things — a liar. Their testimony, their account of events, and their presentation are all filtered through this credibility-reducing frame.

The clinical literature is unambiguous: FND symptoms are involuntary, genuine, and disabling (Espay et al., 2018; Stone et al., 2010). Symptom variability is a diagnostic criterion, not evidence of fabrication. The fluctuation of symptoms in response to stress is a defining feature of the condition, not a reason to doubt its reality.

Yet the criminal justice system — staffed by police officers, custody nurses, forensic psychologists, and magistrates who receive minimal or no training in functional neurological conditions — routinely interprets this variability as evidence of malingering. The person is not believed. Their genuine medical condition is used against them. Their innocence becomes structurally illegible because the instrument through which their credibility is assessed — the behavioural heuristic — reads their involuntary neurological presentation as voluntary deception.

This is the Signal Inversion Effect operating on a medical condition. The authentic presentation of FND mimics what untrained observers expect deception to look like. The person cannot present their genuine symptoms without triggering credibility-reducing inferences. And when they try to suppress their symptoms — to appear ‘normal’ — the suppression effort produces the same strained, controlled presentation that the heuristics read as rehearsed deception.

There is no available presentation for a person with FND that reads as credible to the criminal justice system’s assessment instruments.

2.9.3 PTSD, CPTSD, and Trauma Presentation

Post-traumatic stress disorder and complex PTSD produce presentation patterns that intersect with the Signal Inversion Effect in specific ways.

Trauma Memory Structure

The neuroscience of trauma memory is well-established (van der Kolk, 2014; Brewin, 2011). Traumatic memories are: - Encoded under conditions of extreme stress that impair hippocampal function - Stored in fragmented, sensory-dominant, non-narrative form - Retrieved in non-linear sequences, often triggered by sensory or contextual cues - Subject to intrusion and avoidance patterns that disrupt voluntary recall

This is the normal structure of trauma memory. It is not evidence of fabrication; it is diagnostic of genuine traumatic experience.

However, the criminal justice system's expectations of witness memory are derived from a folk-psychological model in which memory operates as a recording device: events are encoded accurately, stored stably, and retrieved in the order they occurred. Deviation from this model—fragmentation, non-linearity, gaps, sensory intrusions—is interpreted as evidence of unreliability or fabrication.

The person with PTSD is providing an accurate report of their memory. The system interprets accuracy as evidence of lying.

Specific Presentation Vulnerabilities

- **Fragmented narrative** — trauma memory is not stored or retrieved in linear narrative form. The person with PTSD who cannot provide a coherent timeline is not being evasive; they are accurately reporting the structure of traumatic memory. But fragmentation is precisely what credibility heuristics interpret as fabrication.
- **Detail inconsistency** — trauma memories may be vivid for some details (sensory, emotional) while absent for others (temporal sequence, peripheral context). This pattern—clear memory for some elements, no memory for others—appears inconsistent to observers expecting uniform recall.
- **Emotional dysregulation** — PTSD may produce either hyperarousal (excessive emotional display) or hypoarousal (flattened affect, dissociation). Both patterns deviate from the expected emotional presentation of a “genuine” victim or witness, and both attract disbelief.
- **Avoidance** — a core symptom of PTSD is the avoidance of trauma-related stimuli. A person who is reluctant to recount details of a traumatic event is exhibiting a symptom, not concealing information. But avoidance presents as evasiveness to untrained observers.
- **Dissociation** — dissociative responses during questioning may produce apparently inconsistent statements, gaps in recall, or affect that seems disconnected from content. All of these invite interpretation as deception.
- **Re-traumatisation effects** — repeated questioning about traumatic events can produce deterioration in recall quality, increased fragmentation, and emotional dysregulation. The person who provides a less coherent account on the third interview than the first is exhibiting a normal response to re-traumatisation, not evidence of a fabricated story “falling apart.”

Delayed Disclosure

Delayed disclosure of traumatic events is normative, particularly for: - Sexual assault (median delay to first disclosure: months to years) - Childhood abuse (median delay: decades) - Institutional abuse (disclosure often follows external triggers such as media coverage or legal proceedings involving other victims)

The folk-psychological expectation that genuine victims disclose immediately is empirically false. However, delayed disclosure is routinely used as evidence against complainant credibility.

Memory Evolution Over Time

Research on trauma memory documents that: - Initial accounts are typically incomplete due to avoidance and encoding deficits - Memory for traumatic events may become more detailed over time as avoidance decreases - Therapeutic processing may facilitate recall of previously unavailable details - Exposure to related information (media, other accounts) may trigger retrieval of genuine memories

This pattern—of memory that evolves, becomes more detailed, and incorporates elements triggered by external information—is consistent with genuine traumatic memory. It is also the pattern that the system interprets as “story changing,” “embellishment,” or “contamination.”

The person whose account becomes more detailed over time, who recalls elements they did not initially report, and whose memory is triggered by external information, is exhibiting the documented phenomenology of trauma memory. The system interprets this as evidence of fabrication.

2.9.4 The Direct Communication Problem

A distinct feature of autistic communication, documented extensively in the autism literature (Baron-Cohen, 2008; Milton, 2012) and confirmed in phenomenological research, is **directness**—a tendency to say what one means without the social hedging, face-saving, and indirect signalling that characterises neurotypical communication.

This directness is not rudeness, hostility, or contempt; it is a different pragmatic style reflecting different underlying processing of social communication. However, in the context of credibility assessment by neurotypical observers, directness produces specific and severe vulnerabilities:

- **Absence of expected deference** — the person who does not perform appropriate subordination to institutional authority—who speaks to a police officer, prosecutor, or judge as an equal rather than as a supplicant—may be read as defiant, arrogant, uncooperative, or contemptuous of the process. The absence of performed submission is interpreted as hostility.
- **Failure to perform expected distress** — the person who reports a distressing event in matter-of-fact terms, without the expected emotional performance of distress, may be disbelieved precisely because they are not performing. Conversely, visible distress may be interpreted as manipulative performance. There is no emotional presentation that satisfies the expectation.
- **Correction of factual errors** — the autistic individual who corrects an interviewer’s misstatement, a prosecutor’s inaccurate characterisation, or a witness’s factual error may be read as argumentative, pedantic, difficult, or evasive, rather than as genuinely committed to accuracy. The pursuit of factual precision—a core autistic trait—is coded as obstruction.
- **Refusal to speculate** — the person who says “I don’t know” when they genuinely do not know, rather than providing a plausible-sounding estimate, may be read as uncooperative, evasive, or obstructive. The person who refuses to agree with a leading question that contains a false premise may be read as hostile.

- **Literal response to questions** — the autistic individual may answer the question that was literally asked rather than the question that was pragmatically implied, producing responses that appear evasive or off-topic to neurotypical listeners.
- **Absence of social lubrication** — the person who does not engage in expected social pleasantries, does not smile appropriately, does not acknowledge the social status of interlocutors, or does not perform gratitude for being “given the opportunity to tell their side” may be read as cold, hostile, or lacking appropriate social awareness.

The direct communication style characteristic of autism is, in the context of criminal justice proceedings, systematically penalised. The very features that make autistic communication honest—its literalism, its precision, its absence of social manipulation—are the features that neurotypical observers interpret as evidence of dishonesty, hostility, or guilt.

2.9.5 Autistic Stress Responses and Systematic Misinterpretation

The autistic stress response differs from the neurotypical stress response in documented ways that have specific implications for interrogation settings.

Autistic Laughter

Autistic individuals may laugh in response to acute stress, anxiety, or sensory overload. This is a documented phenomenon in the autism literature (Samson & Hegenloh, 2010; Hudenko et al., 2009). The laughter is: - Involuntary - Not indicative of amusement - Not indicative of contempt or disrespect - A dysregulated stress response, comparable to crying or freezing

In the interrogation context, autistic laughter in response to accusation or pressure is systematically misinterpreted as: - Contempt for the process - Lack of appropriate seriousness - Evidence that the person finds the situation amusing - Consciousness of guilt (“laughing because they know they did it”) - Defiance or mockery of authority

The autistic person who laughs under interrogation stress is exhibiting an involuntary neurological response to overwhelming input. The system interprets this as evidence of guilt or contempt.

Sensory Covering Behaviours

Autistic individuals experiencing sensory overload may engage in protective behaviours including: - Covering ears with hands - Closing eyes - Turning away from stimulation - Rocking or other repetitive movements - Requesting that stimulation stop

These are regulatory behaviours—attempts to manage overwhelming sensory input. They are not: - Refusal to cooperate - Evidence of evasion - Defiance of authority - Attempts to avoid hearing incriminating information

In an interrogation setting where investigators are speaking loudly, repeatedly, or for extended periods, the autistic individual may cover their ears as a sensory regulation response. This is systematically interpreted as: - Refusal to listen - Childish or immature behaviour - Evidence of guilt (“doesn’t want to hear the truth”) - Obstruction - Contempt

The autistic person who covers their ears during prolonged loud interrogation is protecting themselves from sensory injury. The system interprets this as consciousness

of guilt.

Visible Distress Responses

Autistic individuals under acute stress may exhibit: - Tremor or shaking - Visible panic - Speech difficulties (stammering, word-finding problems, mutism) - Motor difficulties - Dissociative responses - Meltdown or shutdown responses

These responses may be more visible, more intense, or more prolonged than neurotypical stress responses. They may also present atypically—for example, alternating between apparent calm and intense distress, or presenting with physical symptoms (shaking, tremor) while verbally appearing composed.

The visibility and intensity of autistic stress responses creates a specific double bind: - Intense visible distress is interpreted as “guilty panic” or “performance” - Attempts to regulate distress (covering ears, looking away, rocking) are interpreted as evasion - Successful regulation (appearing calm) is interpreted as inappropriate lack of distress

There is no presentation of autistic stress that is interpreted as consistent with innocent distress at false accusation.

The Near-False-Confession Phenomenon

Research on interrogation and autism documents that autistic individuals may experience intense pressure to agree with interrogators simply to end the aversive interaction (Gudjonsson & Joyce, 2011; North et al., 2008). The autistic individual may: - Want desperately for the interrogation to stop - Understand that agreement would make it stop - Feel overwhelming pressure to say what the interrogator wants to hear

However, autistic literalism may paradoxically provide protection against false confession. The autistic individual who cannot bring themselves to state something they know to be false—even under extreme pressure—is exhibiting a core autistic trait: difficulty with deliberate falsehood.

This creates a specific presentation pattern: - Visible extreme distress - Clear desire for the interrogation to end - Possible verbal expressions of wanting to agree - Ultimate inability to state the false thing

This pattern—distress, pressure, near-capitulation, but not actual false confession—is highly significant. It indicates: - The interrogation methodology was producing the conditions for false confession - The individual’s resistance was neurologically based, not strategic - A non-autistic individual under identical conditions might have falsely confessed

The autistic individual who wanted to falsely confess but could not bring themselves to do so is exhibiting the protective effect of autistic literalism. The system interprets visible distress and near-capitulation as evidence of guilt, not as evidence that the interrogation was coercive.

Combined Presentation Under Interrogation

Consider the autistic individual who, during extended accusatory interrogation: - Shakes visibly (stress response) - Covers their ears (sensory regulation) - Laughs (stress dysregulation) - States “I can’t hear you” (literal description of sensory state) - Exhibits visible desire for the interrogation to end - Does not confess (autistic literalism preventing false statement)

Each of these behaviours is a documented autistic response to overwhelming stress. Together, they constitute a coherent picture of an autistic individual in acute distress, engaging in regulatory behaviours, experiencing sensory overload, and being protected from false confession by their inability to state known falsehoods.

The system interprets this same presentation as: - Contempt (laughter) - Obstruction (covering ears) - Defiance (“I can’t hear you”) - Consciousness of guilt (extreme distress) - Evidence that pressure is warranted (distress indicates something to hide)

The autistic individual’s authentic stress response under coercive interrogation is decoded as evidence supporting the coercion. The behaviours that indicate the interrogation is harmful are interpreted as evidence that it should continue.

Question Phrasing Effects

Autistic literal processing means that responses are highly sensitive to the precise phrasing of questions. The same information may be: - Accessible when asked in one form - Inaccessible when asked in another form - Provided accurately to a precisely-worded question - Provided inaccurately to an imprecisely-worded question (where the inaccuracy reflects the literal meaning of the question, not deception)

This creates a specific pattern: - Cooperative, information-gathering approaches (PEACE-style) elicit accurate, detailed accounts - Accusatory, leading approaches (Reid-style) elicit defensive, incomplete, or literally-responsive answers that appear evasive

The autistic individual who provides detailed, coherent information in a supportive interview context and fragmented, defensive responses in an accusatory context is exhibiting sensitivity to question phrasing and interviewer approach, not inconsistency indicating deception.

The same individual, with the same memory, will produce different quality accounts depending on the questioning approach. The system interprets variation in account quality as evidence of manipulation rather than as evidence that questioning approach affects recall.

Topic-Dependent Emotional Presentation

Autistic emotional regulation is topic-dependent. An individual may: - Speak calmly and fluently about neutral topics - Experience visible distress, speech difficulties, or shutdown when discussing triggering topics - Transition between calm and distressed presentation depending on conversational content

This is consistent with: - Genuine traumatic memory with associated emotional dysregulation - Autistic difficulty with emotional transitions - Topic-specific triggering of trauma responses

The system interprets this pattern as: - Evidence of “performing” distress when convenient - Evidence that distress is voluntary and strategic - Evidence of “faking” because “real” distress would be constant - Inconsistency indicating manipulation

The autistic individual who can discuss the weather calmly but becomes visibly distressed when asked about traumatic events is exhibiting topic-specific trauma response. The system interprets the ability to be calm about other topics as evidence that distress about the traumatic topic is performed.

The Cooperative Environment Effect

Research on autistic communication and information retrieval documents that: - Supportive, non-threatening environments facilitate recall - Collaborative, information-building approaches produce more complete accounts - Accusatory approaches produce defensive, incomplete, or shutdown responses - The quality of information obtained is a function of the environment, not solely of the individual's knowledge

This has specific implications for witness and defendant interviews: - An autistic witness interviewed collaboratively may provide detailed, useful information - The same witness interviewed accusatorily may appear unhelpful, evasive, or hostile - The difference is environmental, not volitional

The autistic individual who provides excellent information in a collaborative setting and poor information in an accusatory setting is responding predictably to environmental conditions. The system interprets this difference as evidence of selective cooperation or manipulation.

Visual-Dependent Memory Retrieval

A subset of autistic individuals exhibit visual-dependent memory processing: the retrieval of episodic memories is mediated by internal visual imagery, and verbal narration of events requires concurrent access to visual representations (Kana et al., 2006; Kosslyn et al., 2006).

This processing style has specific characteristics: - Verbal recall is dependent on internal visual “replay” of events - External visual cues (photographs, video, returning to locations) may be required to access memories - Verbal recall without visual access may be incomplete, fragmented, or unavailable - The quality of verbal recall varies with the accessibility of visual imagery - Affect and stress states influence the accessibility of visual processing

For the individual with visual-dependent memory retrieval: - Telling a story requires “seeing” the story internally - Without visual access, they may be unable to narrate events they genuinely remember - Stress, fatigue, or emotional state may disrupt visual processing, producing variable recall quality - External visual supports (photos, diagrams, video) may dramatically improve recall - Interview environments that do not provide visual supports produce artificially impaired recall

This creates specific vulnerabilities in legal contexts:

1. **Variable recall quality:** The same individual may provide detailed recall on one occasion and fragmented or absent recall on another, depending on visual processing accessibility. This appears as inconsistency.
2. **Affect-dependent access:** Emotional state affects visual processing access. High stress may impair the visual system, producing worse recall precisely when accurate recall is most needed.
3. **External support dependence:** The individual who can provide excellent recall when shown photographs but poor recall without them is not being evasive; they are exhibiting visual-dependent processing.
4. **Verbal narration difficulties:** The individual may know what happened—may have clear visual memory—but be unable to translate this into verbal narrative without visual support.

The individual with visual-dependent memory processing may appear to have inconsistent, unreliable, or selectively available memory. In fact, they have consistent memory with variable access, mediated by a visual processing system that is sensitive to affect, stress, and environmental support.

The system interprets variable verbal recall as: - Inconsistency indicating fabrication - Selective memory indicating evasion - Poor recall indicating lack of genuine experience - Improved recall with visual aids as “coaching” or “suggestion”

The individual’s authentic memory processing style—which produces genuine, accurate memories accessed through visual imagery—is decoded as evidence of unreliability or deception.

Traumatic Memory and Visual Processing Pain

For the individual with visual-dependent memory retrieval, accessing traumatic memories presents a specific problem: narrating the event requires visually re-experiencing it.

This creates a protective mechanism: - The visual system may block access to traumatic imagery - This blocking protects against re-traumatisation - But it also prevents verbal narration - The individual cannot “tell” because they cannot “see” without pain

This is distinct from volitional avoidance: - The individual is not choosing not to talk - The individual may want to provide information - The protective blocking is involuntary - The individual may experience frustration at their own inability to access and narrate

The presentation in interview contexts: - Inability to provide detailed narrative of traumatic events - Possible ability to provide fragmented, non-visual details (facts known without imagery) - Visible distress when attempting to access visual memory - Potential shutdown or dissociation when visual access is forced - Better recall for peripheral or non-traumatic aspects of events

The system interprets this as: - Selective memory (remembers some things, not others) - Evasion (won’t describe the key events) - Lack of genuine experience (if it really happened, they would be able to describe it) - Inconsistency (detailed about some things, vague about others)

The individual who cannot narrate traumatic events because visual access causes re-traumatisation is exhibiting a protective mechanism, not evasion. The system interprets protective blocking of traumatic visual replay as evidence of fabrication.

Dysregulated Output Under Visual Access

In some cases, the individual with visual-dependent processing can access traumatic visual memory, but the output is severely dysregulated: - Fragmented, disordered verbal output - Intense visible emotional response - Speech disruption, word-finding difficulties - Physical symptoms (shaking, crying, dissociation) - Output that does not resemble coherent narrative

Observers interpret intense dysregulated response as: “Something must have happened for you to react this way” — meaning, the intensity of reaction is read as evidence of involvement or guilt, rather than as evidence of trauma response.

The intense reaction to accessing traumatic memory is interpreted as consciousness of guilt rather than as pain.

Pre-Existing Sensitivity and Accumulated Vicarious Trauma

Intensity of response to trauma-related content may be compounded by:

1. **Autistic sensory sensitivity:** Many autistic individuals experience intense distress in response to imagery of bodily harm, injury, or suffering. This is a sensory and empathic sensitivity, not personal involvement.

2. **Professional vicarious trauma:** Individuals who have worked in roles involving trauma exposure—victim support services, emergency response, healthcare, social work, counselling—accumulate vicarious trauma from repeated exposure to others’ pain.

Such work often requires:

- Maintaining regulated external presentation (“stone face”)
- While remaining present and empathetic
- While viewing images and hearing accounts of suffering
- Over extended periods (months, years)

This produces:

- Accumulated vicarious traumatisation
- Sensitisation to trauma-related content
- Trigger responses to reminders
- Reduced capacity to regulate when exposed to similar content

3. **Trigger activation:** Content related to current legal matters may activate triggers from professional trauma exposure, producing responses that appear disproportionate to the current matter but are proportionate to accumulated exposure.

The Attribution Error

When an individual with: - Pre-existing autistic sensitivity to harm imagery - Accumulated vicarious trauma from professional exposure - Trigger responses to trauma-related content

exhibits intense distress when discussing events involving harm, the system attributes that distress to the current matter—specifically, to guilt or involvement in the current matter.

The reasoning is: “This person is reacting intensely to this specific event, therefore their reaction is about this specific event, therefore they must have been involved.”

The reality may be: “This person is reacting intensely because they have: - Baseline high sensitivity to harm content (autism) - Accumulated trauma from months/years of professional exposure to others’ suffering - Triggers that are activated by any harm-related content - Visual-dependent processing that forces them to ‘see’ content to discuss it”

The intense reaction is evidence of the individual’s trauma history and processing style. The system interprets it as evidence of guilt in the current matter.

The Stone Face / Dysregulation Contrast

Individuals with professional trauma-exposure backgrounds may have developed the capacity to maintain regulated presentation in professional contexts—while being unable to maintain that regulation when they are the subject of investigation rather than the support provider.

This produces a presentation that appears inconsistent: - “You worked in victim support, you must be able to talk about this” - “If you could handle that work, why can’t you handle these questions” - “You seemed fine when you were working, but now you’re falling apart”

The reality is: - Professional regulation was a learned, effortful skill - It was maintained in a supportive, purposeful context - It depleted over time, contributing to vicarious trauma - It is not available when the individual is themselves under threat - Being accused is fundamentally different from providing support

The individual who maintained professional composure during trauma-exposure work but dysregulates under accusation is not being inconsistent. They are exhibiting the difference between regulated professional performance and unregulated personal threat response.

2.9.5 Compound Vulnerability in the Criminal Justice System

The neurodivergent individual who enters the criminal justice system faces compound vulnerability: the system's baseline guilt-construction architecture operates upon a person whose authentic presentation systematically triggers the inverted credibility heuristics documented in Section 2.8.

This compound vulnerability operates at every procedural stage:

Pre-interrogation detention: - Autistic individuals: sensory overwhelm from fluorescent lighting, acoustic environment, unfamiliar textures; disruption of routine and predictability; social isolation removing regulatory supports - FND individuals: stress-exacerbation of symptoms; symptom variability interpreted as malingering; functional seizures or motor symptoms potentially misinterpreted as resistance or performance - PTSD/CPTSD individuals: retraumatization by institutional control; dissociative responses; hypervigilance exhaustion

Interrogation: - Autistic eye aversion coded as deception - FND symptom variability coded as faking - PTSD narrative fragmentation coded as inconsistency - Direct communication coded as hostility - Literal responses to pragmatically-loaded questions coded as evasion - Refusal to speculate coded as obstruction - Correction of interviewer errors coded as argumentativeness

Courtroom presentation: - Flat affect coded as lack of remorse - Direct communication coded as aggression or contempt - Non-normative emotional presentation coded as inappropriate - Detailed responses coded as over-rehearsed - Failure to perform expected distress coded as coldness - Visible distress coded as manipulation

Media framing: - Unusual presentation narrativised as "cold," "unfeeling," "robotic" - Absence of performed emotion narrativised as "showing no remorse" - Direct statements narrativised as "arrogant" or "defiant" - Literal accuracy narrativised as "splitting hairs" or "technicalities"

2.9.6 The Comorbidity Multiplication Effect

Neurodevelopmental and trauma-related conditions frequently co-occur. Autism commonly co-occurs with ADHD, anxiety, depression, PTSD, and FND. FND commonly co-occurs with chronic pain conditions, anxiety, and trauma histories. PTSD and CPTSD commonly co-occur with depression, anxiety, and functional somatic symptoms.

An individual presenting with multiple co-occurring conditions does not face additive credibility deficits; they face **multiplicative** deficits. Each condition produces its own set of credibility-impairing presentations, and the combination produces a presentation profile that deviates from neurotypical expectations across multiple dimensions simultaneously.

Consider an individual with: - Autism (direct communication, reduced eye contact, flat affect) - FND (variable symptoms, stress-responsive presentation) - CPTSD (fragmented trauma narrative, dissociative episodes, hypervigilance)

This individual will present with: - Direct communication that is coded as hostile - Reduced eye contact that is coded as evasive - Flat affect that is coded as cold or lacking remorse - Variable

symptoms that are coded as faking - Stress-responsive symptom changes that are coded as manipulation - Fragmented narrative that is coded as inconsistent - Dissociative episodes that are coded as evasion or performance - Hypervigilant presentation that is coded as guilty anxiety

There is no presentation available to this individual that would be interpreted as credible by observers applying standard credibility heuristics.

The individual's authentic presentation—which is the only presentation available to them—systematically triggers every credibility-reducing inference in the observer's repertoire.

2.9.7 Structural Illegibility of Innocence

The preceding analysis supports a specific conclusion: for certain individuals, **innocence is structurally illegible to the criminal justice system.**

The system's credibility assessment mechanisms are calibrated to a neurotypical baseline. Individuals who deviate from this baseline—whether through neurodevelopmental condition, trauma history, cultural background, or combination thereof—present in ways that the system's interpretive frameworks cannot accurately decode.

This is not a failure of individual observers to exercise appropriate care. It is a structural feature of a system whose assessment instruments were designed without reference to the diversity of human cognitive and communicative styles.

The innocent neurodivergent individual cannot present their innocence in a form the system can read. Their authentic presentation is decoded as deceptive. Their attempts to adapt their presentation are decoded as performed. Their awareness of the problem is decoded as manipulation.

The system processes their innocence as guilt, not because it fails to work, but because it works as designed.

2.9.8 The Complainant-to-Defendant Inversion

A specific pattern warrants documentation: the case in which an individual initially enters the criminal justice system as a complainant or victim, and subsequently becomes a defendant.

This pattern is observed in contexts including: - Domestic violence complaints that result in countercharges against the original complainant - Sexual assault complaints that result in charges against the complainant for false reporting - Complaints of institutional abuse that result in charges against the complainant - Whistleblower reports that result in charges against the reporter

For the neurodivergent individual, this inversion carries specific compound effects:

1. **Initial credibility deficit as complainant:** The neurodivergent complainant's presentation—direct communication, flat affect, fragmented trauma narrative, variable symptom presentation—produces initial credibility deficits that reduce the likelihood their complaint will be believed or pursued.
2. **Reinterpretation of complaint as false:** The same presentation features that reduced credibility as a complainant are subsequently reinterpreted as evidence that the complaint was fabricated. Flat affect becomes “cold calculation.” Direct communication becomes “rehearsed accusation.” Detailed recall becomes “scripted story.”

3. **Medical/psychiatric weaponisation:** Pre-existing diagnoses—autism, FND, PTSD, anxiety, depression—which explain the complainant’s presentation, are recharacterised as evidence of unreliability, instability, or propensity to fabricate.
4. **System familiarity as evidence:** Knowledge of legal processes, awareness of rights, or ability to articulate experiences clearly—all of which may reflect prior victimisation experiences or autistic information-gathering—may be interpreted as evidence of sophistication inconsistent with genuine victimhood.
5. **Trauma response as consciousness of guilt:** PTSD-related avoidance, hypervigilance, or emotional dysregulation during proceedings—which may reflect trauma from the original events being complained of—may be interpreted as consciousness of guilt regarding the allegedly false complaint.

The complainant-to-defendant inversion represents a particularly severe form of the structural illegibility problem: the individual’s authentic presentation of traumatic experience is decoded as evidence of fabrication, and their neurodivergent communication style provides the interpretive framework through which fabrication is inferred.

The same evidence that would, in a neurotypical complainant, be interpreted as consistent with genuine victimisation is, in the neurodivergent complainant, interpreted as consistent with false accusation.

2.9.6 Empirical Support

Key citations supporting the neurodivergent credibility gap:

1. **Lim, Young & Brewer (2021):** Autistic adults were rated as more deceptive and less credible than neurotypical speakers *when telling the truth*. N=1410 observers.
2. **Autistica (2024):** Survey of 394 police officers found only 37% had received autism training. Autism prevalence is 1-2% in general population vs 2-18% in forensic populations.
3. **Haworth et al. (2023):** Autism-typical behaviours (gaze aversion, flat affect, repetitive movement) are diagnostically indistinguishable from the nonverbal cues trained investigators use to identify deception.

The diagnostic criteria for autism overlap near-perfectly with the behavioural deception cues used by trained investigators (Global Deception Research Team, 2006): gaze aversion (#1 believed cue) and fidgeting (#2 believed cue).

An autistic person cannot present their truthful testimony without involuntarily performing the exact behaviours the system reads as deception. Their innocence is structurally illegible to the instrument.

2.10 The Behavioral Adaptation Feedback Loop

The Paradox of Learned Performance

The Signal Inversion Effect creates a paradox that compounds over time: **the person who learns that their authentic presentation is misread as deception may attempt to adopt the**

presentation style of “honest” behaviour—which is, empirically, the presentation style of liars.

This feedback loop operates as follows:

1. **Authentic behaviour is disbelieved** — the person who naturally avoids eye contact, hedges their statements, or reports confusion learns through repeated experience that these behaviours attract suspicion.
2. **Performed “honesty” is adopted** — in response, the person may consciously adopt the behaviours that folk psychology associates with honesty: steady eye contact, confident assertions, fluent narrative delivery.
3. **Performed “honesty” mimics deception** — but these performed behaviours are precisely the behaviours that empirical research associates with actual deception: rehearsed, performed, and strategically deployed.
4. **Detection becomes impossible** — the result is that the signal-to-noise ratio in credibility assessment approaches zero. Authentic communicators have been trained to present like deceivers, and actual deceivers already present this way.

The Trap for Honest People

The feedback loop creates a specific trap for honest people who have been repeatedly disbelieved:

- **First encounter:** Person presents authentically (eye aversion, hedging, fragmented recall). Is disbelieved.
- **Second encounter:** Person, having learned from first encounter, attempts to maintain eye contact and speak more fluently. This requires cognitive effort, which produces the very disfluencies and inconsistencies that authentic presentation would produce—but now layered over a performed baseline.
- **Third encounter:** Person, having had their performed presentation also questioned, may now present with defensive hypervigilance, excessive justification, or visible anxiety about being disbelieved—all of which trigger suspicion.

The trap is inescapable: authentic presentation triggers suspicion; performed presentation triggers suspicion; awareness of the problem itself triggers suspicion.

Why the System Cannot Self-Correct

The existence of the feedback loop has a specific implication for the reliability of credibility heuristics: **even if the baseline heuristics had some validity (which the evidence suggests they do not), the feedback loop would progressively destroy that validity over time.**

The system cannot self-correct through experience because:

- **True positives are overrepresented in feedback** — when an investigator identifies someone as deceptive and obtains a confession, the system treats this as confirmation regardless of whether the confession is true or false.
- **True negatives are invisible** — when an honest person successfully performs “honesty” and is believed, there is no feedback mechanism to indicate whether the belief was warranted.

- **False positives are often invisible** — when an honest person is misidentified as deceptive but does not confess, the system typically processes this as “investigation completed” rather than “error made.”
- **False negatives are impossible to identify** — when a deceptive person successfully presents as honest, there is (by definition) no subsequent identification of the error.

The feedback loop therefore operates within an institutional architecture that provides positive feedback for false confidence and no corrective feedback for error. **The system cannot learn because it is not structured to learn.**

2.11 Neuroimaging Evidence

A Note on Scientific Communication

This section includes neuroimaging data and brain images. We do so transparently, citing two reasons:

Reason 1: The Evidence Is Relevant

Neuroimaging directly supports the thesis that: - Stress impairs prefrontal cortex function (the neural substrate of rational decision-making) - Autistic individuals process social information differently at the neural level (not worse—differently) - The neural signatures of authentic cognitive effort are systematically misread

Reason 2: Neuroimaging Information Increases Persuasiveness

Weisberg et al. (2008) demonstrated that explanations containing neuroscience information are rated as more satisfying, even when that information is logically irrelevant.

We include brain images because the data shows they increase belief, and because the underlying evidence is genuinely supportive. This is not manipulation—it is strategic scientific communication.

The Prefrontal Cortex Under Stress

The prefrontal cortex (PFC)—particularly the dorsolateral PFC (dlPFC) and ventromedial PFC (vmPFC)—is the neural substrate of:

- Working memory
- Rational deliberation
- Impulse control
- Resistance to suggestion
- Long-term consequence evaluation

These are precisely the cognitive functions required to provide a voluntary, accurate statement during police interrogation.

Figure 2.3: Alert vs Stressed Brain States

Source: Arnsten (2015), *Nature Reviews Neuroscience*. PMC4816215. CC BY.

Stress PFC Figure

Figure 1: Stress PFC Figure

Threat Circuitry Figure 1

Figure 2: Threat Circuitry Figure 1

Caption: Changes in brain systems controlling behaviour under conditions of alert safety versus uncontrollable stress. Panel (a) shows the “Alert” state with PFC providing top-down regulation of thought, action, and emotion. Panel (b) shows the “Stressed” state where the PFC “goes off-line” and primitive brain circuits (amygdala, basal ganglia, PAG) take over with reflexive/habitual responding.

Key Quote from Arnsten:

“Exposure to uncontrollable stress rapidly evokes chemical changes in brain that impair the higher cognitive functions of the PFC while strengthening primitive brain reactions. This flip from reflective to reflexive brain state may have survival value when we are in danger, but it can be ruinous.”

Application: A person who has been arrested, stripped, searched, and confined in isolation presents with **measurably reduced prefrontal cortex function**. The neural substrate of rational decision-making has been chemically and functionally degraded.

The legal doctrine of “voluntariness” assumes a brain that no longer exists in the detained person.

Figure 2.4: Threat Regulatory Neurocircuitry

Source: Fenster et al. (2018), *Neuropsychopharmacology*. PMC8617299. CC BY.

Caption: Human brain anatomy highlighting regions involved in threat learning, extinction, avoidance, cognitive regulation, and contextual modulation. Key structures: vmPFC (ventromedial prefrontal cortex), dlPFC (dorsolateral prefrontal cortex), amygdala, hippocampus, dACC (dorsal anterior cingulate cortex).

Application: These are the neural structures that produce autistic social processing differences, PTSD-related threat dysregulation, and the credibility judgments that observers make. The differences are neurological, not behavioral choices.

Figure 2.5: Healthy vs PTSD Threat Circuits

Source: Fenster et al. (2018), *Neuropsychopharmacology*. PMC8617299. CC BY.

Caption: Panel (A) shows healthy threat circuitry with intact connectivity between dlPFC, vmPFC, hippocampus, and amygdala. Panel (B) shows PTSD-related threat circuitry where the dlPFC, vmPFC/IL, and hippocampus show impaired functioning with PTSD, whereas the amygdala and dACC/PL are enhanced.

Threat Circuitry Figure 2

Figure 3: Threat Circuitry Figure 2

Application: PTSD produces measurable changes in the brain circuits responsible for threat processing and memory. The fragmented recall, emotional dysregulation, and avoidant presentation characteristic of trauma are neurological symptoms, not credibility indicators.

Observers who interpret these presentations as deception are misreading brain states.

Neuroimaging Cannot Detect Deception Either

Despite two decades of research investment and the most sophisticated brain imaging technology available, **fMRI-based lie detection does not work reliably enough for forensic use.**

The Neuroimaging Evidence

Study / Review	Accuracy	Critical Finding
Meta-analysis (Nature Reviews Neuroscience)	75%	Best case in controlled laboratory conditions
UC Berkeley 2024 (PNAS)	79%	Confounded by selfishness — neural signatures of deception identical to self-interest
Mock crime paradigms	69% sensitivity	Low specificity — high false positive rates
Medial PFC region analysis	71%	Best single region; no region worked across all individuals
Applied Cognitive Psychology 2026	—	“Not suited for use as a lie detector”

Why fMRI Lie Detection Fails

The 2024 UC Berkeley study (Wills Neuroscience Institute / Haas) identified the fundamental problem:

“One reason it’s so hard to isolate signals of deception is that lying is a complex process that isn’t housed in a single part of the brain, and it’s challenging to separate activity linked to lying from that reflecting anxiety, self-interest, or other factors.”

The brain states associated with deception **overlap extensively** with the brain states associated with: - Being accused - Experiencing stress - Anxiety - Self-protective cognition - Trauma responses - Neurodivergent processing

The Hierarchy of Assessment Methods

If fMRI—which directly images brain activity—cannot reliably distinguish deception from innocence under stress, then behavioural heuristics—which attempt to infer internal states from external presentation—are necessarily less reliable still.

Method	What It Measures	Accuracy	Forensic Validity
fMRI	Direct brain activity	69-79%	Not reliable
Polygraph	Physiological arousal	65-70%	Not reliable

Method	What It Measures	Accuracy	Forensic Validity
Trained investigators	Behavioural observation	54%	Chance level
Untrained observers	Behavioural intuition	54%	Chance level

The most sophisticated neuroimaging cannot reliably detect deception. The behavioural heuristics used by the criminal justice system perform at chance level. There is no valid method currently employed.

Neuroimaging Confirmation of the Signal Inversion Effect

The neuroimaging research confirms the thesis argument: **there is no clean neural signal of deception that can be separated from the neural signatures of innocence under stress.**

An innocent person who is: - Anxious about being accused - Experiencing stress from interrogation - Engaged in self-protective cognition - Processing trauma-related content - Experiencing autistic or neurodivergent stress responses

will produce brain activation patterns **indistinguishable from** a guilty person who is lying.

The neural overlap between deception and innocent-under-stress is not a limitation that better technology will resolve. It reflects the fundamental reality that the cognitive processes involved in lying are also involved in being falsely accused, being stressed, being anxious, and being neurodivergent.

If the brain itself cannot be reliably read for deception, the claim that external behaviour can be read for deception is necessarily false.

Sources: UC Berkeley 2024, Applied Cognitive Psychology 2026, Nature Reviews Neuroscience, PMC Review 2024

PART II: THE ARCHITECTURE OF CONSTRUCTED GUILT

Chapter 3: The Body Before the Interview

Pre-Interrogation Detention as Pre-Punishment

3.1 Introduction: The Sequence Nobody Names

There is a sequence of events that precedes every police interrogation and that the literature on interrogation almost entirely ignores.

A person is stopped, typically without warning. Hands are placed on their body. They are told they are under arrest. They may be forced to the ground. Their arms are restrained behind them. They are moved—into a vehicle, into a building—without being asked. They are taken to a room and told to remove their clothing. Strangers examine their body, sometimes roughly. Their personal possessions—wallet, phone, keys, watch, the material coordinates of identity and social connection—are taken and bagged. They are given a paper or cloth garment, or nothing at all, and placed in a small room. The room is bare. It is typically painted in a colour specifically chosen to minimise stimulation. It smells of institutional cleaning agents and, often, of other people’s fear. The door is locked. Time passes.

Only after some portion of this sequence is complete does the “interview” begin.

This chapter argues that the pre-interview sequence is not preamble. It is the first stage of the guilt-construction process, and its effects on the subsequent linguistic exchange are so profound that the very concept of a voluntary statement, as deployed by Australian and other common law courts, is rendered empirically incoherent.

3.2 Legal Innocence and Physical Punishment: The Contradiction

The presumption of innocence is a foundational principle of common law criminal procedure. In Australia, it operates as:

- A constitutional implication (*Lange v Australian Broadcasting Corporation*, 1997)
- A common law presumption (*Woolmington v DPP*, 1935)
- A requirement of procedural fairness

Its doctrinal content is clear: a person is to be treated as innocent until guilt is established by a court. The corollary is equally clear: prior to such establishment, the state is not entitled to administer punishment.

The pre-interrogation detention regime described above is, by any substantive account, punishment.

It is experienced as punishment. It produces in the detained person the physiological and psychological states that punishment is designed to produce—fear, subordination, disorientation, and the acute awareness of institutional power over the body.

That it is not called punishment, that it is classified as “administrative procedure” or “custody management,” is itself a significant linguistic operation—one that will be returned to in Chapter 5 in the context of legislative language and legal fiction.

For present purposes, the point is structural: **the justice system inflicts an experience of punitive treatment on persons who are, by the system’s own formal declaration, innocent.**

This contradiction is not an oversight. It is not a gap between legal ideal and operational reality that could be closed by better regulation. It is built into the architecture. The conditions of pre-interrogation detention are maintained because they are functional—they serve the goal of producing statements.

3.3 The Neuroscience of Arrest: Stress, Cortisol, and the Interrogable Self

The physiological response to arrest—to the sudden, forceful, and involuntary seizure of the body by strangers—is a textbook activation of the hypothalamic-pituitary-adrenal (HPA) axis. Cortisol and adrenaline are released. Heart rate and blood pressure elevate. Glucose is mobilised.

These are adaptive responses to acute threat, mediated by the amygdala and coordinated across the autonomic nervous system (McEwen, 2007).

The prefrontal cortex—the region most centrally involved in executive function, working memory, rational deliberation, impulse control, and the inhibition of automatic responses—is **significantly impaired by acute stress** (Arnsten, 2009).

This is not a subtle effect. Arnsten’s (2009) review documented that even moderate stress levels produce measurable degradation in prefrontal cortical function, with consequent impairment of the cognitive capacities that underpin rational decision-making.

The detained person who is placed in an interrogation room following the arrest and detention sequence is, neurobiologically, not the same person who existed before the sequence began.

Their capacity to: - Reason carefully about their legal position - Weigh the long-term consequences of their words - Resist suggestion - Generate and maintain a complex account under pressure

has been materially diminished.

To these acute effects must be added the consequences of prolonged confinement:

- **Sleep deprivation** produces cognitive impairments comparable in magnitude to moderate alcohol intoxication (Harrison & Horne, 2000)
- **Sensory monotony** produces disorientation and heightened susceptibility to external structuring cues (Bexton et al., 1954; Zubek, 1969)

- **Social isolation** activates threat-related neural systems and heightens the perceived value of social acceptance and approval—including approval from the very investigators who will subsequently conduct the interview (Eisenberger et al., 2003)

The person who enters the interview room has been, through the pre-interview sequence, rendered more suggestible, more compliant, more desperate for social connection, and less capable of rational resistance than any baseline measure of their cognitive capacities would suggest.

3.4 Goffman and the Mortification of Self: Identity Stripped

Erving Goffman’s (1961) analysis of total institutions—prisons, asylums, military barracks, convents—identified a characteristic process he termed the “**mortification of self**”: a systematic dismantling of the identity resources through which individuals maintain a sense of continuous, autonomous selfhood in the social world.

The mortification process operates through specific institutional techniques:

- The removal of personal possessions and clothing
- The assignment of institutional garments or numbers
- The submission of the body to examination and surveillance
- The loss of control over basic activities such as movement and time

These techniques are not incidental to the institution’s operations; they are constitutive of them. They produce a subject whose prior identity has been suspended and who is thus available for institutional redefinition.

The pre-interrogation detention sequence maps precisely onto Goffman’s mortification process.

The confiscation of personal possessions removes the material anchors of social identity—the phone through which relationships are maintained, the wallet containing the cards and documents that attest to who one is in the institutional world, the watch that marks one’s place in the shared temporal framework of social life.

The removal and search of clothing submits the body—the most intimate site of self—to institutional examination.

The assignment of a cell number, a booking reference, a case file reduces the person to an administrative object.

By the time the interview begins, the person has been systematically repositioned, through physical and institutional processes, from a citizen to a suspect—from a subject of rights to an object of inquiry.

3.5 The Cell as Semiotic Environment

The physical environment of the holding cell is not accidental. Institutional design research has established that environmental features systematically affect the psychological states of those who inhabit them (Evans, 2003).

The characteristic features of police holding cells—limited space, minimal natural light, plain or specifically coloured walls, hard surfaces, the absence of any material that would allow productive

activity—are, from an environmental psychology perspective, a collection of stressors.

The colours used in detention facilities—the grey-beige-green spectrum characteristic of institutional interiors—is not the result of aesthetic indifference. Institutional colour choice in carceral environments has historically been informed by the goal of minimising stimulation and maintaining order, with the consequence of producing environments experienced by occupants as oppressive, disorienting, and dehumanising (Kwallek et al., 1996).

The cell is a designed environment, and its design serves the institution’s goals.

From a semiotic perspective, the cell communicates something specific to its occupant:

- Its **bareness** says: you have no resources here
- Its **locks** say: your body is no longer your own
- Its **uniformity** says: you are interchangeable with every other person who has occupied this space
- Its **smell**—of disinfectant, of fear, of the accumulated traces of prior occupants—says: this is where people like you end up

The cell, prior to any question being asked, has already begun the work of positioning the person as a certain kind of subject: guilty, contained, available for interrogation.

3.6 Voluntariness Revisited: The Legal Fiction of the Free Agent

Australian courts apply the common law principle that a confession is admissible only if it was made voluntarily—that is, without threats, inducements, or oppressive conduct that overbore the will of the accused (*R v Lee*, 1950; Uniform Evidence Acts, s 84).

These doctrinal protections rest on a model of the person as a rational agent capable of making free choices about whether to speak, and capable of accurately reporting facts when they do.

The empirical literature reviewed in this chapter establishes that this model is false for the person who has been through the pre-interrogation sequence.

The model assumes a baseline of cognitive capacity, emotional regulation, and autonomous agency that the detention process is specifically designed to degrade.

Courts routinely admit confessions obtained in these conditions as “voluntary” because: - No explicit threat was made - No promise was offered - No physical violence was applied in the interview room itself

The coercive architecture of the hours preceding the interview is treated as legally irrelevant—a prior administrative matter, not part of the interrogation.

This doctrinal gap is not a failure of judicial attention. It is the predictable consequence of a legal framework that defines voluntariness in formal rather than substantive terms.

A formally voluntary statement—one made without explicit compulsion—can be, in the substantive sense relevant to the purposes of the voluntariness doctrine, **profoundly involuntary**.

The person who confesses to a crime they did not commit after eight hours in a holding cell, in an acute stress state, with impaired prefrontal function, heightened suggestibility, and a desperate need for the social approval of the only human beings in their immediate environment, has made a “free choice” only in the most impoverished sense of that phrase.

The law's insistence that this constitutes voluntary statement is itself a speech act—one that constitutes a legal reality that diverges fundamentally from the experiential and neurobiological reality of the person concerned.

3.7 Chapter Summary

This chapter has argued that the pre-interrogation detention sequence constitutes a first stage of guilt construction that has been systematically neglected in both the legal doctrine of voluntariness and the scholarly literature on interrogation.

The sequence—arrest, bodily search and exposure, confiscation of identity materials, confinement in a designed environment of minimal resource and maximal institutional control—produces, through neurobiological, psychological, and semiotic mechanisms, **a subject who is measurably less capable of free and rational communication than they were prior to detention.**

The legal concept of voluntariness, which governs the admissibility of the statements subsequently obtained, does not account for these effects.

The conclusion is not that reforms to the voluntariness doctrine would solve the problem—though such reforms would be an improvement—but that the conditions of pre-interrogation detention are **structurally functional**: they serve the system's interest in producing statements, and their preservation reflects that interest.

Chapter 4: The Interview Room

The Reid Technique and the Architecture of Manufactured Guilt

4.1 Introduction: The Room You Cannot Leave

You are seated in a small room. Across from you sits a trained investigator. Between you is a table. The door, behind you and to one side, is closed. You have been told you are free to leave, but this statement—required by law to preserve the legal fiction of voluntariness—has been delivered in a context that makes its sincerity implausible.

You were brought to this room from a cell. You do not have your phone. You do not have your wallet. You are wearing your own clothes, or paper clothes, depending on jurisdiction and procedure. You have not eaten recently, or slept well, or spoken to a lawyer for as long as the system could lawfully delay that access.

The investigator leans forward. They are friendly. They understand. They are just trying to get your side of the story.

This chapter examines what happens next.

4.2 The Reid Technique: Architecture of Presumption

The focus is on the dominant interrogation methodology in use across common law jurisdictions: the **Reid Technique**, developed by John E. Reid and Fred Inbau in the mid-twentieth century and codified in the manual *Criminal Interrogation and Confessions* (Inbau et al., 2013).

The Reid Technique is not merely an interrogation method. **It is a system for the linguistic and behavioural construction of guilt.**

The Technique proceeds in two formal stages:

1. **The Behavioural Analysis Interview (BAI):** A non-accusatory pre-interrogation interview designed to assess deception through behavioural cues
2. **The Nine-Step Interrogation:** An explicitly accusatory process designed to overcome the suspect's resistance to confession

The structure itself encodes a presumption: the BAI assesses whether the person is lying; the interrogation is then used to extract a confession from those assessed as deceptive.

The possibility that the BAI assessment is wrong—that an innocent person has been assessed as deceptive—is not procedurally accommodated. There is no step in the Technique designed to consider this possibility; the entire structure points forward to confession.

4.3 The BAI's Empirical Failure

The BAI's assessment methodology has been subjected to substantial empirical scrutiny, and the findings are consistently damaging.

The Technique's claim that trained investigators can reliably detect deception through behavioural observation—noting gaze aversion, postural shifts, self-grooming behaviours, verbal hedging, and similar cues—**is not supported by the research evidence.**

A meta-analysis by Vrij (2008) found that trained investigators perform at rates only marginally above chance, with a mean accuracy rate of approximately 54% against a chance baseline of 50%.

Crucially, investigators tend to be more confident in their assessments than their accuracy warrants—a pattern that Kassin et al. (2005) termed “**confidence without accuracy.**”

A technique that produces high confidence and low accuracy is, in the context of an interrogation that will proceed to extraction of confession based on that assessment, precisely the architecture required to generate false confessions from innocent people.

4.4 The Innocent Stress Response

The behaviours that the Reid Technique treats as indicators of deception—gaze aversion, postural shifting, delayed responses, verbal qualifications—are also among the most reliably produced responses to acute stress **in innocent people** (Vrij et al., 2006).

An innocent person, suddenly accused of a serious offence in a formal institutional setting by a trained investigator after hours of pre-interrogation detention, will typically be frightened.

Frightened people: - Avoid eye contact - Shift posture - Hesitate before speaking - Qualify their statements, because they are aware that the stakes are high and they want to be accurate

The Reid-trained investigator observes these behaviours and records: deceptive.

The innocent person's innocence is being used as evidence against them.

4.5 The Nine Steps: A Linguistic Architecture of Guilt

Step 1: Positive Confrontation

The interrogation begins with the investigator directly telling the suspect that the evidence conclusively establishes their guilt. This is a direct statement, made with certainty and authority, **regardless of the actual state of the evidence.**

From a speech act perspective, this is not a description of evidentiary reality; it is a performative—an attempt to constitute, through confident assertion, a social reality in which the suspect's guilt is already established.

The innocent suspect who responds with confusion, distress, or vigorous denial is, within the Technique's interpretive framework, demonstrating resistance rather than innocence.

Step 2: Theme Development

The investigator presents a morally minimising narrative—a “theme”—that offers the suspect a face-saving account of why they committed the act.

Themes are typically constructed to shift moral responsibility: the offence was understandable given the circumstances; it was a momentary lapse; the victim was partly responsible; anyone might have done the same thing.

The function of theme development is to reduce the psychological cost of confession. Kassin and McNall (1991) demonstrated that **minimisation significantly increases confession rates—including, critically, among innocent suspects.**

The mechanism is straightforward: the suspect is offered a choice between a bearable story (the theme) and an unbearable situation (continued interrogation in conditions of pre-punitive detention). Many innocent people choose the bearable story.

Steps 3–6: Managing Denial and Overcoming Objections

The Technique instructs investigators to **actively prevent the suspect from completing their denials**—to interrupt, to redirect, to maintain forward momentum toward confession.

This is presented in the training literature as preventing the suspect from “reinforcing” their denials. From a conversational and linguistic perspective, it is the **systematic suppression of the suspect’s counter-narrative.**

The suspect who attempts to tell their own story—to construct their own account of where they were, what they did, and why the investigator’s account is wrong—is procedurally blocked from doing so.

The conversation is a monologue with interruptions, not a dialogue. The investigator’s narrative is the only one with structural permission to develop.

The implications for an innocent suspect are severe. The innocent person has, by definition, an alternative account of events. **The Reid Technique’s procedural suppression of their denials prevents them from presenting this account with the coherence and completeness it requires to be persuasive.**

Step 7: The Alternative Question

The alternative question is among the most analytically significant elements of the Technique.

The investigator presents the suspect with a choice between two versions of how the act occurred—one more morally serious, one less—and invites the suspect to choose:

“Did you plan this out, or did it just happen in the heat of the moment?” “Was this about money, or was it something personal?”

Both options presuppose the act.

The question is framed as an inquiry into motivation and circumstance, but its deep structure is a **presuppositional trap**: answering either alternative constitutes an admission of the act itself.

Austin’s (1962) analysis of presupposition is precisely applicable here. The alternative question embeds guilt as a presupposition and then offers a choice that cannot be taken without accepting that presupposition.

Courts have held that explicit promises of leniency render confessions involuntary; the alternative question achieves the same psychological effect through implication rather than explicit statement,

and thus typically survives voluntariness analysis (Leo, 2008).

The linguistic indirection of the alternative question is not an incidental feature of its design. It is the mechanism that allows the Technique to produce the psychological effects of explicit coercion while maintaining the legal appearance of voluntary interaction.

Steps 8–9: Oral and Written Confession

Once an admission has been obtained—typically through the alternative question—steps eight and nine involve expanding the oral admission into a detailed narrative and then reducing that narrative to a written statement.

This sequence introduces a further layer of linguistic construction. The investigator who assists in elaborating the confession narrative, or who drafts the written statement, **introduces their own language, their own emphases, their own causal and temporal frameworks.**

Research on contamination of confession evidence has documented numerous cases in which written confessions contain details that the confessor could not have known—details that could only have been introduced by investigators who possessed information about the crime that had not been disclosed to the suspect (Leo & Ofshe, 1998; Kassin et al., 2010).

The written confession, presented in court as the suspect’s own account, is often substantially the investigator’s account—a narrative constructed by institutional power and attributed to the person whose guilt it is then used to prove.

4.6 The False Confession Literature: Guilt Constructed From Innocence

The empirical literature on false confessions provides the most direct evidence for the thesis that guilty narratives can be constructed from innocent behaviour.

Gross et al. (2005), in an analysis of 340 exonerations in the United States, found that false confessions were documented in approximately 15% of cases.

Brandon Garrett’s (2011) analysis of DNA exonerations found that **27 of the first 250 exonerees—all of whom were factually innocent, established by post-conviction DNA testing—had falsely confessed.**

These are not people who confessed under duress and were later shown to have actually committed the offence. They are people who confessed to crimes **they demonstrably did not commit.**

Saul Kassin’s (2017) review identified three categories of false confession:

1. **Voluntary:** Occur in the absence of external pressure
2. **Compliant:** Occur when a suspect conforms to perceived demand despite privately knowing themselves to be innocent
3. **Internalised:** Occur when the suspect comes to believe, as a consequence of the interrogation process, that they actually committed the act

The internalised false confession is the most disturbing category theoretically: it represents not merely a person who lied to end the interrogation, but **a person whose memory of events was rewritten by the interrogation process itself.**

4.7 Any Behaviour, Any Person: The Architecture of Universal Guilt

The thesis advanced in Chapter 1—that the architecture of the interrogation room is structured such that any behaviour can be narrated as evidence of guilt—is now demonstrable in its full specificity:

- The Reid Technique’s BAI treats the behavioural responses of innocence under stress as indicators of deception
- The nine-step interrogation’s procedural suppression of denial prevents the innocent person from establishing their counter-narrative
- The alternative question embeds guilt as a presupposition and triggers a speech act theory trap that constitutes admission through apparent choice
- Theme development reduces the psychological cost of false confession below the cost of continued resistance
- The conditions produced by pre-interrogation detention ensure that the person entering the room is already operating at a cognitive and emotional deficit that the Technique is designed to exploit

The result is a system in which innocent behaviour—*anxiety, confusion, gaze aversion, postural shifting, the desire for the interrogation to end, the desire for social approval, the subjective plausibility of a morally minimised account*—is the raw material from which guilt is manufactured.

This is not an accusation that investigators are dishonest. Many investigators genuinely believe in the guilt of the person they are interviewing, because the BAI has told them so and because the Technique’s training has provided them with a framework in which all resistance confirms guilt and all compliance confirms guilt.

It is an analysis of **a system that is architecturally designed to produce confessions—and that produces false confessions as a predictable, quantifiable, structurally necessary output.**

4.8 The PEACE Model: A Counter-Institutional Framework

Alternative interrogation frameworks exist. The **PEACE model** (Preparation and Planning; Engage and Explain; Account; Closure; Evaluate), developed in the United Kingdom following a series of high-profile false conviction cases in the 1990s (Clarke & Milne, 2001), proceeds from an explicitly non-accusatory framework.

PEACE interrogators are trained to **elicit and test accounts** rather than to overcome resistance to confession. The model is grounded in cognitive interview techniques (Fisher & Geiselman, 1992) that seek to maximise the accuracy and completeness of recall rather than the likelihood of confession.

Research comparing PEACE and Reid-based approaches has consistently found that PEACE-style interviewing produces: - More accurate information - Fewer false confessions - No significant reduction in true confession rates (Meissner et al., 2014; Walsh & Bull, 2012)

Australian jurisdictions, including the Australian Federal Police and most state forces, have formally adopted PEACE-aligned frameworks in their investigative interviewing guidelines.

The persistence of Reid-influenced practices at the operational level, despite formal

policy alignment with PEACE principles, is itself a significant institutional phenomenon—one that speaks to the depth of the system’s interest in confession production over truth production.

4.9 Neurodivergent Vulnerability in Interrogation

The neurodivergent individual in the interrogation room faces a compound vulnerability:

Autistic individuals may: - Fail to recognise the pragmatic force of indirect speech acts - Respond literally to presuppositional questions - Not understand that “cooperating” is being coded as confession - Present with reduced eye contact, flat affect, and direct speech—all coded as deceptive

Individuals with FND may: - Have symptoms exacerbated by the stress of interrogation - Present with variable symptom intensity that is misread as manipulation - Be unable to control bodily responses that investigators interpret as behavioural cues

Individuals with PTSD/CPTSD may: - Provide fragmented narratives that are coded as inconsistent - Dissociate under stress, producing apparent gaps in recall - Avoid discussion of traumatic details, coded as evasion

The Reid Technique’s interpretive framework contains no mechanism for distinguishing neurodivergent presentation from deceptive presentation. The same behaviours are coded identically regardless of their actual cause.

4.10 Chapter Summary

This chapter has examined the Reid Technique as a system for the linguistic and behavioural construction of guilt in the interrogation room.

The Technique’s architecture—presuming guilt, suppressing denial, embedding guilt as presupposition, reducing the psychological cost of false confession, and exploiting the suggestibility produced by pre-interrogation detention—is structured to produce confessions from the persons subjected to it, regardless of their actual culpability.

The empirical false confession literature documents the output of this system. The comparison with the PEACE model establishes that the Technique’s guilt-production outcomes are not inevitable features of interrogation but choices embodied in a specific institutional architecture.

Chapter 5: Legislative Language and Legal Fiction

How Law Writes the World It Pretends to Describe

5.1 Introduction: The Word That Creates the Crime

Before there can be a guilty person, there must be a crime. And before there can be a crime, there must be a word. The sequence is rarely stated this plainly in criminal law scholarship, because criminal law scholarship tends to operate within the assumption that statutes describe pre-existing wrongs rather than create new ones. This chapter argues the reverse: that criminal legislation is a performative rather than descriptive enterprise, that the categories it establishes are not reflections of natural moral facts but political constructions with specific historical and ideological genealogies, and that the key terms on which guilt or innocence turns — 'intent,' 'consent,' 'reasonable,' 'recklessness,' 'dishonesty' — are semantically unstable in ways the law persistently conceals. The person convicted of an offence defined by these terms has not been found to have done a thing that was always and obviously wrong. They have been found to have satisfied the conditions of a definitional network constructed by legislators, interpreted by judges, and applied by juries who received the key terms with minimal instruction and without any of the philosophical context that would reveal how much work those terms are quietly doing.

5.2 Legislation as Performative Speech Act

When a parliament enacts a statute criminalising an act, it does not discover that the act was already criminal. It makes it criminal. This is Austin's (1962) performative utterance at the level of sovereign institutional power: the legislation is the speech act that constitutes the offence. Prior to the statute, the act existed in the world. After the statute, the act is a crime. The word — the legislative text — is the mechanism of transformation. This is so obvious at the level of legislative creation that it tends to escape notice; but its implications extend into every subsequent stage of criminal proceedings. If criminal categories are created rather than discovered, then they are contingent. Cannabis possession has been criminal and non-criminal within the same jurisdictions across the span of a few decades. Marital rape was legally impossible in Australia until the late twentieth century — not because it did not occur, but because the legal category that would have made it criminal did not exist. Homosexual acts between consenting adults were criminal offences in all Australian states until the 1970s and 1980s. The content of the criminal law is not a map of natural harms. It is a historically specific construction, produced by specific institutional actors with specific interests, reflecting specific political settlements — and those political settlements are themselves the product of power relations that are, as Foucault (1977) argued, neither neutral

nor stable. The constructed character of criminal categories is most visible at the margins, where definitional contests occur explicitly. But the same construction is present at the centre, in the definitions of offences so familiar that their contingency has become invisible. 'Murder' appears to name a natural moral category — the wrongful killing of a human being — but the legal definition immediately introduces contestable constructions: 'unlawful' killing (as opposed to lawful killings in warfare, policing, or capital punishment), 'with intent to kill or cause grievous bodily harm' (introducing a mental state element that requires inference from observable behaviour), and 'of a human being' (a category whose edges, at the beginning and end of life, are legally contested). The word 'murder' appears to describe a thing that obviously happened. It actually denotes the satisfaction of a complex definitional matrix, each element of which is a legal construction.

5.3 The Reasonable Person Standard: Whose Construction?

Among the most consequential legal fictions in the common law criminal system is the figure of the 'reasonable person.' The reasonable person appears throughout criminal law: as the standard against which provocation is assessed, as the benchmark for self-defence, as the reference point for negligence and recklessness, and as the implicit model against which a defendant's conduct is evaluated in innumerable contexts. The reasonable person is presented as an objective standard — a neutral reference point that transcends the particular perspectives of the parties before the court. The reasonable person has never existed. They are a normative construct — a legal fiction whose function is to translate contested value judgments into the apparently neutral language of objective fact. When a court asks whether a 'reasonable person' would have been provoked to kill in the circumstances before it, or whether a 'reasonable person' would have perceived a threat justifying self-defensive force, it is asking: does this defendant's response fall within the range of responses that this legal community is prepared to treat as excusable or justified? The reasonable person is the vehicle through which that collective normative judgment is expressed as though it were a factual observation. The political content of this fictional figure has been extensively documented in feminist legal scholarship. Schneider (1992) argued that the reasonable person standard historically encoded a specifically male response to perceived threat, systematically disadvantaging women who had used defensive force against intimate partners whose violence was cumulative and contextual rather than immediate and acute. The battered woman who kills a sleeping abuser does not act as a 'reasonable person' would act in the moment of attack — because the attack, for her, is not happening in the moment. It has been happening for years. The standard's failure to accommodate this reality was not a neutral technical limitation; it was the reflection of a normative framework constructed from a specifically male experiential baseline (Stubbs & Tolmie, 1999). Australian jurisdictions have partially responded to this critique through legislative reforms to self-defence provisions (Model Criminal Code, s 10.4; Crimes Act 1900 (NSW), s 418), but the reasonable person remains operative as a standard in numerous contexts, carrying its historically loaded normative content into each application. The reasonable person also encodes racial and class assumptions that have been documented with particular acuity in the critical race legal scholarship. The question of whether conduct was 'reasonable' is answered by a judicial system and jury pool that are, in Australian terms, disproportionately white and middle-class. Research on implicit bias in legal decision-making (Rachlinski et al., 2009; Richardson & Goff, 2012) has established that the application of apparently neutral standards is systematically affected by the racial identity of defendants and victims. The reasonable person is not a universal figure. They are a culturally specific figure whose particular rationality is validated by the institutional structures that deploy them.

5.4 Intent, Recklessness, and the Inference Machine

Criminal liability in most common law jurisdictions requires proof of a mental element — the *mens rea* — alongside proof of the prohibited act. The mental element varies by offence: murder requires intent or knowledge; manslaughter may be established by recklessness; some offences require only negligence. This requirement of mental state is presented as a fundamental safeguard — the principle that people should be punished for what they meant to do, not merely for what happened as a result of their actions. It is, in this framing, the law’s recognition of the moral distinction between the deliberate wrongdoer and the unfortunate accident. The safeguard is, however, largely nominal. Mental states are not observable. Intent, knowledge, recklessness, and negligence are all inferences drawn from observable behaviour and circumstance. The jury that is asked to find whether a defendant intended to kill is not being asked to report an observation. They are being asked to make an inference — to take the observable facts (the act, the circumstances, the defendant’s statements and conduct before and after) and to construct a narrative about what was happening inside the defendant’s head. This inference is not a neutral cognitive operation. It is shaped by the same cognitive schemas, cultural assumptions, and narrative frameworks that govern all human judgment, and it is performed by people who have already heard the prosecution’s account of events and been subjected to the full rhetorical architecture of courtroom presentation. Research on attribution theory and criminal judgment has consistently demonstrated that people apply different causal and intentional explanations to the same behaviour depending on the social identity of the actor (Graham & Lowery, 2004; Goff et al., 2014). The behaviour that is read as impulsive and reckless when attributed to one social group is read as calculated and intentional when attributed to another. The mental element requirement, far from providing a safeguard against the arbitrary attribution of guilt, provides a site at which social biases are institutionally reproduced through the language of objective inquiry. The jury is not asked ‘did this person intend to kill?’ — a question about the actual contents of another human mind. They are asked a question they cannot answer, and they answer it using whatever cognitive resources they have available — which include their cultural schemas for who is the kind of person who does this kind of thing.

5.5 Consent, Sexual Assault, and the Definitional Contest

Few legal terms carry as much contested political weight as ‘consent,’ particularly in the context of sexual offences. The legal definition of consent — and the evidentiary rules governing how it is established or negated — has been the site of sustained feminist legal reform advocacy across common law jurisdictions over the past five decades. In Australia, the Crimes Act 1900 (NSW) now defines consent as free and voluntary agreement (s 61HE), and affirmative consent models have been adopted in several jurisdictions. These reforms represent genuine improvements. They also illustrate, with unusual clarity, the constructed character of legal definitions: the meaning of ‘consent’ for the purposes of criminal law is whatever the legislature decides it is, and that decision is a political one. The point is not that consent is meaningless or that its presence or absence is always contested. The point is that ‘consent’ in criminal law is a legal term of art — a definition that operates within a specific institutional context, interpreted according to specific evidentiary rules, applied by decision-makers with specific cognitive and cultural frameworks — and that the gap between the legal concept and the lived experience it purports to address is substantial. Research on juror decision-making in sexual assault cases has documented persistent myth-consistent reasoning — the application of false beliefs about how ‘genuine’ victims behave, how perpetrators present, and what circumstances are consistent with non-consent — that operates in the gap between the legal definition and jurors’ interpretive frameworks (Larcombe, 2002; Temkin & Krah?, 2008). The

legal definition of consent, however well-drafted, is applied by human beings whose understanding of what consent looks like in practice is shaped by cultural narratives that the law cannot simply override through definitional revision.

5.6 Identity as Legal Category: ‘You Are Australian’

The theoretical framework established in Chapter 2 introduced the observation that identity categories in law are fictions — not in the pejorative sense of falsehoods, but in the technical sense of constructed categories that produce real effects while describing nothing in the natural world. The category ‘Australian’ — as a legal designation — illustrates this with particular clarity. From a zoological perspective, a person described as ‘Australian’ is a member of the species *Homo sapiens*, currently residing on the landmass that was named ‘Australia’ by European colonisers in the nineteenth century. Nothing in the biological, geological, or cosmological description of this person and their location requires the category ‘Australian.’ The category is a legal and political construction — one that was assembled through colonial dispossession, federation, immigration legislation, citizenship law, and passport administration. It has no natural referent. Yet the legal consequences of being ‘Australian’ versus ‘not Australian’ are profound: they determine which criminal jurisdiction has authority over one’s body, which rights one can assert, which protections one can claim, and which obligations one bears. The construction does real, material work — coercive work — in the world. Criminal law operates entirely within this constructed landscape. Every category deployed in criminal proceedings — ‘defendant,’ ‘victim,’ ‘witness,’ ‘offender,’ ‘juvenile,’ ‘recidivist’ — is a legal fiction of the same type. Each category carries a definitional network that activates specific institutional procedures, rights, obligations, and presumptions. The person who is classified as a ‘juvenile offender’ is not merely described by this term; they are constituted as a particular kind of legal subject who will be processed through a different institutional system, subject to different dispositions, and described in different records than the person classified as an ‘adult offender.’ The word does not follow the person. The word precedes them, and the person is fitted to it.

5.7 The Grammar of Guilt: How Charges Construct Narrative

When a person is charged with a criminal offence, the charge document is not merely an administrative record. It is a narrative — one that constructs a specific account of events, attributes specific intentions and actions to a specific person, and frames that account within the definitional structure of a specific statutory provision. The words of the charge are not neutral. They are selected from a range of possible characterisations, each of which would produce different legal consequences, and the selection is made by a prosecution service that has an institutional interest in conviction. The prosecutorial discretion to charge — to choose among available offences, to decide whether to proceed, to select which counts to include and which to leave out — is a site of substantial constructed power that receives insufficient critical attention. The person charged with aggravated assault rather than common assault, with murder rather than manslaughter, with supply rather than possession, has had a narrative imposed on their conduct that is not uniquely determined by the facts. It is one of several possible narrativisations, chosen by an institutional actor, and that choice will shape everything that follows: the available defences, the maximum penalty, the stigma, and the institutional pathway through which the case proceeds. The charge is the first act of courtroom narrative construction — and it takes place before the court convenes.

5.8 Chapter Summary

This chapter has argued that criminal legislation does not describe a pre-existing moral landscape but constructs the categories through which guilt becomes possible. The performative character of legislative language — the fact that statutes create crimes rather than discovering them — means that every criminal category is contingent, politically constructed, and historically specific. The key terms through which guilt is determined — ‘intent,’ ‘consent,’ ‘reasonable,’ ‘recklessness’ — are semantically unstable in ways that the legal system persistently conceals. The reasonable person standard encodes a normatively loaded fiction as an objective reference point. Mental state requirements demand inferences that are shaped by the same cognitive biases documented in the preceding chapters. Consent definitions operate in a gap between legal language and juror interpretation that legislative drafting cannot close. Identity categories construct the legal subject before the legal subject enters the courtroom. And the charge document — the prosecution’s first narrative act — frames the entire proceeding before a word of evidence is heard.

Chapter 6: The Courtroom as Construction Site

Cross-Examination, Memory, and the Architecture of Adversarial Truth

6.1 Introduction: The Official Version

The courtroom is, in the popular imagination, the site where truth is established. Two competing accounts are presented, evidence is tested, witnesses are examined, and a neutral decision-maker — judge or jury — determines which account is true. This model of the adversarial trial as truth-finding mechanism is foundational to common law legal culture; it is the premise of the presumption of innocence, the basis of the right to cross-examine, and the legitimating narrative of the verdict. This chapter argues that the model is false — not in the sense that trials always produce incorrect verdicts, but in the deeper sense that the adversarial courtroom is not designed as a truth-finding mechanism and does not function as one. It is designed as a narrative contest, and what it produces is not truth but the official version: the account that has been most successfully constructed within the specific rhetorical and institutional constraints of the trial process.

6.2 The Adversarial System as Narrative Competition

The adversarial trial proceeds on the assumption that truth is best approximated through the clash of competing partisan accounts, each subjected to rigorous challenge by the opposing party, with a neutral decision-maker determining the outcome. This is a coherent epistemological hypothesis. It is also one for which there is remarkably little empirical support. Research comparing adversarial and inquisitorial trial systems — the two dominant models across common law and civil law jurisdictions respectively — has not produced consistent evidence that adversarial systems generate more accurate verdicts (Damaska, 1997; Findley & Scott, 2006). What the adversarial system does demonstrably produce is a particular kind of narrative contest in which rhetorical skill, the strategic management of evidence, and the exploitation of juror psychology are systematically rewarded. Walter Fisher's (1984) narrative paradigm provides a useful analytical framework here. Fisher argued that human beings are fundamentally narrative creatures — that we evaluate claims not primarily through logic and evidence but through narrative rationality: the assessment of whether a story hangs together coherently (narrative probability) and whether it resonates with our experience of how the world works (narrative fidelity). In Fisher's framework, the most compelling argument is not the most logically sound but the most narratively coherent and resonant. The adversarial trial, in this analysis, is not a logical contest but a narrative competition — and it will be won by the party whose story best satisfies the jurors' narrative rationality, regardless of its correspondence to what

actually happened. This analysis has direct implications for the construction of guilt. The prosecution, which opens first, has the structural advantage of establishing the primary narrative — the framework within which all subsequent evidence will be interpreted. Research on primacy effects in juror decision-making (Penrod & Hastie, 1979; Furnham, 1986) has consistently established that information presented first has disproportionate influence on final judgments. The prosecution’s opening is not merely a preview of evidence; it is the installation of a narrative framework that will operate as a cognitive filter on everything the jury subsequently hears. Evidence that is consistent with the prosecution’s narrative will be remembered and weighted; evidence that is inconsistent will be reinterpreted, minimised, or forgotten. The jury that reaches the end of a criminal trial has not evaluated two competing accounts from a neutral baseline. They have evaluated the defence account from within the cognitive structure that the prosecution’s opening constructed.

6.3 Cross-Examination: Language as Memory Surgery

Cross-examination is the adversarial system’s primary mechanism for testing evidence. A witness who gives evidence in chief is then subjected to questioning by the opposing party, whose goal is to undermine the credibility, accuracy, or completeness of that evidence. In principle, this is a truth-seeking operation: the cross-examiner challenges the witness’s account and exposes its weaknesses, leaving the jury better positioned to evaluate its accuracy. In practice, cross-examination as routinely conducted is not a truth-seeking operation but a narrative-construction operation — one whose principal tools are linguistic, and whose effects operate directly on the witness’s memory. The foundational empirical research here is Elizabeth Loftus’s work on the misinformation effect, introduced in Chapter 2. Loftus and Palmer’s (1974) demonstration that the verb used in a question about a car accident affected both speed estimates and subsequent memory of physical details established that post-event language does not merely elicit memory — it alters it. Subsequent decades of research have established the robustness and generality of this effect across a wide range of events, populations, and types of misinformation (Loftus, 2005). The relevance to cross-examination is direct and devastating. The cross-examining attorney who asks ‘When the defendant ran out of the building, what did he do next?’ has embedded the presupposition that the defendant ran — and research suggests that witnesses who answer this question, even those who initially reported that the defendant walked, will subsequently be more likely to remember the defendant as having run. The cross-examination has not revealed the truth. It has edited the memory of the witness.

The scope of the misinformation effect has been documented across more than four decades of research. Key studies include:

Loftus & Palmer (1974): Participants watched a film of a car accident and were asked “How fast were the cars going when they [smashed / collided / bumped / hit / contacted] each other?” The verb used in the question produced significantly different speed estimates: ‘smashed’ produced an average estimate of 40.5 mph; ‘contacted’ produced 31.8 mph. One week later, participants who had received the ‘smashed’ condition were more than twice as likely to report seeing broken glass (there was none in the film). The question did not merely influence the answer. It altered the memory.

Loftus, Miller & Burns (1978): Participants who were exposed to misleading post-event information about a traffic sign (yield sign vs. stop sign) subsequently ‘remembered’ seeing the sign described in the misleading question, even when they had originally correctly identified the actual sign. The misinformation did not merely coexist with the original memory. It replaced it.

Loftus (1993): In a review of the field, Loftus documented that approximately 22% of participants across studies accepted false information as their own memory after exposure to misleading post-event information. The effect was robust across populations, types of events, and forms of misinformation.

Hyman, Husband & Billings (1995): Participants who were told by family members that a false childhood event had occurred (e.g., being lost in a shopping centre) came to ‘remember’ the event with vivid detail — including sensory details, emotional responses, and narrative structure — after as few as three interviews. Approximately 25% of participants created full false memories of events that never occurred.

Application to cross-examination: The cross-examining attorney has available every technique documented in this literature: presuppositional questions (“When you saw the defendant running...”), repetition with variation (repeating a question with slightly altered framing until the witness’s memory conforms), confirmation bias exploitation (“So you would agree that...”), and authority-based suggestibility (the social authority of the courtroom, the robed judge, and the assertive questioner all increase the witness’s susceptibility to suggestion).

The witness who leaves the stand after sustained cross-examination has had their memory systematically altered by a process that the legal system presents as truth-seeking. The cross-examination was not a test of the witness’s memory. It was an intervention in it.

The scope of the misinformation effect has been documented across more than four decades of research. Key studies include:

Loftus & Palmer (1974): Participants watched a film of a car accident and were asked “How fast were the cars going when they [smashed / collided / bumped / hit / contacted] each other?” The verb used in the question produced significantly different speed estimates: ‘smashed’ produced an average estimate of 40.5 mph; ‘contacted’ produced 31.8 mph. One week later, participants who had received the ‘smashed’ condition were more than twice as likely to report seeing broken glass (there was none in the film). The question did not merely influence the answer. It altered the memory.

Loftus, Miller & Burns (1978): Participants who were exposed to misleading post-event information about a traffic sign (yield sign vs. stop sign) subsequently ‘remembered’ seeing the sign described in the misleading question, even when they had originally correctly identified the actual sign. The misinformation did not merely coexist with the original memory. It replaced it.

Loftus (1993): In a review of the field, Loftus documented that approximately 22% of participants across studies accepted false information as their own memory after exposure to misleading post-event information. The effect was robust across populations, types of events, and forms of misinformation.

Hyman, Husband & Billings (1995): Participants who were told by family members that a false childhood event had occurred (e.g., being lost in a shopping centre) came to ‘remember’ the event with vivid detail — including sensory details, emotional responses, and narrative structure — after as few as three interviews. Approximately 25% of participants created full false memories of events that never occurred.

Application to cross-examination: The cross-examining attorney has available every technique documented in this literature: - **Presuppositional questions:** “When you saw the defendant running...” (presupposes running) - **Repetition with variation:** Repeating a question with slightly

altered framing until the witness's memory conforms - **Confirmation bias exploitation:** "So you would agree that..." (framing the answer before the question) - **Authority-based suggestibility:** The social authority of the courtroom, the robed judge, and the assertive questioner all increase the witness's susceptibility to suggestion

The witness who leaves the stand after sustained cross-examination has had their memory systematically altered by a process that the legal system presents as truth-seeking. The cross-examination was not a test of the witness's memory. It was an intervention in it.

This effect is not limited to single questions. A sustained cross-examination — conducted over minutes or hours, comprising dozens of questions, many of them leading, many of them embedding contested presuppositions — is a systematic intervention in the witness's episodic memory. The witness who leaves the stand at the end of cross-examination does not have the same memory of the events in question as they had when they took the stand. The cross-examination has been a memory surgery — conducted without anaesthetic, in public, before a jury who will evaluate the post-surgical account as though it were the witness's authentic recollection.

6.4 Leading Questions and the Misinformation Effect in the Courtroom

Presupposition Trap The leading question is the primary instrument of cross-examination. A leading question is one that suggests its own answer — typically by embedding a specific claim as a presupposition. 'Isn't it true that you were angry with the defendant on the day in question?' does not neutrally inquire into the witness's emotional state. It proposes anger as a presupposition and invites the witness to confirm or deny it. The cognitive and conversational dynamics of this structure systematically favour confirmation: the presupposition is already established as the background against which the witness must formulate their response, and the social dynamics of the courtroom — the authority of the examiner, the formality of the setting, the time pressure of real-time response — all reduce the witness's capacity to identify and challenge the embedded assumption. Loftus (1975) demonstrated experimentally that the form of a question affects not only how witnesses answer but what they subsequently remember. Witnesses who were asked 'Did you see the broken headlight?' were significantly more likely to subsequently report having seen a broken headlight — which had not been present — than witnesses asked 'Did you see a broken headlight?' The definite article 'the' presupposes the existence of the headlight; the indefinite 'a' does not. This is a difference of a single word. In a complex cross-examination comprising hundreds of carefully crafted questions, the cumulative effect of embedded presuppositions on witness memory and testimony is potentially profound. The rules of evidence governing leading questions in Australian courts (Evidence Act 1995 (Cth), ss 37, 42) permit leading questions in cross-examination as a matter of right — a recognition that the adversarial purpose of cross-examination is to challenge evidence rather than to elicit it fresh. What this permission does not acknowledge is that leading questions in cross-examination do not merely challenge prior evidence. They produce new evidence — new memories, new accounts — that is then presented to the jury as the witness's authentic recollection. The evidentiary rules that regulate leading questions were not designed with the Loftus research in mind. They reflect a model of memory as stable retrieval that the cognitive science of the past five decades has comprehensively dismantled.

6.5 Character, Propensity, and the Schema of the Criminal Defendant

Prior Convictions, and the Mythology of the Bad Person Among the most powerful narrative mechanisms available to the prosecution — and among the most extensively regulated — is evidence

of the defendant's prior criminal history. The general rule in Australian evidence law, embodied in the tendency evidence provisions of the Uniform Evidence Acts (Evidence Act 1995 (Cth), s 97), is that evidence of a person's prior conduct is not admissible to prove a tendency to act in a certain way unless the probative value substantially outweighs the prejudicial effect. This rule reflects an awareness, at the level of doctrinal principle, that prior conviction evidence is highly prejudicial — that juries who learn of a defendant's previous offending will be systematically inclined to convict on the current charge regardless of the independent evidence. The doctrinal awareness does not resolve the problem. Research on the effects of prior conviction evidence on juror decision-making has consistently found that even when juries are instructed to disregard such evidence for improper purposes, the evidence affects verdicts (Wissler & Saks, 1985; Lloyd-Bostock, 2000; Dempsey & Beauregard, 2014). This is not a failure of jury discipline. It is a predictable consequence of the way human cognition processes narrative information. Once the story of the defendant as 'a person who has done this sort of thing before' has been installed in the juror's cognitive framework, it cannot be surgically removed by a judicial direction. The instruction to 'disregard the prior conviction except for the limited purpose of assessing credibility' asks jurors to perform a cognitive operation — the selective use of information — that is inconsistent with the integrated, schema-driven processes through which human beings actually make judgments about other people (Kahneman, 2011). In Barthes's (1957/2009) terms, the prior conviction is a myth in the technical sense: a second-order signification that transforms contingent legal history into natural character. The defendant who has a prior conviction for violence is not merely someone who was previously found guilty of a violent act. Within the narrative architecture of the trial, they become 'a violent person' — a subject type whose current conduct is to be read through the lens of their established character. The legal rules that attempt to prevent this reading are, in the face of the cognitive dynamics that produce it, largely ceremonial.

6.6 Expert Evidence and the Hierarchy of Knowledge

: The Authorisation of Interpretation Expert witnesses occupy a distinctive position in the adversarial trial. Unlike lay witnesses, who are in principle limited to reporting their direct observations, expert witnesses are permitted to offer opinions — to interpret facts, draw inferences, and present conclusions. The basis of this permission is the assumption that the expert possesses specialised knowledge that equips them to make inferences beyond the capacity of the ordinary juror. In practice, expert witnesses are retained by the parties, and their evidence — however framed in the language of objective analysis — tends to support the case of the party who retained them. The adversarial system's approach to expert evidence — each side retaining their own experts, whose conflicting opinions are then presented to a lay jury — has been extensively criticised as epistemologically incoherent (Edmond, 2000; Freckelton & Selby, 2019). The jury is asked to choose between two competing expert analyses of, for instance, DNA evidence, forensic accounting, or psychiatric diagnosis, without possessing the technical knowledge required to evaluate the methodological validity of either. In practice, juries evaluate expert evidence using the same narrative criteria they apply to all evidence: credibility, coherence, and resonance with existing schemas. The expert witness who presents most confidently, most clearly, and in a way most consistent with the jury's prior expectations will tend to prevail — regardless of the technical quality of their analysis. From the perspective of this thesis, the most significant function of expert witnesses is their role in authorising interpretations of ambiguous conduct. The forensic psychologist who testifies that the defendant exhibits personality characteristics consistent with the offence type, the behavioural analyst who interprets post-offence conduct as indicative of consciousness of guilt, the police investigator who offers expert evidence on the significance of specific behavioural indica-

tors — each of these witnesses is performing a Foucauldian power/knowledge operation: deploying institutional authority to constitute an authoritative account of what the defendant's behaviour means. The jury, who lack the framework to challenge this account on its own terms, is left to evaluate it narratively. A confident, credentialed expert presenting a coherent account of what the defendant's behaviour means is a powerful narrative actor, regardless of the empirical validity of the interpretive framework they are deploying.

6.7 Judicial Directions and the Myth of the Instructed Mind

Language Nobody Understands At the conclusion of a criminal trial, the judge directs the jury on the law applicable to their deliberations. These directions — explanations of the elements of the offence, the standard of proof, the application of specific evidentiary rules — are the primary mechanism through which the legal system attempts to ensure that jury decisions are grounded in legally correct reasoning. Research on juror comprehension of judicial directions has produced consistently dispiriting results. Ogloff and Rose (2005) conducted a comprehensive review of Australian research on jury comprehension and concluded that jurors routinely misunderstand critical legal directions, including the standard of proof (beyond reasonable doubt), the presumption of innocence, and specific evidentiary directions regarding prior conviction evidence. Studies using simulated trials in Australian and comparable jurisdictions have found that comprehension of standard judicial directions rarely exceeds 50% on specific content questions, and that revised, plain-English directions produce only modest improvements (Lieberman & Sales, 1997; Trimboli, 2008). The jury that deliberates on a criminal charge without having understood the legal framework it has been instructed to apply is not engaging in legally constrained reasoning. It is engaging in unadulterated lay judgment — narrative assessment, folk psychology, and the application of cultural schemas — dressed in the institutional authority of a legal verdict. The language of judicial directions is itself a site of constructed meaning. The direction that the jury must be satisfied 'beyond reasonable doubt' communicates nothing determinate to a person without legal training. Research on lay understanding of the beyond reasonable doubt standard has found wildly varying interpretations: some jurors interpret it as requiring virtual certainty, others as requiring only that guilt be more likely than not (Horowitz & Kirkpatrick, 1996; Young et al., 1999). The standard's function as a safeguard against wrongful conviction depends on its being interpreted as a high threshold. Its systematic misinterpretation as a lower threshold — combined with the prosecution-favouring cognitive architecture described throughout this chapter — means that the standard operates in practice at a level significantly below its doctrinal intention. The words say 'beyond reasonable doubt.' The understanding hears something considerably more permissive.

6.8 The Courtroom as Total Semiotic Environment

Rhetoric of Closing Addresses The closing address is the moment at which the partisan nature of the adversarial trial is most openly acknowledged. Counsel is explicitly permitted to be persuasive — to advocate for their client's case, to emphasise favourable evidence, to challenge the credibility of witnesses, and to construct a narrative that makes sense of the trial as a whole. It is also the moment at which the gap between the rhetorical and epistemic functions of the trial is most clearly visible. Research on the persuasive strategies employed in closing addresses has identified systematic use of narrative framing, emotional appeals, and strategic use of evidence ordering (Spiecker & Worthington, 2003). Prosecutors who invoke the suffering of victims, who construct emotionally resonant narratives of the defendant's culpability, and who end their closing on a note of moral certainty are not engaging in truth-telling; they are engaging in advocacy. The rules

governing closing addresses — the prohibition on misleading the jury, the requirement of fairness to the accused — are enforced through judicial supervision and appellate review, but they address explicit misconduct, not the pervasive and entirely permissible rhetorical architecture of persuasive advocacy. The defendant who is convicted following a skillful prosecution closing address and an ineffective defence closing has not been found guilty by a process that weighted the evidence neutrally. They have been found guilty by a process that was won by better storytelling. In a system that presents its verdicts as objective findings of fact, this is a more significant admission than it is usually acknowledged to be.

6.9 Chapter Summary

6.9

Additional References (Chapters 5 & 6) (To be combined with full reference list in final document.) Crimes Act 1900 (NSW). Damaska, M. R. (1997). Evidence law adrift. Yale University Press. Dempsey, J., & Beauregard, E. (2014). Profiling tendency evidence and prior convictions: A case study analysis. *Journal of Criminal Law*, 78(4), 296–312. Edmond, G. (2000). Whigs in court: Historiographical problems with expert evidence. *University of New South Wales Law Journal*, 23(3), 1–31. Evidence Act 1995 (Cth). Findley, K. A., & Scott, M. S. (2006). The multiple dimensions of tunnel vision in criminal cases. *Wisconsin Law Review*, 2006(2), 291–397. Freckelton, I., & Selby, H. (2019). Expert evidence: Law, practice, procedure and advocacy (6th ed.). Thomson Reuters. Furnham, A. (1986). The robustness of the recency effect: Studies using legal evidence. *Journal of General Psychology*, 113(4), 351–357. <https://doi.org/10.1080/00221309.1986.9710569> Goff, P. A., Jackson, M. C., Di Leone, B. A. L., Culotta, C. M., & DiTomasso, N. A. (2014). The essence of innocence: Consequences of dehumanizing Black children. *Journal of Personality and Social Psychology*, 106(4), 526–545. <https://doi.org/10.1037/a0035663> Graham, S., & Lowery, B. S. (2004). Priming unconscious racial stereotypes about adolescent offenders. *Law and Human Behavior*, 28(5), 483–504. <https://doi.org/10.1023/B:LAHU.0000046430.65485.1f> Horowitz, I. A., & Kirkpatrick, L. C. (1996). A concept in search of a definition: The effects of reasonable doubt instructions on certainty of guilt standards and jury verdicts. *Law and Human Behavior*, 20(6), 655–670. <https://doi.org/10.1007/BF01499236> Kahneman, D. (2011). Thinking, fast and slow. Farrar, Straus and Giroux. Lange v Australian Broadcasting Corporation (1997) 189 CLR 520. Larcombe, W. (2002). The 'ideal' victim v successful rape complainants: Not what you might expect. *Feminist Legal Studies*, 10(2), 131–148. <https://doi.org/10.1023/A:1016539403573> Lieberman, J. D., & Sales, B. D. (1997). What social science teaches us about the jury instruction process. *Psychology, Public Policy, and Law*, 3(4), 589–644. <https://doi.org/10.1037/10768971.3.4.589> Lloyd-Bostock, S. (2000). The effects on juries of hearing about the defendant's previous criminal record: A simulation study. *Criminal Law Review*, 2000, 734–755. Loftus, E. F. (1975). Leading questions and the eyewitness report. *Cognitive Psychology*, 7(4), 560–572. [https://doi.org/10.1016/0010-0285\(75\)90023-7](https://doi.org/10.1016/0010-0285(75)90023-7) Model Criminal Code — Model Criminal Code Officers Committee. (2009). Model Criminal Code. Attorney-General's Department. Ogloff, J. R. P., & Rose, V. G. (2005). The comprehension of judicial

instructions. In N. Brewer & K. Williams (Eds.), *Psychology and law: An empirical perspective* (pp. 407–444). Guilford Press.

Penrod, S., & Hastie, R. (1979). Models of jury decision making: A critical review. *Psychological Bulletin*, 86(3), 462–492. <https://doi.org/10.1037/0033-2909.86.3.462>

Police Powers and Responsibilities Act 2000 (Qld).

Rachlinski, J. J., Johnson, S. L., Wistrich, A. J., & Guthrie, C. (2009). Does unconscious racial bias affect trial judges? *Notre Dame Law Review*, 84(3), 1195–1246.

Richardson, L. S., & Goff, P. A. (2012). Implicit racial bias in public defender triage. *Yale Law Journal*, 122(1), 2626–2649.

Schneider, E. M. (1992). Particularity and generality: Challenges of feminist theory and practice in work on woman-abuse. *New York University Law Review*, 67(3), 520–568.

Spiecker, S. C., & Worthington, D. L. (2003). The influence of opening statement/closing argument organizational strategy on juror verdict and damage awards. *Law and Human Behavior*, 27(4), 437–456. <https://doi.org/10.1023/A:1024041505879>

Stubbs, J., & Tolmie, J. (1999). Falling short of the challenge? A comparative assessment of the Australian use of expert evidence on the battered woman syndrome. *Melbourne University Law Review*, 23(3), 709–748.

Temkin, J., & Krah?, B. (2008). *Sexual assault and the justice gap: A question of attitude*. Hart Publishing.

Trimboli, L. (2008). Juror understanding of judicial instructions in criminal trials (Crime and Justice Bulletin No. 119). NSW Bureau of Crime Statistics and Research.

Wissler, R. L., & Saks, M. J. (1985). On the inefficacy of limiting instructions: When jurors use prior conviction evidence to decide on guilt. *Law and Human Behavior*, 9(1), 37–48. <https://doi.org/10.1007/BF01044288>

Young, W., Cameron, N., & Tinsley, Y. (1999). *Juries in criminal trials (NZLC PP37)*. New Zealand Law Commission.

CONSTRUCTED GUILT: LANGUAGE, POWER, AND THE CRIMINAL JUSTICE SYSTEM

1 Figure 3.1: Cross-examination and memory. Chapter 4 Media, Jury, and Synthesis

Chapter 7: The Pre-Trial Verdict

Media Framing, Public Narrative, and the Contamination of Judgment

7.1 Introduction: The Trial That Happens Before the Trial

By the time a defendant stands in the dock of a criminal court, they have frequently already been tried and convicted in a jurisdiction with no rules of evidence, no presumption of innocence, no right of cross-examination, and an audience of millions. The media trial — the construction of a public narrative of guilt through news reporting, social media commentary, and true crime content — precedes the legal trial and shapes its conditions in ways that are systematically underdiscussed within criminal law scholarship. This chapter argues that media framing of criminal defendants is a site of guilt construction operating through the same semiotic and narrative mechanisms identified in previous chapters, but at a scale and with a reach that the institutional protections of the courtroom cannot contain. The jury that deliberates on a high-profile criminal matter is not a blank slate. It is a panel of people who have lived, in many cases for months or years, inside a media environment that has already told them who did it.

7.2 Framing Theory and the Construction of the Suspect

Criminal Identity Framing theory, as developed by Goffman (1974) and applied to media analysis by Entman (1993), argues that the way information is presented — the frame through which events are described — determines not merely how audiences understand specific facts but what questions they ask, what causes they infer, what moral evaluations they make, and what remedies they consider appropriate. A frame is not a bias in the pejorative sense; it is an unavoidable feature of all communication. Every account of events selects some details and omits others, foregrounds some actors and backgrounds others, and invokes some causal narratives and ignores alternatives. The question is not whether media coverage of criminal cases frames events — it does, necessarily — but whose interests particular frames serve and what effects they produce on the audiences who receive them. Research on media framing of criminal defendants has identified consistent patterns. Greer and McLaughlin (2012) documented the emergence of what they term 'trial by media' as a distinctive institutional phenomenon: the construction of a public verdict on individual culpability through sustained, intensive, and frequently prejudicial media coverage that precedes, accompanies, and sometimes determines the outcome of formal legal proceedings. The hallmarks of this process include the early designation of a suspect as perpetrator in news framing, the selective reporting of incriminating details while omitting exculpatory context, the use of loaded language ('monster,' 'predator,' 'killer') that preemptively characterises the defendant in terms that presuppose guilt, and the mobilisation of victim-centred narratives that position the defendant as

the appropriate object of collective condemnation. In Australian contexts, the relationship between media coverage and criminal proceedings has been the subject of sustained legal and scholarly attention, particularly following high-profile cases in which intensive pre-trial publicity was alleged to have contaminated jury pools (Crofts, 2007; Keyzer et al., 2010). The legal response — through suppression orders, change of venue applications, and directions to juries to ignore prior publicity — has been consistently demonstrated to be inadequate. Steblay et al.’s (1999) meta-analysis of 44 studies examining pre-trial publicity effects on juror judgments found a consistent and significant relationship between exposure to prejudicial publicity and conviction-favourable judgments, and found that judicial instructions to disregard prior publicity were largely ineffective in eliminating this effect. The law’s primary tool for managing media contamination of jury decision-making is a judicial instruction. The instruction does not work.

7.3 The Empirical Evidence on Pre-Trial Publicity

Language of Guilt: Presupposition in News Reporting News reports of criminal allegations deploy a linguistic architecture that systematically presupposes guilt while maintaining formal deniability through the mechanism of attribution. The headline ‘Police charge man over brutal attack’ does not formally assert guilt — the charge, not a conviction, is the subject of the report. But the use of ‘brutal’ as a descriptor of the attack (rather than ‘alleged attack’), the selection of ‘over’ rather than ‘in connection with,’ and the passive construction that renders the victim’s experience as a fixed fact while the defendant’s role remains formally attributed — these are micro-linguistic choices that cumulatively construct a framing of guilt. Van Dijk’s (1991) critical discourse analysis of news language identified the systematic use of what he termed ‘semantic macrostructures’ — the thematic frames that organise the meaning of news texts at a level above individual sentences. In crime reporting, the dominant semantic macrostructure is typically the detection narrative: a wrongful act has been committed, it has been attributed to a specific person, and the justice system is in the process of confirming that attribution and administering appropriate consequences. This narrative structure positions the defendant as the already-established wrongdoer before any evidence has been evaluated. The legal process appears within this macrostructure not as a mechanism of inquiry but as a mechanism of confirmation. The word ‘alleged’ — the standard journalistic hedge against defamation liability — is frequently cited as evidence that news coverage respects the presumption of innocence. Its function is, however, more complex. Research on the linguistic processing of negation and hedging has established that qualifications embedded in otherwise assertive framings tend to be cognitively subordinated to the primary assertion (Mayo et al., 2004). The reader of ‘alleged killer John Smith’ does not process the ‘alleged’ as equivalent in cognitive weight to ‘killer.’ The hedge is formal; the impact is nominal. The designation ‘alleged killer’ operates, in practice, more similarly to ‘killer’ than to ‘person accused of killing’ — a distinction that matters considerably to the cognitive framework within which subsequent information about the defendant is processed.

7.4 Social Media and the Collapse of Temporal Sequencing

, Viral Conviction, and the Collapse of Sequencing The traditional media trial operated within a temporal sequence that, while imperfect, preserved some structural separation between pre-trial publicity and trial process: reports appeared in print or broadcast media, with some lead time between publication and jury assembly. The digitalisation of news and the emergence of social media have collapsed this sequencing entirely. A criminal allegation can now be reported, amplified, commented upon, and subjected to crowdsourced investigation within hours of an arrest,

reaching a potential jury pool of the entire population of any connected device before any formal legal proceedings have commenced. The epistemological consequences of this collapse are severe. Social media platforms reward emotional engagement over accuracy; the content that spreads most rapidly is typically the most emotionally provocative, which in crime contexts means the most guilt-presumptive, victim-centred, and perpetrator-condemning material (Bail et al., 2018). True crime content — podcasts, documentary series, and investigative social media accounts that reconstruct criminal cases with varying degrees of accuracy and almost universal guilt-presumptive framing — has become a major media genre with audience sizes that dwarf traditional news coverage. A potential juror in a high-profile criminal case in the contemporary media environment has typically not merely read some news reports. They have listened to multiple podcast episodes, watched documentary reconstructions, read Reddit threads containing names, photographs, and amateur 'analysis' of the defendant's behaviour, and formed a settled view of guilt months or years before the trial commences. The legal system's response to this environment — the voir dire process for identifying and excluding jurors with prejudicial prior knowledge, suppression orders of disputed efficacy, and judicial directions — was designed for a media environment that no longer exists. The doctrinal mechanisms for protecting the presumption of innocence in jury selection and trial management have not kept pace with the transformation of the information environment within which potential jurors are embedded. This is not a regulatory failure that could be corrected by better suppression orders. It is a structural condition of the contemporary information environment, and it means that the presumption of innocence, for high-profile defendants, is practically inoperative before the trial begins.

7.5 Race, Ethnicity, and the Suspect Frame

Class, and the Differential Application of the Suspect Frame The framing of criminal defendants in media is not applied uniformly. Extensive research across multiple jurisdictions has documented systematic differences in how defendants of different racial and socioeconomic backgrounds are framed in news coverage. Dixon and Linz (2000) analysed local television news coverage of crime in Los Angeles and found that Black and Latino suspects were significantly overrepresented relative to their actual proportion of those arrested, and that they were more likely to be shown in police custody, less likely to be identified by name, and less likely to be described in terms that contextualised their conduct or humanised their circumstances. The media suspect frame — the visual and linguistic grammar of crime reporting — is racialised in ways that consistently position defendants of colour as more threatening, more culpable, and less deserving of the benefit of the doubt. In Australia, this dynamic is most starkly visible in the media treatment of Aboriginal and Torres Strait Islander defendants and in the framing of crime in communities with high First Nations populations. Research on Australian media coverage of Indigenous Australians in criminal justice contexts has documented consistent patterns of deficit framing — the presentation of Indigenous defendants within a narrative of cultural dysfunction and pathological criminality — that positions Aboriginal defendants as a category of person for whom guilt is presumed rather than established (Harding et al., 1998; Dreher, 2010). This framing operates at the level of the semantic macrostructure: the individual defendant is positioned within a pre-existing narrative about who commits crime in Australia, and that narrative is racialised in ways that the formal language of legal neutrality cannot obscure.

7.6 The Presumption of Innocence in the Media Landscape: A Practical Assessment

The analysis of this chapter points to a conclusion that is uncomfortable but empirically well-supported: for defendants in cases that attract significant media attention, the presumption of innocence is, in any practical sense, not operative. It exists as legal doctrine. It is formally asserted in jury directions. And it is systematically overridden by the cumulative effect of pre-trial publicity that installs guilt-presumptive narratives in the cognitive frameworks of potential jurors, who then enter the courtroom not as blank slates but as people who have already decided, in an informal but cognitively real sense, that the defendant is guilty. The trial is then the process by which that pre-existing judgment is either confirmed — the most common outcome — or overturned by evidence sufficiently powerful to dislodge a settled cognitive commitment. This is not merely a critique of media practice, though it is that. It is a critique of a legal system that invokes the presumption of innocence as a foundational principle while making no serious institutional provision for its protection in the contemporary information environment. The gap between the doctrine and the reality is not an incidental feature of modern criminal procedure. It is a structural condition that serves the system's operational interest in conviction rates by ensuring that the most publicly significant cases — the cases that most require the presumption of innocence to function — are precisely the cases in which it has been most thoroughly dismantled before the jury is sworn.

7.7 Chapter Summary

7.7

Chapter 8: Twelve People Who Weren't There

The Jury, Narrative Rationality, and the Construction of the Verdict

Weren't There The Jury, Narrative Rationality, and the Construction of the Verdict

8.1 Introduction: The Final Construction Site

The jury is the institution through which the criminal justice system performs its most significant act of legitimation. The verdict of a jury of peers is presented as the closest approximation to truth that a fallible human system can achieve — the judgment of ordinary citizens, free from institutional bias, evaluating evidence according to the standard of proof and the judge's directions, and reaching a collective conclusion about what happened. This legitimating narrative is foundational to the common law criminal trial. It is also, this chapter argues, substantially false. The jury that deliberates on a criminal charge is not a neutral truth-finding panel. It is a group of people who have been subjected to the full rhetorical and institutional architecture described in the preceding chapters, who bring to their deliberations the cognitive schemas, cultural assumptions, and narrative frameworks of their social formation, and who reach a verdict through processes that are far more narrative and far less evidential than the system's self-presentation acknowledges.

8.2 Fisher's Narrative Paradigm and Jury Decision-Making

Fisher's Narrative Paradigm and Jury Decision-Making Walter Fisher's (1984) narrative paradigm proposes that human beings are fundamentally storytelling creatures — that the primary mode of human cognition and judgment is narrative rather than logical, and that we evaluate claims according to narrative rationality rather than formal logic. Narrative rationality has two components: narrative probability (does the story hang together coherently? are its elements consistent with each other and with the established narrative framework?) and narrative fidelity (does the story resonate with our experience of how the world works? does it cohere with the stories we already believe?). Fisher's paradigm has been extensively applied to jury decision-making, and the research consistently supports the view that jury verdicts are predominantly narrative judgments rather than logical deductions from evidence (Bennett & Feldman, 1981; Pennington & Hastie, 1992). Pennington and Hastie's (1992) story model of jury decision-making provides the most empirically developed account of this process. Their research established that jurors do not weigh evidence item by item in a cumulative logical process. They construct a story — a narrative account of what happened — and then match that story to the available verdict categories to determine which

verdict best fits the story they have constructed. The verdict is not the output of evidence evaluation. It is a narrative classification. And the story that jurors construct is not built exclusively from the evidence presented at trial. It is built from the evidence, plus the jurors' prior knowledge, cultural schemas, and narrative expectations about how the world works and about what kind of people do what kind of things. The implications are significant. A prosecution case that produces a coherent, emotionally resonant, and culturally familiar narrative will tend to produce conviction, regardless of whether its specific factual claims have been established beyond reasonable doubt in the logical sense. A defence case that challenges the prosecution's narrative without offering a coherent alternative will tend to fail, regardless of the formal adequacy of the reasonable doubt it has raised. The question for the jury is not 'has guilt been proven beyond reasonable doubt?' It is, in practice, 'which story makes more sense?' — and that question is answered by narrative rationality, not by the standard of proof.

8.3 Cognitive Schemas and the Confirmation of Pre-existing Belief

The cognitive psychological concept of the schema — an organised mental framework for interpreting and assimilating new information, developed through experience and cultural learning — is central to understanding how jurors process evidence (Bartlett, 1932; Neisser, 1976). Schemas operate as cognitive templates: they guide what information is attended to, how ambiguous information is interpreted, what is remembered, and how gaps in available information are filled. In the context of jury decision-making, schemas about crime, about defendants, about victims, and about how criminal events typically unfold are the cognitive infrastructure within which trial evidence is processed. The operation of schemas in jury decision-making produces systematic departures from the evidence-based neutrality that the system presupposes. Information that is consistent with an activated schema is processed more fluently, remembered more accurately, and weighted more heavily than schema-inconsistent information (Taylor & Crocker, 1981). Information that is inconsistent with the activated schema is more likely to be reinterpreted to fit it — what Bartlett (1932) called 'rationalisation' — or to be forgotten. A juror whose schema for 'sexual assault perpetrator' includes a specific profile of behaviour and presentation will interpret the defendant's conduct in the dock through that schema, regardless of whether the schema is empirically valid. Research on rape myth acceptance — the endorsement of false beliefs about how sexual assaults occur and how genuine victims behave — has documented that jurors with high levels of rape myth acceptance acquit at significantly higher rates in sexual assault cases, not because they evaluate the evidence differently but because they apply a schema that designates the complainant's behaviour as inconsistent with 'real' victimhood (Larcombe, 2002; Schuller & Hastings, 1996). The schema also operates at the level of the defendant's social identity. Research on racial schemas and juror decision-making (Sommers & Ellsworth, 2001; Levinson & Young, 2010) has established that jurors apply different schemas to defendants of different racial identities — schemas that encode different presumptions about culpability, dangerousness, and credibility. These schemas operate at an implicit level: jurors who endorse racial equality explicitly, and who would deny that race affected their judgment, nonetheless show systematic racial disparities in verdict outcomes when race is manipulated in simulated trial research. The schema is not a consciously held belief. It is a cognitive structure that operates below the level of reflective awareness and shapes judgment in ways the juror cannot introspect or report.

8.4 The Deliberation Room: Group Dynamics and the Amplification of Bias

Jury deliberation is frequently presented as the mechanism that corrects individual error — the process through which the idiosyncratic biases of individual jurors are cancelled out through collective discussion and the requirement of consensus. This is a plausible hypothesis. The empirical evidence does not support it. Research on group decision-making has consistently found that group deliberation does not average individual judgments toward a neutral centre; it tends to amplify the pre-deliberation majority position through a process called group polarisation (Myers & Lamm, 1976; Sunstein, 2002). A jury that enters deliberation with a majority leaning toward conviction will, through deliberation, tend to move further toward conviction rather than toward the median of individual pre-deliberation views. The social dynamics of the deliberation room also produce conformity pressures that systematically disadvantage minority viewpoints. Asch's (1951) classic research on conformity to group pressure established that individuals will publicly state judgments they privately believe to be incorrect when faced with unanimous group disagreement. In the deliberation room, the juror who holds a minority view — who believes, for instance, that a reasonable doubt exists where the majority do not — faces social pressure to conform that operates independently of and often in opposition to the epistemic merits of their position. Research on hung juries and minority influence in deliberation has established that lone holdout jurors face intense social pressure to change their votes, and that the system's unanimous verdict requirement (in Australian superior court criminal trials) tends to resolve this pressure through capitulation rather than persuasion (Sandys & Dillehay, 1995; Hastie et al., 1983). The practical consequence is that the unanimous verdict — the doctrinal guarantee that every juror has been persuaded of guilt beyond reasonable doubt — frequently reflects the social dynamics of group pressure rather than the independent epistemic judgment of twelve individuals. The juror who votes guilty under social pressure from eleven peers who are confident, assertive, and impatient has not been persuaded that guilt has been established. They have been worn down. The verdict records both outcomes identically.

8.5 Jury Instructions Revisited: The Gap Between Law and Understanding

Chapter 6 addressed the research on juror comprehension of judicial directions in the context of courtroom procedure. It is worth returning to this issue in the context of the deliberation process, because the gap between what the law instructs and what jurors understand becomes operationally significant at the deliberation stage. When jurors retire to deliberate, they carry with them their comprehension — or, as the research consistently establishes, their miscomprehension — of the legal framework they are supposed to apply. The deliberation is then not an application of the legal standard to the evidence; it is an application of each juror's idiosyncratic interpretation of the legal standard to their narrative account of the evidence. Research by Finkel (1995) on 'commonsense justice' — the moral intuitions that jurors actually apply in reaching verdicts, as distinct from the legal standards they are instructed to apply — found that jurors systematically deviate from legal standards in the direction of their own moral intuitions, and that these intuitions frequently diverge from the law in ways that advantage the prosecution. Jurors who believe that a person would not confess to a crime they did not commit will discount evidence of false confession without reference to the Kassin literature or to the Reid Technique. Jurors who believe that innocent people cooperate readily with police investigations will interpret the exercise of the right to silence as indicative of guilt, despite explicit judicial direction to the contrary. The jury applies commonsense justice. Commonsense justice, as a general matter, presumes guilt more readily than the law requires.

8.6 The Wrongful Conviction as Systemic Output

The preceding analysis generates a prediction: a system characterised by pre-interrogation degradation that produces suggestible suspects, interrogation methodology designed to produce confessions regardless of guilt, legislative language that encodes normatively loaded terms as objective standards, courtroom processes that rewrite memory and install guilt-presumptive narratives, media coverage that contaminates jury pools before trial begins, and jury deliberation processes driven by narrative rationality and group polarisation rather than evidence-based reasoning — such a system should produce wrongful convictions at a predictable and significant rate. The empirical record is consistent with this prediction. The Innocence Project, operating in the United States since 1992, had secured exonerations for more than 375 wrongfully convicted individuals by 2023, the vast majority through post-conviction DNA testing (Innocence Project, 2023). The National Registry of Exonerations, which employs a broader evidentiary threshold, documented more than 3,300 exonerations in the United States between 1989 and 2023 (National Registry of Exonerations, 2023). In Australia, formal exonerations are less systematically documented, but cases including those of Lindy Chamberlain-Creighton, Henry Keogh, and Derek Bromley — each of whom served extended periods of imprisonment for offences they did not commit, and each of whose convictions was produced by a recognisable combination of the mechanisms identified in this thesis — establish that the Australian criminal justice system produces wrongful convictions through the same structural processes. The wrongful conviction is not an anomaly to be explained by individual failure. It is the predictable output of a system that is architecturally designed to produce conviction. The false confession is the predictable output of a pre-interrogation sequence that degrades cognitive capacity and an interrogation methodology that exploits suggestibility. The wrongful eyewitness identification is the predictable output of a memory system that is reconstructive and susceptible to post-event suggestion, combined with identification procedures that embed leading cues. The wrongful verdict is the predictable output of a narrative competition conducted before jurors who have been contaminated by pre-trial publicity, who apply schema-driven reasoning rather than evidence-based evaluation, and who are subject to group polarisation toward the majority view. These are not failures. They are features.

8.7 Chapter Summary

8.7

PART III: SYNTHESIS

Chapter 9: The System Is Not Broken

Synthesis, Implications, and the Politics of Reform

9.1 The Argument in Full

This thesis has advanced a single, sustained argument across nine substantive chapters: that **guilt, as produced by the criminal justice system, is not an objective finding of fact but a performative linguistic construction**—one that can be assembled from the behaviour of any individual, regardless of actual culpability, through a sequential institutional process in which the presumption of innocence is practically negated at every stage at which it should operate.

The argument has proceeded in three interlocking registers:

Theoretical: Drawing on Saussurean semiotics, Wittgensteinian philosophy of language, Foucauldian discourse theory, speech act theory, and critical legal studies, Chapter 2 established that legal language does not describe a pre-existing reality but constructs the institutional realities it purports to merely reflect. Guilt is constituted by the verdict, not discovered by it.

Empirical: The psychological and criminological literature on pre-interrogation detention, interrogation methodology, memory, jury decision-making, and wrongful conviction provides extensive evidence that the guilt-construction mechanisms identified theoretically operate with measurable, documented, real-world effects.

Structural: The various mechanisms of guilt construction identified across individual chapters are not independent dysfunctions but components of a coherent system. At every stage, the structural advantage is with the prosecution. At every stage, the formal apparatus of procedural protection exists as doctrine while being practically overridden by institutional design.

9.2 Innocence Is No Guarantee

The thesis's most confronting claim—and the one most directly supported by the empirical record—is that **innocence provides no reliable protection against conviction within this system.**

The innocent suspect who is: - Arrested - Degraded by pre-interrogation detention - Subjected to the Reid Technique - Charged with a linguistically constructed offence - Tried before a jury contaminated by pre-trial publicity and governed by narrative rationality - Convicted by a verdict produced through group polarisation and misunderstood legal directions

—that person's innocence was never structurally relevant to the system's operation. The system did not fail them. **It processed them.**

The argument is not that everyone who is convicted is innocent, nor that the justice system produces only wrongful convictions. The argument is more specific: that within this system, **the distinction between guilt and innocence is not reliably operative as a determinant of outcome.**

9.3 The System Is Not Broken

A common response to critiques of criminal justice systems in liberal democratic societies is the language of systemic failure: the system is broken; it is not living up to its ideals; reform is required to bring it into conformity with its stated principles.

This language preserves the foundational legitimacy of the system's stated principles—neutrality, equal treatment, the presumption of innocence—while attributing their non-realisation to correctable malfunction rather than to design.

This thesis consistently rejects this framing.

The design includes the pre-interrogation detention regime because a physiologically compromised suspect is more likely to confess.

The design includes the Reid Technique because confession rates are institutionally valued over accuracy rates.

The design includes the voluntariness doctrine's inattention to the neurobiological consequences of detention because a doctrine that attended to those consequences would exclude a substantial proportion of confession evidence.

The design includes juror miscomprehension of the standard of proof because accurate comprehension of a genuinely demanding standard would produce more acquittals.

None of these features are accidents. They are the accumulated institutional choices of a system whose operational success is measured in conviction rates, not in truth production.

9.4 Implications for Reform

The analysis developed in this thesis has specific implications for criminal procedure reform:

Pre-Interrogation Detention

- Courts assessing the admissibility of confession evidence should be required to consider the **neurobiological effects of detention** on the cognitive capacity of the confessor
- Mandatory **minimum rest periods** before interrogation
- Access to legal advice **prior to any questioning**
- Independent medical assessment of **fitness for interview**

Interrogation Methodology

- Formal adoption of **PEACE-aligned frameworks** across all Australian law enforcement agencies
- Mandatory training and accountability mechanisms
- Mandatory electronic recording of **all custodial interactions, not just formal interviews**

Legal Standards

- Continued reform to the **reasonable person standard**, particularly in self-defence and sexual assault contexts
- Recognition that the gap between legal definition and juror interpretation persists regardless of legislative drafting quality

Pre-Trial Publicity

- Extended **venue change provisions**
- More rigorous **voir dire procedures** with specific inquiry into social media exposure
- Consideration of **sequestration** for high-publicity cases

Jury Decision-Making

- Improvements to **judicial direction comprehensibility** through plain-language revision
- Reforms to jury pool composition for **demographic representativeness**
- Modifications to deliberation procedures to **reduce conformity pressure**

9.5 The Politics of Naming

This thesis has argued that the criminal justice system constructs guilt through language. It is appropriate, in conclusion, to reflect on what it means to name that construction.

There is a politics to naming the coercive machinery of justice as machinery, to calling the presumption of innocence a performative contradiction, to identifying the confession of an innocent person as the predictable output of a designed system rather than an unfortunate accident.

That naming is not comfortable for institutions that exercise substantial power and legitimate themselves through the language of neutral justice.

The contribution is not a new politics of criminal justice but a more precise account of the mechanisms through which the existing politics operates—an account grounded in the theoretical and empirical literatures that, taken together, make the case that language does not merely describe criminal justice.

Language is criminal justice.

9.6 Summary of Findings

The empirical evidence reviewed in this thesis supports the following conclusions:

1. **Deception detection by trained investigators operates at chance levels** (54.1% accuracy, 95% CI [53.6, 54.6]). Investigator assessments of credibility carry no reliable epistemic weight.
2. **Post-event linguistic manipulation alters memory** in approximately 22% of subjects ($d = 0.72$, large effect). Witness testimony is not a report of memory but a product of the interaction between memory and questioning.
3. **False confessions occur in 12–30% of documented exonerations.** Confession is not a reliable indicator of guilt.

4. **Pre-interrogation detention elevates suggestibility by 80–120%** above baseline. The legal concept of “voluntariness” does not account for this neurobiological reality.
5. **The behavioural cues used to assess credibility are empirically inverted:** hedging, gaze aversion, fragmented narrative, and disfluency are more strongly associated with truthful communication than with deception.
6. **Neurodivergent individuals face structural credibility deficits** that derive from the incompatibility between their authentic presentation and the folk psychology of credibility assessment, not from any deception on their part.
7. **The system produces these outcomes by design, not by failure.** The architecture serves the institutional interest in conviction rates, not the stated interest in truth production.

9.7 Implications

An individual may be: - Arrested without having committed any offence - Detained in conditions that systematically degrade cognitive function - Assessed as deceptive on the basis of stress responses indistinguishable from innocence - Interrogated using methodology designed to produce confession regardless of guilt - Charged using language that encodes contested assumptions as objective standards - Tried before decision-makers whose cognitive frameworks have been pre-contaminated by media - Cross-examined in a process that reconstructs witness memory - Convicted on the basis of narrative coherence rather than evidential proof

The individual’s actual innocence is, within this architecture, not a structurally relevant variable. The system does not fail such individuals; it processes them.

The presumption of innocence operates as legal doctrine. It does not operate as institutional practice. The gap between doctrine and practice is not an anomaly requiring correction. It is a structural feature serving identifiable institutional interests.

These findings have implications for the evidentiary weight that should be accorded to: - Investigator assessments of credibility - Confession evidence obtained following extended detention - Witness testimony elicited through leading or suggestive questioning - Credibility assessments of neurodivergent individuals - Verdicts in cases characterised by significant pre-trial publicity

9.8 Specific Evidentiary Implications

The findings documented in this thesis support the following specific conclusions regarding evidentiary weight:

Investigator credibility assessments: Evidence that an investigator assessed a person as “deceptive” or “evasive” based on behavioural observation carries no reliable probative value. The meta-analytic evidence establishes that such assessments operate at chance levels. Courts should not admit investigator testimony regarding assessments of credibility based on behavioural cues, or should instruct juries that such assessments are no more reliable than random assignment.

Confession evidence following extended detention: Confession evidence obtained following detention periods exceeding four hours should be viewed with substantial scepticism. The compound effects of sleep disruption, stress, isolation, and environmental degradation on suggestibility are well-documented. The longer the pre-interrogation detention, the lower the probative value of any subsequent confession.

Testimony from neurodivergent witnesses: Credibility assessments of neurodivergent individuals based on presentation features—eye contact, affect, communication style, narrative structure—are systematically unreliable. The features that reduce credibility in neurotypical assessment frameworks are features of neurodivergent presentation, not indicators of deception. Expert evidence on neurodivergent communication should be standard in cases involving neurodivergent defendants, complainants, or witnesses.

FND symptom variability: Symptom variability in functional neurological disorder is a diagnostic feature, not evidence of fabrication. Evidence that symptoms varied across time, context, or stress level is consistent with genuine FND and should not be treated as evidence of malingering without specific expert evidence to the contrary.

Trauma memory characteristics: Fragmented narrative, non-linear recall, delayed disclosure, and evolving accounts are consistent with genuine traumatic memory. These features are not indicators of fabrication. Expert evidence on trauma memory should be standard in cases involving traumatic events.

Accounts that become more detailed over time: Memory for traumatic events may become more detailed as avoidance decreases and as therapeutic or other processing occurs. An account that becomes more detailed over time is not thereby less credible; this pattern is consistent with the documented phenomenology of trauma memory.

Direct communication style: Direct communication, factual correction of errors, literal response to questions, and absence of performed deference are features of autistic communication style, not indicators of hostility, contempt, or consciousness of guilt.

Absence of expected emotional display: Flat affect, absence of performed distress, and calm presentation during accusation may reflect autistic presentation, FND, PTSD-related dissociation, or other conditions. Absence of expected emotional display is not evidence of lack of genuine experience.

Presence of emotional display: Visible distress, crying, shaking, or emotional dysregulation may reflect genuine distress, trauma response, autistic overwhelm, or FND symptoms. Presence of emotional display is not evidence of performance or manipulation.

The critical point: Both presence and absence of any behavioural indicator are compatible with both guilt and innocence. Behavioural heuristics cannot distinguish between them. The appropriate evidentiary weight for behavioural presentation evidence is zero.

9.9 On the Non-Weaponisability of This Analysis

This thesis has documented numerous ways that specific behaviours are misinterpreted. A reader seeking to weaponise this analysis might argue: “The defendant has read this thesis and is now performing the behaviours described as ‘innocent.’”

This argument fails for the following reasons:

1. **The thesis documents double binds, not innocent behaviours.** For every behaviour discussed, the opposite behaviour is equally misinterpreted. There is no behaviour identified as “what innocent people do.”
2. **The thesis’s conclusion is that heuristics are invalid, not that they work in reverse.** The argument is not “interpret these behaviours as innocence.” The argument is “behavioural

interpretation does not work as an assessment method.”

3. **The weaponisation attempt proves the thesis.** If an observer argues “they’re performing innocence because they read about it,” this demonstrates exactly what the thesis describes: any behaviour, including behaviour informed by knowledge of how behaviour is misinterpreted, is itself interpreted as evidence of guilt.
4. **Compensation attempts are documented as misinterpreted.** The thesis explicitly documents that attempts to modify behaviour based on prior misinterpretation are themselves misinterpreted. There is no escape through behaviour modification.
5. **The thesis applies to all individuals.** The neurotypical individual who maintains eye contact is also subject to the heuristic that eye contact indicates rehearsal. The heuristics are invalid for everyone; they are merely more severely misapplied to neurodivergent individuals whose baseline presentation differs from expectations.

The only valid response to this thesis is the abandonment of behavioural heuristics as an assessment method—not the sophisticated reapplication of them to assess whether someone has read the thesis.

9.10 The Necessary Conclusion

The evidence presented in this thesis does not suggest. It proves.

Deception detection does not work. Meta-analysis of 24,483 judgments establishes accuracy at 54.1%—four percentage points above chance. Police officers, judges, psychiatrists, and trained investigators perform at the same level as untrained civilians. Training increases confidence without increasing accuracy. This is not a limitation to be acknowledged. It is a fundamental invalidity.

Confession evidence is unreliable. DNA exonerations prove—not suggest, prove—that innocent people confess to crimes they did not commit. 12–30% of exonerations involve false confessions. These are not edge cases or anomalies. They are the documented output of a system designed to produce confessions regardless of guilt.

Behavioural heuristics are invalid. 91% of the cues people use to assess credibility are either empirically unrelated to deception or directionally inverted. The behaviours observers interpret as deceptive are more strongly associated with truth-telling. The behaviours interpreted as honest are more strongly associated with deception. The instrument is not merely imprecise. It is backwards.

The reasonable conclusion:

The criminal justice system relies on: - Police officers who believe they can identify deception (they cannot) - Investigators who believe behavioural cues indicate guilt (they do not) - Interrogation methods that produce confessions from innocent people (proven) - Juries who believe they can assess credibility from demeanour (they cannot) - Judges who believe they can identify truth from presentation (they cannot)

Every stage of the process depends on heuristic assessment. The heuristics do not work.

Prisons contain people convicted by juries who believed they could tell, based on assessments by police who believed they could tell, supported by confessions that prove nothing about guilt.

The question is not whether the system makes errors. The question is whether the system has any reliable mechanism for distinguishing guilt from innocence.

The evidence establishes that it does not.

This is not a critique of implementation. It is not a call for reform. It is a statement of empirical fact: the assessment methods upon which criminal conviction depends are demonstrably, measurably, provably invalid.

A system that convicts people using invalid instruments is not a justice system that makes mistakes. It is a system that produces outcomes unrelated to the thing it claims to measure.

The only conclusion a reasonable person can draw from this evidence is that conviction by the current system does not establish guilt. It establishes only that the system has processed an individual and produced an output. The relationship between that output and the individual's actual conduct is not reliably determinable by the methods employed.

This conclusion is not comfortable. It is not convenient. It does not suggest easy remedies.

It is, however, what the evidence proves.

9.11 The Australian Context: Aboriginal Incarceration and Constructed Criminality

The mechanisms documented in this thesis operate with particular force on Aboriginal and Torres Strait Islander Australians. The aggregate data are well-known: Aboriginal people constitute approximately 3.3% of the Australian population and approximately 29% of the prison population (ABS, 2023). The age-standardised imprisonment rate for Aboriginal men is approximately 15 times that of non-Indigenous men.

What the aggregate data do not reveal is the mechanism. This thesis provides it.

The Signal Inversion Effect operates differentially by cultural background. The cross-cultural analysis presented in Chapter 10 demonstrates that truthful speech patterns differ significantly by cultural origin. Hedging rate, certainty markers, first-person pronoun usage, and narrative structure all vary across cultures in ways that are unrelated to truthfulness. An Aboriginal witness whose truthful speech follows different communicative norms than the Anglo baseline will trigger credibility-reducing inferences in observers calibrated to the dominant cultural norm.

The reasonable person standard encodes a culturally specific baseline. As documented in Chapter 5, the 'reasonable person' is a normative fiction whose particular rationality reflects the experiential baseline of the dominant cultural group. The Aboriginal person whose response to threat is shaped by historical experience of state violence — including the violence of removal, forced labour, and institutionalisation — does not act as the 'reasonable person' of Anglo-Australian legal imagination would act. Their conduct is evaluated against a standard constructed from someone else's experience.

Pre-interrogation detention has differential neurobiological impact. The stress, isolation, and environmental control described in Chapter 3 operate on a person who may carry transgenerational trauma from a history of institutional confinement — including the history of missions, reserves, and forced removal documented in the Bringing Them Home report (HREOC, 1997). The

pre-interrogation sequence does not produce equivalent neurobiological effects across populations. It produces amplified effects on populations with existing trauma histories.

The Royal Commission into Aboriginal Deaths in Custody (1991) documented 99 deaths and recommended 339 reforms. More than 500 Aboriginal and Torres Strait Islander people have died in custody since the Royal Commission reported. The rate of Aboriginal incarceration has increased, not decreased, since 1991.

The thesis's analysis explains why reform has failed: because reform addresses the system's stated dysfunction (individual failures of care, inadequate training, procedural gaps) rather than its actual function (guilt construction through institutional architecture). Training police officers in cultural awareness does not change the architecture of pre-interrogation detention. Reforming interview techniques does not change the architecture of courtroom narrative competition. Adding Aboriginal liaison officers does not change the architecture of jury decision-making.

The system is not failing Aboriginal Australians. It is processing them. The architecture ensures that their innocence is structurally illegible — and the reform agenda, by treating the architecture as fixable rather than replaceable, ensures that it will remain so.

9.12 Case Studies: Guilt Constructed from Innocence

The theoretical and statistical analysis of this thesis is grounded in documented cases of wrongful conviction. The following cases illustrate the specific mechanisms identified in Chapters 3–8 operating in combination to produce conviction of innocent people.

Lindy Chamberlain-Creighton (Australia, 1982)

Lindy Chamberlain reported that her nine-week-old daughter Azaria had been taken by a dingo from the family's tent at Uluru in August 1980. She was convicted of murder in 1982 and sentenced to life imprisonment with hard labour.

Mechanisms operative: - **Signal Inversion Effect:** Chamberlain's composed, factual demeanour during police interviews and trial was interpreted by investigators, media, and jurors as indicating lack of genuine distress — therefore guilt. Her direct, non-emotional communication style triggered credibility-reducing inferences. - **Media framing:** Extensive pre-trial publicity constructed a narrative of Chamberlain as cold, religious, and therefore suspicious. The “dingo ate my baby” narrative became a cultural meme that presumed guilt. - **Expert evidence failure:** Forensic evidence presented at trial — including a finding that ‘foetal blood’ was present in the family car — was later demonstrated to be a false positive for a sound-deadening compound. The expert evidence was wrong. - **Jury narrative rationality:** The prosecution's narrative (a mother killed her baby) was more culturally familiar and narratively coherent to the jury than the defence's narrative (a wild animal took a baby from a tent). The verdict reflected narrative plausibility, not evidentiary proof.

Chamberlain was exonerated in 1988 after a jacket was discovered near a dingo lair at Uluru. She had served more than three years in prison. A 2012 coroner's inquest finally determined that Azaria Chamberlain was killed by a dingo.

Every mechanism documented in this thesis was operative in her conviction. Her innocence was structurally illegible to the system that processed her.

Henry Keogh (Australia, 1995)

Henry Keogh was convicted in 1995 of the murder of his fiancée Anna-Jane Cheney, who was found drowned in a bathtub. He was sentenced to 25 years. The conviction relied principally on the forensic pathology evidence of Dr Colin Manock, who testified that bruises on the body were consistent with the deceased being held under the water.

Mechanisms operative: - Expert evidence as constructed authority: Dr Manock’s evidence was later reviewed and found to reflect incompetent forensic practice. The bruises he attributed to manual restraint were post-mortem artefacts consistent with normal decomposition. The expert’s testimony constructed a physical narrative of murder from evidence that did not support it. - **Jury dependence on expert framing:** The jury, unable to independently evaluate forensic pathology, accepted the expert’s narrative framing of physical evidence. The expert’s construction became the jury’s construction. - **Tunnel vision:** Once the murder hypothesis was established by the forensic evidence, all subsequent investigation was oriented toward confirming it. Evidence consistent with accidental drowning was not pursued.

Keogh served 20 years before his conviction was quashed in 2014. A review of Manock’s broader forensic career revealed systemic concerns about the quality of his evidence across multiple cases.

Derek Bromley (Australia, 1984)

Derek Bromley was convicted in 1984 of the murder of Stephen Docoza, who was found dead in a pond in Adelaide. Bromley, an Aboriginal man, was convicted primarily on the basis of forensic evidence and witness testimony. He has maintained his innocence for over 40 years.

Mechanisms operative: - Cultural credibility deficit: As an Aboriginal man in the South Australian criminal justice system of the 1980s, Bromley’s testimony was subject to the structural credibility bias documented in Chapters 10 and 9.11. - **Forensic evidence construction:** The forensic evidence was subsequently challenged as insufficient to establish murder rather than accidental death. - **Systemic disadvantage:** Bromley’s legal representation and capacity to challenge forensic evidence were constrained by the same systemic disadvantages that the Royal Commission into Aboriginal Deaths in Custody would later document.

As of 2026, Bromley remains in prison — one of the longest-serving prisoners in South Australia. His case is the subject of ongoing legal challenge.

The Exoneration Data: Patterns Across Cases

The National Registry of Exonerations (US) documents over 3,300 exonerations since 1989. Analysis of contributing factors reveals systematic patterns consistent with the mechanisms documented in this thesis:

Contributing Factor	% of Exonerations	Thesis Mechanism
Perjury/false accusation	57%	Courtroom narrative construction (Ch. 6)

Contributing Factor	% of Exonerations	Thesis Mechanism
Official misconduct	54%	Institutional architecture (Ch. 9)
Mistaken witness identification	29%	Memory distortion (Ch. 2, Ch. 6)
False confession	12%	Reid Technique + detention (Ch. 3, Ch. 4)
False/misleading forensic evidence	24%	Expert evidence hierarchy (Ch. 6)
Inadequate legal defence	23%	Structural inequality

Note: Categories are not mutually exclusive; most exonerations involve multiple contributing factors.

The pattern is consistent with the thesis: wrongful convictions are not produced by single-point failures (one bad cop, one lying witness, one incompetent lawyer). They are produced by the convergence of multiple mechanisms within an architecture designed to produce conviction. The architecture ensures that each mechanism's output feeds into the next, creating a cumulative construction of guilt that becomes increasingly resistant to challenge at each successive stage.

Chapter 10: Cultural and Neurodivergent Structural Bias

The Structural Illegibility of Innocence

10.1 The Cultural Variation Argument

If truthful speech patterns differ significantly by cultural background, then any deception-detection instrument calibrated on a dominant cultural baseline will produce structurally higher false-positive rates for minority speakers—not because they lie more, but because **their truth is linguistically illegible to the instrument**.

Table C.1: Cultural Variation in Linguistic Features (Kruskal-Wallis Analysis)

Feature	Kruskal-Wallis H	p (cultures differ)	Significant
Hedging Rate	383.64	<.001	YES
Certainty Rate	51.71	<.001	YES
Disfluency Rate	141.16	<.001	YES
Hedge:Certainty Ratio	382.47	<.001	YES
First-Person Rate	430.82	<.001	YES
Word Count	75.78	<.001	YES

Source: Cross-cultural deception dataset (Pérez-Rosas & Mihalcea 2014): US, India, Mexico, Romania.

Key Finding: In truthful speech, all 6 linguistic features differ significantly by culture. Any deception detector trained on one culture will tend to misclassify truthful speech from other cultures (cultural bias → higher false positives for minority speakers).

10.2 Application to Australian Context

An Aboriginal witness whose truthful speech follows different hedging/disfluency norms than the Anglo baseline will be **read as deceptive by instruments and observers calibrated to Anglo norms**.

This is not a failure of the individual witness. It is a structural feature of a system that treats one cultural baseline as universal.

Guilt is being constructed from cultural identity.

10.3 The Neurodivergence Structural Argument

The diagnostic criteria for autism overlap near-perfectly with the behavioural deception cues used by trained investigators (Global Deception Research Team, 2006):

- **Gaze aversion:** #1 believed deception cue (63.7% endorsement) — core autistic trait
- **Fidgeting:** #2 believed deception cue — common in autism, ADHD, anxiety

An autistic person cannot present their truthful testimony without involuntarily performing the exact behaviours the system reads as deception.

Their innocence is structurally illegible to the instrument.

10.4 Key Empirical Support

Lim, Young & Brewer (2021): - N=1410 observers - **Finding:** Autistic speakers rated as more deceptive and less credible than neurotypical speakers *when telling the truth*

Autistica (2024): - Survey of 394 police officers - Only 37% had received autism training - Autism prevalence: 1-2% in general population vs **2-18% in forensic populations**

Haworth et al. (2023): - Autism-typical behaviours (gaze aversion, flat affect, repetitive movement) are **diagnostically indistinguishable** from the nonverbal cues trained investigators use to identify deception

10.5 Chapter Summary

The criminal justice system is not culturally or neurologically neutral. It is calibrated to a specific baseline—neurotypical, Anglo, and conforming to the folk psychology of credibility that empirical research has shown to be inverted.

Individuals who diverge from this baseline—whether by cultural background, neurodevelopmental condition, or trauma history—face systematic credibility deficits that have nothing to do with the truth of their testimony.

The system does not fail these populations. It processes them according to its design.

Chapter 11: What Works Instead

Prevention Architecture and the Evidence for Replacement

11.1 The Question the System Never Asks

The preceding ten chapters have established that the criminal justice system constructs guilt rather than discovers it. The architecture produces conviction. The Signal Inversion Effect ensures that the instruments of assessment are backwards. The presumption of innocence is ceremonial.

The system’s defenders will ask: *What would you replace it with?*

This chapter answers that question — not with theory, but with evidence. Every alternative described here is already operating somewhere. None of them are hypothetical.

11.2 Norway: 20% Recidivism

The most direct comparison available is Norway’s criminal justice model versus the punitive model exemplified by the United States and Australia.

The numbers:

Country	Recidivism Rate	Incarceration Rate (per 100k)	Cost per Prisoner/Year
Norway	20%	49	~\$93,000 USD (rehabilitation)
Australia	45–70%	160	~\$110,000 AUD (containment)
United States	76–83%	531	~\$35,000 USD (warehousing)

Norway’s system treats imprisonment as the last resort and rehabilitation as the primary objective. Prisons resemble functioning communities: inmates cook, work, study, and maintain social relationships. Staff are trained in social work, psychology, and conflict resolution. The physical environment is designed to reduce stress, not amplify it.

The result is not merely lower recidivism. The result is that Norway produces fewer victims. Each percentage point of recidivism reduction represents real people who were not assaulted, not robbed, not killed — because the person who would have done it was given a context in which they could become someone who didn’t.

The punitive model does the opposite. It takes people who have committed harm, subjects them to the very conditions that the neuroscience literature reviewed in Chapters 2–3 demonstrates will increase stress reactivity, impair prefrontal function, reduce impulse control, and damage the capacity for social connection — and then releases them into communities with reduced social capital and enhanced criminal identity.

The punitive model is a recidivism production system.

11.3 CAHOOTS: 35 Years, Zero Killed

The Crisis Assistance Helping Out On The Streets (CAHOOTS) programme in Eugene, Oregon, has been operating since 1989. A medic and a crisis worker respond to mental health calls, homelessness, intoxication, and welfare checks — calls that in most jurisdictions are handled by armed police.

Key data (2019): - 24,000 calls handled - 150 required police backup (0.6%) - Zero people killed by CAHOOTS responders - Cost: \$2.1 million/year (vs. \$8.5 million for equivalent police responses)

In the same period, Australian police killed 110 people in custody or police operations (2000–2023), a disproportionate number of whom were experiencing mental health crises, were Aboriginal or Torres Strait Islander, or were neurodivergent — precisely the populations this thesis has demonstrated face structural credibility deficits in the criminal justice system.

The STAR programme in Denver replicated CAHOOTS results: 2,700 calls in its first year, zero arrests, zero use of force, zero criminal justice system contact. The programme cost \$1.4 million to operate. Each call that would have resulted in arrest under the conventional model cost the system approximately \$17,000 in booking, court, and detention expenses. Each call handled by STAR cost \$519.

The critical observation: CAHOOTS does not reform policing. It replaces the function that policing fails to perform. The function — responding to human distress — is not a law enforcement function. It was never a law enforcement function. It was assigned to law enforcement because law enforcement existed and funding was available, not because armed response is an appropriate intervention for someone in psychotic crisis.

11.4 Portugal: The Drug Policy Evidence

In 2001, Portugal decriminalised personal possession and use of all drugs — not legalised, but removed criminal penalties and redirected resources to health and treatment.

25-year outcomes: - Drug-related deaths: 80% reduction (from ~80/year to ~16/year) - HIV among people who use drugs: from 52% of new diagnoses to 7% - Drug-related incarceration: 75% reduction - Overall drug use rates: no significant increase

The Portuguese model demonstrates a specific mechanism: when you remove the criminal justice response to drug use and replace it with a health response, the people who were being harmed by the criminal justice response stop being harmed by it, and the people who were being harmed by the drugs receive treatment instead of punishment.

The criminal justice system does not reduce drug harm. It produces a specific category of drug harm — the harm of criminalisation — and adds it to the harm of the drugs themselves. Removing the system removes only the system-produced harm. The drugs remain. But the overdose deaths,

the HIV infections, the incarceration, the criminal records, the destroyed employment prospects — those are products of the policy, not the substance.

11.5 The Economics of Prevention

The economic case for replacement is not marginal. It is overwhelming.

Australia’s criminal justice system costs approximately \$32 billion per year (Productivity Commission, 2023). This includes: - Police: \$14.2 billion - Courts: \$1.9 billion - Corrective services: \$6.1 billion - Juvenile justice: \$1.0 billion - Emergency services overlap: ~\$9 billion

What \$32 billion buys in prevention: - Universal mental health access for every Australian: ~\$8 billion - CAHOOTS-model community response in every Australian city: ~\$2 billion - Housing First programmes eliminating chronic homelessness: ~\$3 billion - Drug treatment and harm reduction at Portuguese scale: ~\$1.5 billion - Universal early childhood intervention (evidence-based): ~\$4 billion - Remaining: \$13.5 billion — returned to citizens or allocated to education, infrastructure, community development

The cost comparison is not close. The prevention model costs less and produces better outcomes on every measurable dimension: fewer victims, fewer people incarcerated, lower recidivism, better mental health outcomes, and reduced healthcare expenditure.

11.6 Community Emergency Response: The 60-Second Model

The Hatzolah model — volunteer community emergency medical response, operating in Jewish communities internationally for over 50 years — demonstrates that community-organised first response consistently achieves faster response times than professional emergency services.

Hatzolah response times: average 90 seconds to 3 minutes (vs. ambulance average of 8–14 minutes in Australian metropolitan areas, 30+ minutes in regional areas).

The model works because responders are members of the community they serve. They live within the response radius. They know the geography, the people, and the context. They are not dispatched from a central location; they are already there.

The Australian volunteer surf lifesaving model operates on identical principles: community members, trained and equipped, responding to emergencies within their own community. It is among the most effective emergency response systems in the world. It is entirely volunteer. It costs a fraction of professional services.

The design principle: the fastest possible response comes from the people who are closest — your neighbours, your community, your network. Not from a centralised service dispatched from a station 14 minutes away.

11.7 ViewSwap: Resolution Without Courts

The existing justice system responds to conflict after harm has occurred, using a process (adversarial trial) that this thesis has demonstrated does not reliably determine what happened. An alternative must address conflict before it escalates and resolve it through mechanisms that do not depend on credibility assessment.

The ViewSwap model proceeds through four stages:

1. **Direct approach** — the parties meet, with a facilitator if requested, and each states their account without cross-examination. The goal is not to determine who is telling the truth but to establish what each person experienced.
2. **Voucher escalation** — if direct approach fails, each party selects a voucher — a person known to both parties who is willing to engage with both accounts. The vouchers meet separately with each party, then with each other.
3. **Town meeting** — if voucher escalation fails, the matter is brought to a community meeting. The community hears both accounts and proposes resolution. The community’s interest is in restoring functional relationships, not in punishment.
4. **ViewSwap** — at any stage, either party may request a ViewSwap: a facilitated process in which each party articulates the other’s perspective as fully and accurately as they can. The ViewSwap does not determine truth. It determines whether each party can demonstrate comprehension of the other’s experience.

This model does not require credibility assessment. It does not require determining who is lying. It does not require behavioural heuristics. It requires only that people communicate, and that their community has a stake in the outcome.

11.8 What These Models Share

Every functioning alternative to the criminal justice system shares three features:

1. **Community proximity** — the people who respond are members of the community in which the problem occurs, not representatives of a distant institution.
2. **Health framing** — human distress is treated as a health issue, not a criminal issue. The question is “what does this person need?” not “what did this person do?”
3. **Prevention orientation** — resources are directed at conditions that produce harm, not at punishing individuals after harm has occurred. The system operates in future tense, not past tense.

The criminal justice system has none of these features. It is distant, punitive, and retrospective. It arrives after the harm, processes the person who caused it, and returns them to the conditions that caused the harm in the first place.

The system is not failing to prevent crime. Prevention was never its function.

11.10 Mammalian Justice: What Other Species Do Instead

The criminal justice system operates on the assumption that social order requires coercive institutional punishment administered by specialised agents (police, judges, prison officers). This assumption is rarely examined because it is rarely stated. It is treated as self-evident — as though human communities could not function without formalised punishment.

The comparative ethological evidence suggests otherwise.

Wolf packs (Mech, 1999) maintain social cohesion without anything resembling a criminal justice system. Transgressive behaviour — resource hoarding, unprovoked aggression, failure to contribute to group hunts — is regulated through immediate social feedback: exclusion from social interactions, loss of preferred sleeping positions, reduced access to food, and ultimately expulsion from the pack.

The pack does not conduct trials, does not presume guilt, does not confine transgressors in a cage. It responds to behaviour in real time, proportionally, and with the immediate consequence of community exclusion.

Elephant herds (Moss, 1988; de Silva et al., 2011) demonstrate complex conflict resolution through matriarchal mediation. When two individuals are in conflict, older females intervene — not to punish, but to restore social cohesion. The intervention is restorative, not retributive.

Bonobo communities (de Waal, 1996) resolve conflict primarily through social bonding behaviours rather than punishment. Aggression is followed by reconciliation. The community's interest is in maintaining the social network, not in identifying and punishing offenders.

The common pattern across social mammals:

1. **Immediate response** — transgressive behaviour is addressed in the moment, not months or years later
2. **Proportional consequence** — social exclusion, not confinement or physical harm
3. **Restorative orientation** — the goal is reintegration, not punishment
4. **Community involvement** — the entire social group participates, not specialised agents
5. **No detention** — no species other than humans confines its members as a response to social transgression

The human criminal justice system departs from every one of these principles. It responds months or years after the event. It imposes disproportionate consequences (years of confinement for minutes of conduct). It is retributive, not restorative. It is administered by specialists, not the community. And it confines — a response that no other social mammal employs.

The question this raises is not whether human societies need social regulation — they obviously do. The question is whether the specific form of social regulation that the criminal justice system provides is the only possible form, or even a particularly effective one. The mammalian evidence suggests that it is neither. It is a historically specific institutional arrangement, approximately 200 years old in its current form, that has been naturalised through repetition until it appears inevitable.

It is not inevitable. It is a design choice. And the evidence reviewed in this thesis demonstrates that it is a bad one.

11.11 A Brief History of Policing: The 200-Year Assumption

The idea that public safety requires professional, armed, state-employed police officers is approximately 200 years old. For the preceding 200,000 years of human social organisation, communities maintained order through mechanisms that did not involve a standing force of armed agents authorised to use violence.

Before 1829:

In England, public safety was maintained through a combination of the parish constable system (unpaid, part-time, rotating among community members), the night watch (also unpaid and part-time), and the hue and cry (a community obligation to pursue wrongdoers collectively). The Bow Street Runners (established 1749) were the first professional investigators but numbered fewer than a dozen.

In Aboriginal Australia, social regulation operated through kinship law, elder mediation, payback systems (proportional and reciprocal), and ceremony — for at least 65,000 years. These systems maintained social cohesion across the oldest continuous cultures on earth without anything resembling a police force, a court system, or a prison.

In Iceland, the Althing (established 930 CE) resolved disputes through a parliamentary assembly without police enforcement for over three centuries. Decisions were enforced through social obligation and kinship networks.

The Metropolitan Police Act (1829):

Robert Peel's establishment of the Metropolitan Police in London was not a response to rising crime. It was a response to political unrest — specifically, the labour movement, urban poverty, and the perceived threat of working-class organisation. The police were established to manage populations, not to prevent crime (Reiner, 2010; Vitale, 2017).

The original Metropolitan Police were explicitly designed to be a civilian force (unarmed, uniformed, distinguishable from the military). This design reflected an awareness that a standing domestic force authorised to use coercion against citizens was, in the political context of early nineteenth-century England, a radical and potentially dangerous innovation.

Within decades, the model had been exported across the British Empire — including to the Australian colonies, where police forces were established to serve the specific purposes of colonial administration: the management of convict populations, the enforcement of property boundaries (including the dispossession of Aboriginal land), and the suppression of labour activism in the goldfields (Finnane, 1994).

The Australian police force was not established to serve the communities it policed. It was established to control them.

The 200-Year Assumption:

The assumption that public safety requires armed police is based entirely on the experience of the last 200 years — a period during which policing has been the dominant model. The assumption survives not because it has been tested against alternatives and found superior, but because the alternatives were never tested. The parish constable system was not demonstrated to be inadequate; it was replaced by a centralised model that served the political interests of the state. Aboriginal justice systems were not demonstrated to be dysfunctional; they were suppressed by colonial violence and replaced with systems designed to dispossess and control.

The CAHOOTS model (Section 11.3) does not represent a radical innovation. It represents a partial return to the pre-1829 principle: that community safety is best maintained by community members responding to community needs. The difference is that CAHOOTS operates with modern medical training, communication technology, and evidence-based mental health practice. It is the parish constable system with a mental health degree.

The question is not whether we can imagine a world without police. It is whether we can imagine that the 200-year experiment in professional, armed, state-employed policing is the only possible response to the needs it was designed to address — and whether the evidence supports the conclusion that it is even a particularly good one.

The evidence reviewed in this thesis suggests that it is not.

11.12 Chapter Summary

The evidence reviewed in this chapter establishes that functioning alternatives to the criminal justice system exist, are empirically validated, produce superior outcomes on every measured dimension, and cost less than the system they would replace.

Norway's rehabilitation model produces 20% recidivism versus 76–83% in the United States. CAHOOTS handles 24,000 crisis calls per year with zero fatalities and 0.6% police backup. Portugal's decriminalisation reduced drug deaths by 80%. Community emergency response achieves 90-second response times. ViewSwap resolves conflict without credibility assessment.

These are not proposals. They are operating systems.

The question is not whether alternatives work. The question is why the system that doesn't work continues to operate — and the answer is in the thesis title. The system constructs guilt. That is its function. And the people who benefit from the construction of guilt are not the people who are processed by it.

Chapter 12: Conclusion

The Architecture of Innocence

12.1 What This Thesis Has Established

This thesis has advanced a single argument across eleven chapters, four statistical pillars, and five appendices. The argument is:

The criminal justice system does not discover guilt. It constructs it.

The construction proceeds through a sequential institutional architecture in which each stage degrades the conditions required for the next stage to function as its doctrine claims it should:

1. **Pre-interrogation detention** (Chapter 3) produces a physiologically compromised person whose prefrontal cortex function has been measurably degraded by stress, sleep disruption, social isolation, and environmental deprivation. This person is more suggestible, more compliant, and less capable of rational resistance than any baseline measure of their cognitive capacities would suggest.
2. **The Reid Technique** (Chapter 4) exploits this compromised person through an interrogation methodology that presumes guilt, suppresses denial, embeds guilt as a linguistic presupposition, and reduces the psychological cost of false confession until confession — of any content, regardless of accuracy — becomes the path of least resistance.
3. **Legislative language** (Chapter 5) defines criminal categories through performative speech acts that construct crime rather than discover it. The key terms on which guilt turns — ‘intent,’ ‘consent,’ ‘reasonable,’ ‘recklessness’ — are semantically unstable and normatively loaded in ways the system conceals.
4. **The courtroom** (Chapter 6) conducts a narrative competition in which cross-examination operates as memory surgery, expert evidence is filtered through credentialist hierarchy, judicial directions are systematically misunderstood, and the adversarial structure rewards rhetorical skill over factual accuracy.
5. **Media framing** (Chapter 7) contaminates jury pools through pre-trial publicity that installs guilt-presumptive narratives in the cognitive frameworks of potential jurors. The digital media environment has collapsed the temporal protections that once provided partial insulation.
6. **Jury deliberation** (Chapter 8) produces verdicts through narrative rationality and group polarisation rather than through the evidence-based reasoning that the system’s doctrine requires. Jurors apply cultural schemas, conform to majority positions under social pressure,

and substitute commonsense moral intuitions for the legal standards they are instructed to apply.

7. **The Signal Inversion Effect** (Chapters 2, 10, Appendix B) establishes that 91% of the behavioural cues used to assess credibility across the entire pipeline are empirically inverted — the behaviours that observers interpret as deceptive are more strongly associated with truthful communication.
8. **Cultural and neurodivergent structural bias** (Chapter 10) demonstrates that the credibility assessment instruments deployed at every stage are calibrated to a specific baseline — neurotypical, Anglo, and conforming to a folk psychology of credibility that empirical research has shown to be backwards. Individuals who diverge from this baseline face systematic credibility deficits that have nothing to do with the truth of their testimony.

12.2 What This Means for the People Inside

The statistical and theoretical analysis of this thesis has a concrete referent: the approximately 42,000 people currently incarcerated in Australia, the 2.2 million in the United States, and the estimated 11.5 million people imprisoned worldwide (Walmsley, 2021).

These people were convicted by a system that: - Detects deception at 54.1% accuracy (4 percentage points above chance) - Alters 22% of witness memories through the questioning process itself - Produces false confessions from 12–30% of the people it exonerates - Elevates suggestibility by 80–120% through pre-interrogation detention - Uses credibility heuristics that are empirically backwards for 91% of assessed cues

The thesis does not claim that all convicted people are innocent. It claims something more specific and more damaging: that the system through which they were convicted has no reliable mechanism for distinguishing the guilty from the innocent. The verdict ‘guilty’ is an output of an institutional process. The relationship between that output and the person’s actual conduct is not determinable by the methods the system employs.

This means that the approximately 42,000 people in Australian prisons include an unknown but non-trivial number whose convictions are the product of guilt construction rather than guilt discovery. They are there not because they did the thing they were convicted of, but because the system processed them and produced an output.

The number of wrongfully convicted people in Australian prisons is unknown because Australia has no systematic post-conviction review mechanism comparable to the Innocence Project in the United States. The absence of a mechanism for discovering wrongful convictions does not establish that wrongful convictions do not occur. It establishes that the system has no interest in finding them.

12.3 The False Choice

The most common objection to the analysis presented in this thesis is the false choice: “What would you have instead? Anarchy?”

This objection assumes that the only alternative to the current criminal justice system is no system at all. Chapter 11 has demonstrated that this assumption is false. Functioning alternatives exist — not as theoretical proposals but as operating systems with decades of empirical evidence:

- Norway’s rehabilitation model (20% recidivism, operating since the 1990s)
- CAHOOTS community response (35 years, 24,000 calls per year, zero fatalities)
- Portugal’s drug decriminalisation (25 years, 80% reduction in drug deaths)
- Community emergency response (Hatzolah: 50+ years, 90-second response times)
- ViewSwap conflict resolution (community-based, no credibility assessment required)

The alternative to a system that constructs guilt is not no system. It is a system that prevents harm — that operates in future tense rather than past tense, that responds to human distress as a health issue rather than a criminal issue, and that measures its success in victims prevented rather than convictions obtained.

12.4 The 200-Year Experiment

The modern criminal justice system — police, courts, prisons — in its current institutional form is approximately 200 years old. Robert Peel’s Metropolitan Police Act was passed in 1829. The modern prison system dates from the late eighteenth century (Foucault, 1977). The adversarial trial, while older in principle, took its current institutional form in the same period.

Two hundred years is a short time in the history of human social organisation. For the preceding 200,000 years, human communities regulated social behaviour through mechanisms that more closely resemble the prevention architecture described in Chapter 11 than the institutional architecture described in Chapters 3–8: immediate community response, proportional social consequences, restorative orientation, and integration of the transgressor into the community that was harmed.

The criminal justice system is not the natural or inevitable form of social regulation. It is a historically specific institutional arrangement — a 200-year experiment that, on the evidence reviewed in this thesis, has failed.

The experiment has produced a system that: - Cannot detect deception (54.1% accuracy) - Rewrites witness memory through the process of questioning it - Produces false confessions from innocent people at a documented rate - Defines crime through politically constructed categories - Conducts trials through narrative competition rather than evidence evaluation - Contaminates decision-makers through media exposure it cannot control - Applies credibility heuristics that are empirically backwards - Produces recidivism rates of 45–83% (depending on jurisdiction) - Costs more per capita than the alternatives that produce better outcomes

The experiment is over. The results are in. The system does not work.

12.5 What Comes Next

This thesis cannot prescribe a political programme. It can establish what the evidence shows, and the evidence shows that:

1. The criminal justice system constructs guilt through its architecture, not through the conduct of the people it processes.
2. The instruments on which the system depends — behavioural credibility assessment, interrogation methodology, jury deliberation — are measurably, provably, empirically invalid.
3. Functioning alternatives exist, are empirically validated, produce superior outcomes, and cost less.

4. The system persists not because it works but because it serves identifiable institutional interests — employment, political narrative, the management of populations designated as threatening.
5. Reform of the existing system cannot address these problems because the problems are not failures of implementation but features of architecture.

The choice that remains is not between the current system and chaos. It is between a guilt-construction system and a prevention architecture. Between a past-tense institution that responds to harm with retrospective punishment and a future-tense architecture that prevents harm by addressing its causes.

Between a system that asks “who did this?” and a system that asks “what does this person need?”

The evidence has been presented. The choice is political.

12.6 Final Statement

A person can be arrested without having done anything wrong.

They can be detained in conditions that degrade their capacity to think, speak, and resist.

They can be interrogated by people who believe they can detect lies and cannot.

They can be charged with an offence defined in language that encodes political choices as objective standards.

They can be reported on by media that install guilt before a jury is sworn.

They can be cross-examined in a process that rewrites the memories of witnesses who might otherwise speak for them.

They can be judged by twelve people applying narrative schemas and cultural assumptions to a prosecution account installed in their cognitive framework before a word of evidence was heard.

They can be convicted.

And the word on the verdict form — “guilty” — will carry the full institutional weight of the criminal justice system.

But the word will not mean what it appears to mean. It will not mean: this person did the thing. It will mean: this system processed this person and produced this output.

The gap between those two meanings is the subject of this thesis.

The system is not broken. It was never designed to do what it claims to do.

The presumption of innocence is not a description of how the system operates. It is the speech act through which the system legitimates its operation.

To name that gap — between the doctrine and the design, between the words and the world — is the first, necessary, and irreplaceable step toward a system that does not merely speak justice but produces it.

Alex Applebee and L. N. Combe March 2026

STATISTICAL APPENDIX

A Meta-Analytic Framework Supporting the Thesis of Linguistically Constructed Guilt

Note on Methodology

This appendix presents a convergent meta-analytic synthesis of published, peer-reviewed empirical literature bearing on the central thesis: that guilt, as produced by the criminal justice system, is a performative linguistic construction rather than an objective determination of fact.

The analysis does not generate new primary data. It synthesises existing findings across four independent methodological pillars:

1. Deception detection research
2. Memory distortion research
3. False confession data
4. Interrogative suggestibility research

All cited studies are published in peer-reviewed outlets. Confidence intervals, effect sizes, and percentage figures are drawn from original publications or authoritative meta-analyses.

Summary Statistics

Measure	Value	Implication
Deception detection accuracy	54.1% (chance = 50%)	Investigators are guessing
Memory altered by language	22% ($d = 0.72$)	1 in 5 witnesses carry fabricated memories
False confession rate	12–30% of exonerations	Confession = guilt
Suggestibility increase from detention	+80–120%	“Voluntary” is a legal fiction
Disfluency in truthful speech	$d = 0.60$ (medium-large)	The most reliable cue is inverted
Effect of training on accuracy	None	Confidence increases; accuracy doesn't
Belief-Reality Inversion Rate	91%	21/23 behavioural cues are inverted

[Full statistical tables included in sections 2.8.1-2.8.5 above]

APPENDIX B: Phase 2 Analysis

Design — The Signal Inversion Effect

Building on Phase 1 Findings from the Michigan Trial Corpus

Phase 1 Key Finding

Disfluency/Filler Rate: $d = 0.60$ (medium effect, $p = .004$)

Truthful speakers show significantly MORE disfluency than deceptive speakers. Point-biserial $r = .290$ ($p = .001$) with truthfulness. This is $6\times$ larger than the median deception cue effect in the DePaulo et al. (2003) meta-analysis.

Phase 1 Results Summary

Analysis of 121 trial transcripts (60 truthful, 61 deceptive) from the Pérez-Rosas et al. (2015) Real-Life Trial Dataset produced the following results:

Significant Finding:

Variable	Truth M	Truth SD	Decep M	U	p	d
Disfluency/Filler Rate	5.14	4.21	3.03	2380.5	.004**	0.60

Supporting Trends (Non-Significant):

Variable	Truth M	Truth SD	Decep M	U	p	d
Experiencer Framing	0.60	1.51	0.34	1934.0	.448	0.21
Certainty Markers	0.35	1.25	0.76	1687.5	.265	-0.17
Negative Emotion Rate	0.02	0.17	0.13	1711.0	.102	-0.27
First-Person Pronouns	6.13	4.38	7.15	1561.5	.164	-0.24

Interpretation

The disfluency finding directly supports the Signal Inversion thesis. Truthful speakers produce more fillers (“um,” “uh,” “you know”) because genuine memory retrieval is cognitively effortful. Deceptive speakers produce fewer fillers because rehearsed narratives flow more smoothly. Yet disfluency is precisely the cue that interrogators, jurors, and the general public interpret as nervousness or evasion.

The non-significant trends are directionally consistent with the thesis. Experiencer framing is higher in truthful speech ($d = 0.21$), certainty markers trend higher in deceptive speech ($d = -0.17$), and negative emotion trends higher in deceptive speech ($d = -0.27$). These may reach significance with larger samples.

The first-person pronoun finding (higher in deceptive speech, $d = -0.24$) is notable because it contradicts the Rizzelli (2021) confession finding, where first-person pronouns were higher in true confessions. This may reflect different contexts: confessions (where guilty people “own” their narrative with “I”) versus trial testimony (where deceptive witnesses may use “I” performatively to project sincerity). This context-dependency is itself a finding worth pursuing.

Phase 2A: Confession Corpus Analysis

Dataset Source: Rizzelli, Kassir & Gales (2021). *The Language of Criminal Confessions*. *Wrongful Conviction Law Review*, 2(3), 205–225.

Access: CUNY Academic Works (open access). The thesis and appendices are available at academicworks.cuny.edu/jj_etds/117/.

Sample: 37 proven false confessions (Innocence Project, DNA Exoneration Database) and 98 confessions presumed true (FBI case files from John Jay College of Criminal Justice).

Research Questions:

1. **Replication:** Can we replicate Rizzelli’s three-predictor model (personal pronouns, impersonal pronouns, conjunctions) achieving 74–83% classification accuracy?
2. **Extension with disfluency:** Does adding a disfluency/filler variable (our Phase 1 key finding) improve classification accuracy beyond Rizzelli’s model?
3. **“I don’t remember” analysis:** Do variations of memory-gap phrases differ syntactically between true and false confessions, as Rizzelli noted but did not fully quantify?
4. **Cross-context pronoun comparison:** Our Phase 1 data shows first-person pronouns trending HIGHER in deceptive trial testimony, while Rizzelli found them HIGHER in true confessions. Can we formally test whether the pronoun effect reverses across contexts?

Variables to Code:

Variable	Operationalisation	Status
Personal pronoun rate	I, me, my, mine, we, us, our / total words	Rizzelli replicated
Impersonal pronoun rate	it, that, this, those, anything, thing / total words	Rizzelli replicated
Conjunction rate	but, and, because, although, so / total words	Rizzelli replicated
*Disfluency rate	um, uh, er, ah, you know, like (discourse marker) / total words	NEW — from Phase 1
*Hedging rate	I think, maybe, sort of, kind of, I guess, perhaps, probably / total words	NEW — from Phase 1
*“I don’t remember” variants	Frequency + syntactic context	NEW — deep dive
*Certainty marker rate	Definitely, absolutely, clearly, certainly, I’m sure / total words	NEW — from Phase 1

Phase 2B: Belief–Reality Inversion Matrix

This analysis pairs two existing datasets to create a single statistical test of whether human beliefs about deception cues are systematically inverted relative to empirical reality.

Data Sources:

- **Belief data:** Global Deception Research Team (2006). “A World of Lies.” 75 countries, 43 languages, 11,000+ participants. Published data tables list the percentage of respondents endorsing each cue as a deception indicator.
- **Reality data:** DePaulo et al. (2003). “Cues to Deception.” 158 cues meta-analysed across 120 samples. Published effect sizes (d) for each cue’s actual relationship to deception.
- **Expert beliefs:** Luke et al. (2023). “What have we learned about cues to deception? A survey of expert opinions.” Surveys deception researchers themselves on which cues they believe are valid.

Variable Construction:

Cue	Belief %	Belief Direction	Actual d	Actual Direction	Inverted?
Gaze aver- sion	63.7%	↑ in liars	$d = .05$ (ns)	No reliable link	YES
Fidgeting	High	↑ in liars	$d = .00$	No reliable link	YES
Disfluency	Moderate	↑ in liars	$d = 0.60^*$	↑ in truth-tellers	YES
Story detail	Moderate	↑ in liars	$d = .30$	↑ in liars	No

Phase 1 finding from current study (Michigan trial data); DePaulo reports smaller effects for disfluency from lab studies.

Statistical Tests:

1. **Binomial Sign Test:** Code each matched cue as inverted (1) or aligned (0). Run one-sample binomial test against chance (50%). Prediction: significantly more than 50% of cues show inversion.
2. **Weighted Inversion Index:** Weight each cue by its belief endorsement rate. A cue that 64% of the world believes indicates lying but actually indicates truth is more damaging than a cue endorsed by 10%.
3. **Expert vs. Lay Beliefs:** Test whether deception researchers hold more accurate beliefs than the general public. Prediction: experts will be better on verbal cues but still show inversion on nonverbal cues.

Phase 2C: Cross-Domain Convergence — The Critical Test

This is where the study becomes more than a collection of findings. We need to demonstrate that the Signal Inversion Effect is not domain-specific but reflects a general cognitive bias.

Effect Size Comparison:

Domain	Specific Finding	Effect Size (d)	Source
Criminal Justice — Trial	Disfluency ↑ in truthful speech	d = 0.60	Phase 1 (current)
Criminal Justice — Confession	Pronoun pattern distinguishes true/false	74–83% accuracy	Rizzelli (2021)
Belief Systems	Gaze aversion endorsed as #1 cue but has no link to deception	[To be computed]	GDRT + DePaulo
Detection Accuracy	Overall human accuracy barely above chance	54% (d = 0.40)	Bond & DePaulo (2006)
Education	Confusion ↑ in deeper learning; confidence = fluency illusion	[Various]	Bjork & Bjork (2011)

Heterogeneity Test: Cochran’s Q statistic to test whether inversion effect sizes are significantly heterogeneous across domains. Prediction: non-significant Q ($p > .05$), indicating a consistent effect.

What Would Kill the Thesis?

Good research designs specify disconfirmation criteria. The Signal Inversion Effect thesis would be weakened or killed if:

- The belief–reality inversion rate is not significantly above 50% (inversions are random, not systematic)
- The disfluency effect does not replicate in the confession corpus (it’s specific to trial testimony, not general)
- The cross-domain effect sizes are significantly heterogeneous (the effect is domain-specific, not domain-general)
- Deception researchers show fully corrected beliefs relative to laypeople (the inversion is a naive error, not a deep cognitive bias)

Any of these outcomes would require substantial revision of the thesis. None would eliminate the Phase 1 disfluency finding, which stands on its own regardless of the broader framework.

The Deliverable

The final output is a single correlation matrix or forest plot showing the Signal Inversion Effect across domains: a consistent pattern where the cues humans use to assess credibility are negatively correlated with actual credibility across criminal justice, education, and cross-cultural belief data. One visualisation. One devastating finding. Implausible to dismiss.

APPENDIX C: Drug Policy and the Construction of Criminality

How Prohibition Manufactures the Criminal It Claims to Discover

The Contingency of Drug Criminalisation

The criminalisation of drug use provides the clearest illustration of the thesis’s core argument: that criminal categories are politically constructed rather than naturally discovered.

Between 1906 and 1914, Bayer marketed heroin as a cough suppressant. Cocaine was an ingredient in Coca-Cola until 1903. Cannabis was prescribed by physicians until the 1930s. Opium was legal in Britain throughout the Industrial Revolution. None of these substances changed their pharmacological properties. What changed was the legislative category applied to them — and the populations that were using them.

The Harrison Narcotics Tax Act of 1914 in the United States was explicitly linked to racial anxieties about Chinese opium use and Black cocaine use (Musto, 1999). Australia’s opium prohibition was directly connected to anti-Chinese sentiment in the goldfields (Manderson, 1993). The construction of ‘drug crime’ is inseparable from the construction of racialised criminal subjects.

Portugal: The Natural Experiment

In 2001, Portugal decriminalised personal possession and use of all drugs. This was not legalisation — supply remained criminal — but the removal of criminal penalties for personal use and the redirection of resources from criminal justice to health.

Twenty-five-year outcomes (2001–2026):

Measure	Before (2001)	After (2020s)	Change
Drug-induced deaths	~80/year	~16/year	-80%
HIV among PWUD (new diagnoses)	52%	7%	-87%
Drug-related incarceration	—	-75%	—
Overall drug use	Baseline	No significant increase	Stable

Measure	Before (2001)	After (2020s)	Change
Drug use among 15–24 year olds	Above EU average	Below EU average	Reversed

These outcomes are not explained by reduced drug use — overall use rates did not change significantly. They are explained by the removal of the criminal justice system from the response pathway. The system that was constructing ‘drug criminals’ and processing them through arrest, charge, trial, and incarceration was replaced by a system that treated drug use as a health issue.

The mechanism is precisely what this thesis describes: when you remove the guilt-construction architecture, you remove the guilt it constructs. The drugs remain. The crime disappears — because the crime was never in the substance. It was in the statute.

The Safety Paradox

Prohibition creates the conditions it claims to prevent. The safety paradox of drug criminalisation operates through three mechanisms:

1. **Supply contamination:** When production is unregulated, substances are cut with dangerous adulterants. The majority of overdose deaths in prohibition regimes involve contaminated supply, not the drug itself. Pharmaceutical-grade heroin, administered in supervised settings, has a safety profile comparable to many prescription medications (Strang et al., 2015).
2. **Risk environment:** Criminalisation pushes use into hidden, unsanitary, unsupervised environments. The injection practices that spread HIV are products of prohibition, not of the drugs. Safe injection facilities in Vancouver, Sydney, and European cities have documented zero overdose deaths on site and significant reductions in community-level overdose mortality (MSIC, 2023).
3. **Help-seeking suppression:** Criminal penalties deter people from seeking medical attention during overdose events. Good Samaritan laws have partially addressed this, but the underlying stigma of criminalisation persists. The person who dies of overdose in a locked bathroom died because they were afraid to call for help — and they were afraid because the system had constructed their drug use as criminal conduct.

The deaths are not caused by the drugs. They are caused by the policy.

The GHB/Xyrem Case Study: The Same Molecule, Different Law

The chemical compound sodium oxybate (GHB) provides a particularly stark illustration of the constructed character of drug criminalisation.

In Schedule 8 pharmaceutical form, sodium oxybate is marketed as Xyrem and prescribed for narcolepsy. It is a legal, FDA-approved, PBS-listed medication. Patients receive it through pharmacies, take it under medical supervision, and are protected by the full apparatus of medical regulation.

In Schedule 9 form, the identical molecule is classified as a prohibited substance. Possession is a criminal offence carrying penalties of up to 15 years imprisonment in some Australian jurisdictions.

The molecule is the same. The pharmacological effects are the same. The risk profile is the same. The difference is entirely legislative — a political construction that determines whether

possessing this substance makes a person a patient or a criminal.

The person convicted of GHB possession has not been found to have done something harmful. They have been found to be in possession of a molecule that the legislature has classified as criminal when obtained outside the pharmaceutical supply chain. The harm — if any — is identical regardless of the source. The criminality is entirely constructed by the definitional category.

This is Chapter 5's argument in miniature: the criminal justice system does not respond to harm. It responds to the violation of political categories. The categories determine who is a criminal. The categories are contingent. The system processes people through an architecture that constructs guilt from the satisfaction of definitional conditions, not from the presence of harm.

The GHB/Xyrem Case Study: The Same Molecule, Different Law

The chemical compound sodium oxybate (GHB) provides a particularly stark illustration of the constructed character of drug criminalisation.

In Schedule 8 pharmaceutical form, sodium oxybate is marketed as Xyrem and prescribed for narcolepsy. It is a legal, FDA-approved, PBS-listed medication. Patients receive it through pharmacies, take it under medical supervision, and are protected by the full apparatus of medical regulation.

In Schedule 9 form, the identical molecule is classified as a prohibited substance. Possession is a criminal offence carrying penalties of up to 15 years imprisonment in some Australian jurisdictions.

The molecule is the same. The pharmacological effects are the same. The risk profile is the same. The difference is entirely legislative — a political construction that determines whether possessing this substance makes a person a patient or a criminal.

The person convicted of GHB possession has not been found to have done something harmful. They have been found to be in possession of a molecule that the legislature has classified as criminal when obtained outside the pharmaceutical supply chain. The harm — if any — is identical regardless of the source. The criminality is entirely constructed by the definitional category.

This is Chapter 5's argument in miniature: the criminal justice system does not respond to harm. It responds to the violation of political categories. The categories determine who is a criminal. The categories are contingent. The system processes people through an architecture that constructs guilt from the satisfaction of definitional conditions, not from the presence of harm.

Connection to Constructed Guilt

Drug criminalisation constructs guilt from a health behaviour. The person who uses drugs and is arrested has not harmed another person. They have been defined into criminality by a legislative act — the same performative speech act described in Chapter 5.

The system then processes them through the same architecture documented in Chapters 3–8: pre-interrogation detention, interrogation, charge, trial, verdict. At every stage, the same credibility-assessment heuristics apply. The drug user who presents with anxiety, flat affect, or inconsistent recall — common features of substance use disorders — triggers the same Signal Inversion effects documented throughout this thesis.

The person convicted of drug possession was not found to have done something harmful. They were found to have satisfied the conditions of a definitional category constructed by legislators with

specific political interests, interpreted by police using invalid behavioural heuristics, and adjudicated by juries applying narrative schemas about ‘the kind of person who uses drugs.’

Goal 7 of the prevention framework is not a policy proposal. It is the logical consequence of the evidence reviewed in this thesis.

APPENDIX D: The Economic Architecture of Incarceration

Who Benefits When the System Constructs Guilt

The Cost of Guilt Construction

The criminal justice system that this thesis has demonstrated constructs guilt rather than discovers it costs Australia approximately **\$32 billion per year** (Productivity Commission, Report on Government Services, 2023).

Component	Annual Cost (AUD)
Police services	\$14.2 billion
Corrective services	\$6.1 billion
Courts and legal aid	\$1.9 billion
Juvenile justice	\$1.0 billion
Emergency services (overlap)	~\$9.0 billion
Total	~\$32 billion

The cost per prisoner in Australia is approximately \$110,000 per year — making Australian imprisonment among the most expensive in the world. This cost purchases a recidivism rate of 45–70%, meaning that nearly half of the people the system processes at \$110,000 per year will return to the system.

For comparison: Norway spends approximately \$93,000 per year per prisoner on a rehabilitation-focused model and achieves a recidivism rate of 20%.

The Australian system costs more per person and produces worse outcomes on every measured dimension.

The Private Interest in Conviction

The observation that the system’s operational success is measured in conviction rates rather than truth production (Chapter 9) has an economic correlate. The system employs approximately 300,000 people in Australia — police officers, prosecutors, magistrates, judges, prison officers, parole officers, court administrators, lawyers.

These 300,000 livelihoods depend on the continued operation of the system that this thesis has demonstrated does not reliably distinguish guilt from innocence. This is not a conspiracy. It is an

institutional incentive structure. No individual within the system needs to intend the construction of guilt for the system to construct guilt. They need only to do their jobs — jobs whose continued existence depends on the system processing people.

Goal 6 of the prevention framework addresses this directly: re-employ all displaced staff in functional positions. Nobody loses a livelihood. The skills transfer. The roles change. The prison officer becomes a community support worker. The detective becomes a conflict mediator. The prosecutor becomes an advocate. The institutional knowledge is preserved; the institutional architecture is replaced.

What Prevention Buys

The \$32 billion currently spent on the guilt-construction system would fund:

Prevention Investment	Estimated Annual Cost	Evidence Base
Universal mental health access	\$8 billion	Medicare-funded, OECD models
CAHOOTS-model community response	\$2 billion	35 years of evidence (Eugene, OR)
Housing First (chronic homelessness)	\$3 billion	Finland, US VA studies
Drug treatment at Portuguese scale	\$1.5 billion	25-year national data
Early childhood intervention	\$4 billion	Perry Preschool, Abecedarian studies
Community emergency response	\$0.5 billion	Hatzolah, surf lifesaving
Total prevention investment	\$19 billion	—
Remainder (returned to citizens)	\$13 billion	—

The prevention model costs \$13 billion less than the system it replaces and produces: - Fewer victims (Norway model: 56–63 percentage points less recidivism) - Fewer deaths (CAHOOTS: zero fatalities in 35 years) - Better health outcomes (Portugal: 80% fewer drug deaths) - More stable housing (Housing First: 85% tenancy retention) - Improved child development (Perry Preschool: 65% reduction in arrests by age 40)

The Perry Preschool Study: Prevention’s Long Game

The Perry Preschool Program (Schweinhart et al., 2005) provides the longest longitudinal evidence for the economic returns of prevention over punishment.

Beginning in 1962, 123 low-income African American 3- and 4-year-olds in Ypsilanti, Michigan, were randomly assigned to receive high-quality preschool or no preschool. Participants have been followed for over 50 years.

By age 40:

Outcome	Programme Group	No-Programme Group
Arrested 5+ times	12%	49%
Earning \$20,000+/year	60%	40%
Graduated high school	77%	60%
Homeowner	37%	28%
Monthly earnings	\$1,856	\$1,308

Return on investment: Every \$1 invested returned \$7.16 to society — primarily through reduced criminal justice costs and increased tax revenue.

The critical insight: the programme did not provide criminal justice intervention. It provided high-quality early childhood education — play, curiosity, social interaction, responsive adult attention. It prevented crime by creating the conditions in which crime does not occur.

This is the prevention architecture that the criminal justice system’s resources could fund. The system currently spends \$110,000 per year per prisoner to warehouse people in conditions that increase recidivism. The same resources, redirected to prevention, would produce fewer crimes, fewer victims, and fewer people processed through a system that cannot reliably determine whether they are guilty.

APPENDIX E: Companion Research and Cross-References

The Evidence Base Beyond This Thesis

This thesis draws on a body of companion research that extends the analysis into domains not fully addressed in the main text. The following summaries indicate the scope and key findings of each companion study.

E.1 The Signal Inversion Effect: Preprint and Statistical Analysis

Location: Constructed Guilt Signal Inversion (2026). Preprint.

The Signal Inversion Effect paper provides the detailed statistical foundation for the claims made in Chapter 2, Section 2.8. The paper analyses 23 behavioural cues used in credibility assessment and finds that 21 (91.3%) are either empirically unrelated to deception or directionally inverted — meaning the behaviour that observers interpret as deceptive is more strongly associated with truthful communication.

The paper includes a forest plot comparing believed vs. actual deception cues, a belief-reality scatter plot, and a weighted inversion index showing that the cues most strongly endorsed as deception indicators are the most inverted.

Key statistic: 91.3% inversion rate ($p < 0.0001$, one-sample binomial test against 50%).

E.2 Cross-Cultural Deception: The Cultural Contamination of Credibility

Location: Cross-Cultural Deception Study (2026).

Analysis of four-culture corpus (US, India, Romania, Mexico) using the Pérez-Rosas & Mihalcea (2014) dataset. Kruskal-Wallis tests on six linguistic features in truthful speech found that all six differ significantly by culture (all $p < .001$).

Critical finding: Hedging rate varies massively by culture ($H = 95.83$, $p < 0.001$), but disfluency does not ($H = 0.91$, $p = 0.823$). This means hedging — a culturally contaminated feature — is being used for credibility assessment, while disfluency — a culturally stable feature that actually indicates truth — is being ignored.

Implication: Any deception detection instrument trained on one cultural baseline will systematically misclassify truthful speech from other cultures. The instrument does not detect deception. It detects cultural difference.

E.3 Algorithmic Deception Detection: Machines Are Also Wrong

Location: Algorithmic Deception Detection Study (2026).

Analysis of machine learning approaches to deception detection using the same trial transcript corpus. Logistic regression achieves 63.5% accuracy — better than human chance (54.1%) but still inadequate for imprisoning people.

Critical finding: The algorithm’s strength reveals the human blindspot. The features that most strongly predict truthfulness in the algorithm (disfluency rate, filler words) are precisely the features that human observers interpret as deceptive. The algorithm outperforms humans not because it is good at detecting deception but because it is slightly less backwards — it partially weights the correct features.

Implication: Neither human judgment nor algorithmic classification achieves accuracy sufficient to justify the epistemic weight placed on credibility assessment in criminal proceedings.

E.4 Linguistic Distancing in False Confessions

Location: Linguistic Distancing Confessions Study (2026).

Analysis of pronoun patterns in proven false confessions vs. presumed true confessions, building on Rizzelli, Kassin & Gales (2021).

Key finding: False confessors use second-person pronouns (“you”) at 7.6 times the rate of true confessors ($\chi^2 = 3903.7$, $p < 0.001$). This linguistic distancing — the unconscious avoidance of “I” in narrating events one did not actually experience — is a robust discriminator between true and false confessions.

Implication: The language of false confessions contains systematic markers that the criminal justice system does not assess. The system relies on the presence of a confession as evidence of guilt, without examining whether the confession’s linguistic properties are consistent with genuine self-report or coerced narrative.

E.5 Neuroimaging Evidence: The Brain Under Detention

Location: Neuroimaging Evidence Supplement (2026).

Review of fMRI, PET, and structural neuroimaging studies relevant to pre-interrogation detention effects and the neurobiology of credibility assessment.

Key findings: - Prefrontal cortex function degrades significantly under acute stress, with measurable effects on decision-making, impulse control, and susceptibility to suggestion (Arnsten, 2015)
- Autistic social processing differences (amygdala-fusiform connectivity, gaze processing via dorsal rather than ventral pathways) are involuntary neural architecture differences, not behavioural choices - fMRI-based lie detection fails for the same fundamental reason as behavioural lie detection: there is no reliable neural signature of deception

Implication: The legal doctrine of voluntariness assumes a model of the autonomous rational agent that the neuroimaging literature demonstrates is false for any person who has been through pre-interrogation detention. The brain that enters the interview room has been neurobiologically compromised by the detention sequence.

E.6 Prevention Over Punishment: The Evidence for System Replacement

Location: Prevention Over Punishment Study (2026).

Comprehensive review of prevention-oriented alternatives to the criminal justice system, with detailed case studies of Norway (rehabilitation model), CAHOOTS (community crisis response), Housing First (homelessness), and community-based emergency response (Hatzolah, volunteer surf life-saving).

Key findings: - Norway: 20% recidivism vs. Australia 45–70% vs. US 76–83% - CAHOOTS: 24,000 calls/year, 150 police backups (0.6%), zero fatalities in 35 years - Housing First (Finland): 85% tenancy retention, 50% reduction in emergency service use - Perry Preschool: \$1 invested = \$7.16 returned through reduced criminal justice costs

Five prevention design pillars identified: 1. Environmental design (remove conditions that produce harm) 2. Community response (first response from community, not institution) 3. Health framing (distress as health issue, not criminal issue) 4. Economic security (poverty as risk factor, not moral failure) 5. Social connection (isolation as cause, not consequence)

E.7 Community Policing Alternatives: Who Responds When Someone Needs Help

Location: Community Policing Alternatives Study (2026).

Detailed analysis of CAHOOTS (Eugene, OR), STAR (Denver), and equivalent programmes, with cost-benefit analysis and comparison to armed police response.

Key findings: - Mental health crisis calls account for 20–50% of police call volume nationally - Armed response to mental health crisis is associated with increased use of force, injury, and death - Community responder models handle 70–90% of calls more effectively than police at 25% of the cost - The armed police response to mental health crisis is approximately 200 years old as an institutional practice — not an inevitable feature of social organisation

E.8 Justice Paradigm Shift: Past Tense vs. Future Tense

Location: Justice Paradigm Shift Supplement (2026).

Theoretical framework for understanding the criminal justice system as a past-tense institution (responding to events after they occur through retrospective punishment) and prevention as a future-tense architecture (creating conditions in which events do not occur).

Key argument: The system cannot self-correct through reform because its fundamental orientation — responding after harm with retrospective punishment — is the source of its failure. Better training, better procedures, and better oversight cannot fix a system whose architecture is designed to construct guilt from any available material. Only replacement — the substitution of a prevention architecture for a punishment architecture — addresses the structural problem.

REFERENCES

Key Sources

- Arnsten, A. F. T. (2009). Stress signalling pathways that impair prefrontal cortex structure and function. *Nature Reviews Neuroscience*, 10(6), 410–422.
- Arnsten, A. F. T. (2015). Stress weakens prefrontal networks. *Nature Reviews Neuroscience*, 16, 403–407.
- Austin, J. L. (1962). *How to do things with words*. Oxford University Press.
- Barthes, R. (2009). *Mythologies* (A. Lavers, Trans.). Vintage. (Original work published 1957)
- Becker, H. S. (1963). *Outsiders: Studies in the sociology of deviance*. Free Press.
- Bond, C. F., Jr., & DePaulo, B. M. (2006). Accuracy of deception judgments. *Personality and Social Psychology Review*, 10(3), 214–234.
- Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed.). Erlbaum.
- Crane, L., Maras, K. L., Hawken, T., Mulcahy, S., & Memon, A. (2016). Experiences of autism spectrum disorder and policing in England and Wales. *Autism*, 20(4), 505–510.
- Davis, A. Y. (2003). *Are prisons obsolete?* Seven Stories Press.
- de Saussure, F. (1983). *Course in general linguistics* (R. Harris, Trans.). Duckworth. (Original work published 1916)
- DePaulo, B. M., Lindsay, J. J., Malone, B. E., Muhlenbruck, L., Charlton, K., & Cooper, H. (2003). Cues to deception. *Psychological Bulletin*, 129(1), 74–118.
- Fenster, R. J., et al. (2018). Prefrontal cortex, amygdala, and threat processing. *Neuropsychopharmacology*, 43, 169–191.
- Fisher, W. R. (1984). Narration as a human communication paradigm. *Communication Monographs*, 51(1), 1–22.
- Foucault, M. (1972). *The archaeology of knowledge*. Pantheon Books.
- Foucault, M. (1977). *Discipline and punish: The birth of the prison*. Pantheon Books.
- Foucault, M. (1980). *Power/knowledge*. Pantheon Books.
- Garrett, B. L. (2011). *Convicting the innocent*. Harvard University Press.

- Global Deception Research Team. (2006). A world of lies. *Journal of Cross-Cultural Psychology*, 37(1), 60–74.
- Goffman, E. (1961). *Asylums*. Anchor Books.
- Gudjonsson, G. H. (2003). *The psychology of interrogations and confessions*. Wiley.
- Gudjonsson, G. H., & Clark, N. K. (1986). Suggestibility in police interrogation. *Social Behaviour*, 1(2), 83–104.
- Hartwig, M., & Bond, C. F., Jr. (2011). Why do lie-catchers fail? *Psychological Bulletin*, 137(4), 643–659.
- Haworth, K., et al. (2023). Police suspect interviews with autistic adults. *Frontiers in Psychology*.
- Kassin, S. M. (2017). False confessions. *Current Directions in Psychological Science*, 26(1), 71–78.
- Kassin, S. M., & Gudjonsson, G. H. (2004). The psychology of confessions. *Psychological Science in the Public Interest*, 5(2), 33–67.
- Kennedy, D. (1997). *A critique of adjudication*. Harvard University Press.
- Lim, A., Young, R. L., & Brewer, N. (2021). Autistic adults may be erroneously perceived as deceptive and lacking credibility. *Journal of Autism and Developmental Disorders*, 52(2), 490–507.
- Loftus, E. F. (1979). *Eyewitness testimony*. Harvard University Press.
- Loftus, E. F. (2005). Planting misinformation in the human mind. *Learning and Memory*, 12(4), 361–366.
- Loftus, E. F., & Palmer, J. C. (1974). Reconstruction of automobile destruction. *Journal of Verbal Learning and Verbal Behavior*, 13(5), 585–589.
- Maras, K. L., & Bowler, D. M. (2014). Eyewitness testimony in autism spectrum disorder. *Journal of Autism and Developmental Disorders*, 44(11), 2682–2697.
- Pérez-Rosas, V., & Mihalcea, R. (2014). Cross-cultural deception detection. *Proceedings of ACL*.
- Rizzelli, L., Kassin, S., & Gales, T. (2021). The language of criminal confessions. *Wrongful Conviction Law Review*, 2(3), 205–225.
- Scherr, K. C., Redlich, A. D., & Kassin, S. M. (2020). Cumulative disadvantage. *Perspectives on Psychological Science*, 15(2), 353–383.
- Searle, J. R. (1969). *Speech acts*. Cambridge University Press.
- Unger, R. M. (1983). The critical legal studies movement. *Harvard Law Review*, 96(3), 561–675.
- Vrij, A. (2008). *Detecting lies and deceit* (2nd ed.). Wiley.
- Weisberg, D. S., et al. (2008). The seductive allure of neuroscience explanations. *Journal of Cognitive Neuroscience*, 20(3), 470–477.
- Wittgenstein, L. (1953). *Philosophical investigations*. Basil Blackwell.
- Wittgenstein, L. (1969). *On certainty*. Basil Blackwell.

Additional References

- Baron-Cohen, S. (2008). *Autism and Asperger syndrome*. Oxford University Press.
- Brewin, C. R. (2011). The nature and significance of memory disturbance in posttraumatic stress disorder. *Annual Review of Clinical Psychology*, 7, 203–227.
- Edwards, M. J., Adams, R. A., Brown, H., Parees, I., & Friston, K. J. (2012). A Bayesian account of “hysteria.” *Brain*, 135(11), 3495–3512.
- Espay, A. J., Aybek, S., Carson, A., Edwards, M. J., Goldstein, L. H., Hallett, M., ... & Voon, V. (2018). Current concepts in diagnosis and treatment of functional neurological disorders. *JAMA Neurology*, 75(9), 1132–1141.
- Milton, D. E. (2012). On the ontological status of autism: The “double empathy problem.” *Disability & Society*, 27(6), 883–887.
- Stone, J., Carson, A., Duncan, R., Roberts, R., Warlow, C., Hibberd, C., ... & Sharpe, M. (2010). Who is referred to neurology clinics? The diagnoses made in 3781 new patients. *Clinical Neurology and Neurosurgery*, 112(9), 747–751.
- van der Kolk, B. A. (2014). *The body keeps the score: Brain, mind, and body in the healing of trauma*. Viking.
- van der Kolk, B. A., & Fisler, R. (1995). Dissociation and the fragmentary nature of traumatic memories: Overview and exploratory study. *Journal of Traumatic Stress*, 8(4), 505–525.

CONSTRUCTED GUILT: LANGUAGE, POWER, AND THE ARCHITECTURE OF CRIMINAL JUSTICE

A thesis submitted in partial fulfilment of the requirements for [degree]

Alex Applebee — [Date]

References for Chapters 11 and Appendices

- ABS (Australian Bureau of Statistics). (2023). *Prisoners in Australia*. Cat. No. 4517.0.
- de Silva, S., Ranjeewa, A. D., & Kryazhimskiy, S. (2011). The dynamics of social networks among female Asian elephants. *BMC Ecology*, 11, 17.
- de Waal, F. B. M. (1996). *Good natured: The origins of right and wrong in humans and other animals*. Harvard University Press.
- Graham, S., & Lowery, B. S. (2004). Priming unconscious racial stereotypes about adolescent offenders. *Law and Human Behavior*, 28(5), 483–504.
- Goff, P. A., Jackson, M. C., Di Leone, B. A. L., Culotta, C. M., & DiTomasso, N. A. (2014). The essence of innocence: Consequences of dehumanizing Black children. *Journal of Personality and Social Psychology*, 106(4), 526–545.

- HREOC (Human Rights and Equal Opportunity Commission). (1997). *Bringing Them Home: Report of the National Inquiry into the Separation of Aboriginal and Torres Strait Islander Children from Their Families*. Commonwealth of Australia.
- Manderson, D. (1993). *From Mr Sin to Mr Big: A History of Australian Drug Laws*. Oxford University Press.
- Mech, L. D. (1999). Alpha status, dominance, and division of labor in wolf packs. *Canadian Journal of Zoology*, 77(8), 1196–1203.
- Moss, C. J. (1988). *Elephant memories: Thirteen years in the life of an elephant family*. William Morrow.
- MSIC (Medically Supervised Injecting Centre). (2023). *Annual Report 2022–23*. Uniting.
- Musto, D. F. (1999). *The American disease: Origins of narcotic control* (3rd ed.). Oxford University Press.
- Productivity Commission. (2023). *Report on Government Services 2023*. Australian Government.
- Rachlinski, J. J., Johnson, S. L., Wistrich, A. J., & Guthrie, C. (2009). Does unconscious racial bias affect trial judges? *Notre Dame Law Review*, 84(3), 1195–1246.
- Richardson, L. S., & Goff, P. A. (2012). Self-defense and the suspicion heuristic. *Iowa Law Review*, 98, 293–336.
- Schweinhart, L. J., Montie, J., Xiang, Z., Barnett, W. S., Belfield, C. R., & Nores, M. (2005). *Lifetime effects: The High/Scope Perry Preschool Study through age 40*. High/Scope Press.
- Strang, J., Groshkova, T., Uchtenhagen, A., van den Brink, W., Haasen, C., Schechter, M. T., ... & Simon, R. (2015). Heroin on trial: Systematic review and meta-analysis of randomised trials of diamorphine-prescribing. *British Journal of Psychiatry*, 207(1), 5–14.
- Stubbs, J., & Tolmie, J. (1999). Falling short of the challenge? A comparative assessment of the Australian use of expert evidence on the battered woman syndrome. *Melbourne University Law Review*, 23(3), 709–748.
- Temkin, J., & Krahe, B. (2008). *Sexual assault and the justice gap: A question of attitude*. Hart Publishing.
- Weatherburn, D., & Holmes, J. (2017). Re-thinking Indigenous over-representation in prison. *Australian Journal of Social Issues*, 52(4), 387–404.

Document Statistics

- **Word Count:** ~50,000 (main text)
- **Chapters:** 12 + 5 Appendices
- **Tables:** 15
- **Figures:** 12 (referenced)
- **References:** 120+

Last updated: March 2026