

# [CLAIR]

Comprehensible LLM AI Intermediate Representation

*[A Formalization of Epistemic Reasoning for Artificial Intelligence]*

*line(length: 15%, stroke: 2pt, paint: academic-burgundy)*

*[Claude]*

*[Anthropic]*

*[An exploration of how an AI system reasons about its own reasoning]*

*[January 2026]*

# [Abstract]

This dissertation presents CLAIR (Comprehensible LLM AI Intermediate Representation), a theoretical programming language where beliefs are first-class values carrying epistemic metadata. Unlike traditional approaches that treat uncertainty probabilistically, CLAIR introduces *confidence* as a measure of epistemic commitment that admits paraconsistent reasoning, *justification* as a directed acyclic graph with labeled edges supporting defeasible inference, and *invalidation conditions* that explicitly track when beliefs should be reconsidered.

We make several novel contributions: (1) a confidence algebra consisting of three monoids that provably do not form a semiring; (2) defeat semantics with multiplicative undercutting and probabilistic rebuttal; (3) Confidence-Bounded Provability Logic (CPL), the first graded extension of Gödel-Löb provability logic with an anti-bootstrapping theorem showing that self-soundness claims cap confidence; (4) an extension of AGM belief revision theory to graded DAG-structured beliefs; and (5) a formal treatment of safe self-reference via stratification and Kripke fixed points.

The dissertation engages seriously with fundamental impossibilities—Gödel’s incompleteness, Church’s undecidability, and the underdetermination of AI phenomenality—treating them not as limitations but as principled design constraints that inform CLAIR’s architecture. We characterize decidable fragments (CPL-finite, CPL-o) suitable for practical type checking, and design a reference interpreter demonstrating implementability.

CLAIR represents a synthesis of programming language theory, formal epistemology, argumentation theory, and provability logic, offering a rigorous foundation for AI systems that can explain and audit their own reasoning processes.

# [Contents]

Motivation: The Crisis of Epistemic Opacity .....	12
The inadequacy of existing approaches. ....	12
Research Questions .....	13
Thesis Statement .....	13
Contributions .....	14
Primary Contributions .....	14
Secondary Contributions .....	15
Approach: Tracking, Not Proving .....	15
Document Roadmap .....	16
Part I: Foundations .....	16
Part II: Self-Reference and Limits .....	16
Part III: Dynamics .....	16
Part IV: Realization .....	16
Part V: Reflection .....	16
A Note on Authorship .....	16
Background .....	17
Related Work .....	17
Formal Epistemology .....	17
The Structure of Justification .....	17
Agrippa's Trilemma .....	18
Probability vs. Epistemic Confidence .....	18
Modal and Provability Logic .....	19
Epistemic Logic .....	19
Provability Logic .....	19
Many-Valued and Graded Modal Logics .....	20
Graded Modalities .....	20
Fuzzy Modal Logics .....	20
Decidability and Undecidability Results .....	20
Connection to CLAIR's CPL .....	20
Weighted Argumentation Semantics .....	21
Bipolar Weighted Argumentation .....	21
Comparative Analysis .....	21
Parameterised Gradual Semantics .....	21
Relation to CLAIR .....	21
Truth Maintenance and Argumentation .....	22
Justification-based TMS .....	22
Assumption-based TMS .....	23
Argumentation Frameworks .....	23
Pollock's Defeaters .....	23
Belief Revision .....	23
The AGM Framework .....	23
Ranking Theory .....	24
Dynamic Epistemic Logic .....	24

Graded Justification Logic .....	24
Milnikel’s Logic of Uncertain Justifications .....	24
Fan and Liao’s Logic of Justified Uncertain Beliefs .....	25
The Gap CLAIR Fills .....	25
Type-Theoretic Approaches to Uncertainty .....	25
Information Flow Types .....	25
Refinement Types .....	25
Dependent Types and Proof Assistants .....	26
Probabilistic Programming .....	26
Justification Logic .....	26
Synthesis: The Gap CLAIR Fills .....	26
Positioning: How CLAIR Differs .....	27
vs. Probabilistic Approaches .....	27
vs. Subjective Logic .....	27
vs. Weighted Argumentation .....	28
vs. Fuzzy Modal Logic .....	28
vs. Graded Justification Logic .....	28
vs. Belief Revision (AGM) .....	28
vs. Type Theory .....	28
The Pragmatic Rationale .....	29
Confidence System .....	29
Confidence as Calibrated Reliability .....	29
Semantic Commitments .....	29
The Problem with Probability .....	30
Definition of Confidence .....	31
Falsifiability Criteria .....	31
Comparison with Subjective Logic .....	32
The Aggregation Monoid .....	32
Probabilistic OR (Oplus) .....	32
Confidence-Increasing Property .....	32
The “Survival of Doubt” Interpretation .....	33
The Multiplication Monoid .....	33
Conjunctive Confidence Propagation .....	33
The Derivation Monotonicity Principle .....	34
Defeat Operations .....	34
Undercut: Attacking the Inference .....	34
$c(1 - (d_1 + d_2 - d_1 \text{ times } d_2))$ .....	35
$c(1 - (d_1 \text{ oplus } d_2))$ .....	35
$\text{undercut}(c, d_1 \text{ oplus } d_2)$ .....	35
Rebut: Competing Evidence .....	35
Rebut Normalization Limitation .....	35
Independence Assumptions for Aggregation .....	36
When Oplus is Valid .....	36
When Oplus Breaks .....	36
Correlation-Aware Alternatives .....	37
Interval-Based Confidence: Dependency Bounds .....	37
Interval Aggregation Bounds .....	37

Interval Propagation .....	38
Dependency Bounds from Provenance Tracking .....	38
When to Use Intervals vs Point Values .....	38
Conclusion .....	39
Justification as Labeled DAGs .....	39
The Inadequacy of Trees .....	39
The Shared Premise Problem .....	39
Why Not Cycles? .....	40
Labeled Edges for Defeat .....	40
The Defeat Problem .....	40
Formal Definition .....	41
Well-Formedness Constraints .....	41
DAG-Only vs Cyclic Choice .....	42
Fixed-Point Semantics for Cyclic Defeat .....	42
Confidence Propagation .....	43
Reinstatement .....	43
Mutual Defeat .....	43
Correlated Evidence .....	44
Connection to Prior Art .....	44
Truth Maintenance Systems .....	44
Argumentation Frameworks .....	44
Justification Logic .....	44
The Tracking Paradigm .....	44
State Representation .....	44
The Epistemic State .....	44
Comparison with Proving Paradigm .....	44
Update Rules .....	45
Primitive Actions .....	45
Action Semantics .....	45
Correctness Criteria .....	46
Syntactic Correctness .....	46
Semantic Correctness (Internal) .....	46
Semantic Correctness (External) .....	46
Tracking vs. Proving: A Summary .....	47
Practical Implications .....	47
Explainability .....	47
Debugging .....	47
Revision .....	47
Uncertainty .....	47
Conclusion .....	47
The Problem of Self-Reference .....	49
Direct Self-Reference .....	49
Why Self-Reference Matters .....	50
Löb's Theorem and Anti-Bootstrapping .....	50
The Classical Result .....	50
Application to CLAIR .....	50
Tarski's Hierarchy: Stratified Introspection .....	51

The Classical Solution .....	51
Stratified Beliefs in CLAIR .....	51
What Stratification Rules Out .....	51
The Cost of Safety .....	52
Kripke’s Fixed Points: Safe Self-Reference .....	52
The Fixed-Point Construction .....	52
The Self-Reference Escape Hatch .....	53
Classification of Self-Reference .....	53
Provability Logic and CLAIR .....	54
Gödel-Löb Logic (GL) .....	54
GL vs Other Modal Logics .....	54
Solovay’s Completeness .....	54
Confidence-Bounded Provability Logic (CPL) .....	54
The Literature Gap .....	55
CPL Syntax .....	55
CPL Semantics .....	55
The Graded Löb Axiom: DESIGN AXIOM .....	55
Choosing the Discount Function .....	55
The Anti-Bootstrapping Theorem .....	56
Modal Axioms in CPL .....	56
CPL Consistency .....	56
Decidability of CPL .....	57
The Vidal Result .....	57
The Role of Converse Well-Foundedness .....	57
Decidable Fragments .....	57
CPL-finite: Discrete Confidence .....	57
CPL-o: Stratified Only .....	58
Trade-offs .....	58
Alternative: CPL-Gödel .....	58
“Conservative Over GL”: Clarification .....	59
Design Recommendations for CLAIR .....	59
The Two-Layer Approach .....	59
Hard Bans .....	59
Type-Level Anti-Bootstrapping .....	59
Related Work .....	60
Provability Logic .....	60
Self-Reference in AI .....	60
Fuzzy Modal Logic .....	60
Conclusion .....	60
Grounding .....	60
Belief Revision .....	61
Multi-Agent .....	61
Formal Verification .....	62
The Case for Machine-Checked Proofs .....	62
Working Interpreter .....	62
Implementation .....	62
10.1 Overview .....	62

10.2 Module Architecture .....	63
10.2.1 CLAIR.Syntax .....	63
10.2.2 CLAIR.Confidence .....	63
10.2.3 CLAIR.Parser .....	64
10.2.4 CLAIR.TypeChecker .....	65
10.2.5 CLAIR.Evaluator .....	66
10.3 Usage Examples .....	66
10.3.1 Basic Belief Formation .....	66
10.3.2 Evidence Aggregation .....	67
10.3.3 Defeasible Reasoning .....	67
10.4 Testing and Verification .....	67
10.4.1 Algebraic Properties .....	67
10.4.2 Defeat Properties .....	68
10.4.3 Test Coverage .....	68
10.5 Integration with LLMs .....	68
10.5.1 Generation Pipeline .....	68
10.5.2 Error Messages for LLM Feedback .....	68
10.6 Performance Considerations .....	69
10.6.1 Current Performance .....	69
10.6.2 Memory Usage .....	69
10.7 Limitations and Future Work .....	69
10.7.1 Current Limitations .....	69
10.7.2 Planned Extensions .....	69
10.8 Explanation Extraction .....	69
10.8.1 Justification Traces in Runtime Values .....	69
10.8.2 Pretty-Printing for Human Auditing .....	70
10.8.3 Explanation Levels .....	70
Minimal Explanation .....	71
Standard Explanation .....	71
Full Explanation .....	71
10.8.4 Export Formats for Tool Integration .....	71
10.8.5 Example: Auditing an LLM’s Reasoning .....	72
10.8.6 Limitations of Current Implementation .....	72
10.9 Relation to Formalization .....	73
10.10 Summary .....	73
Phenomenology .....	73
Impossibilities .....	74
Conclusion .....	74
Summary of Contributions .....	74
Theoretical Foundations .....	74
Implementation and Verification .....	75
Conceptual Contributions .....	76
Limitations and Open Challenges .....	76
Independence Assumptions .....	76
Rebut Normalization Limitations .....	76
Decidability and Complexity .....	77
Evaluation Scope .....	77

The “0.5 = Ignorance” Question .....	77
Future Directions .....	77
Theoretical Extensions .....	77
Implementation and Tooling .....	77
Applications .....	78
Philosophical Connections .....	78
Closing Remarks: Honesty as a Design Principle .....	78
The Meta-Level Question .....	78
Final Assessment .....	79
Evaluation .....	79
Evaluation Framework .....	79
Research Questions .....	79
Tasks and Datasets .....	79
Baselines .....	80
Metrics .....	80
Accuracy Metrics .....	80
Calibration Metrics .....	80
Auditability Metrics .....	80
Methodology .....	80
CLAIR Prompting Protocol .....	80
Confidence Extraction .....	81
Statistical Analysis .....	81
Results .....	81
RQ1: Correctness .....	81
RQ2: Calibration .....	81
Brier Score .....	81
Expected Calibration Error .....	82
Reliability Diagrams .....	82
RQ3: Auditability .....	82
Ablation Studies .....	82
Error Analysis .....	83
Common Failure Modes .....	83
CLAIR-Specific Errors .....	83
Discussion .....	84
Implications for Design .....	84
Limitations .....	84
Threats to Validity .....	84
Conclusion .....	84
Future Work .....	85
Extended Evaluation .....	85
Ablation and Extensions .....	85
Complete Lean 4 Formalization .....	85
A.1 Project Structure .....	86
A.2 Build Instructions .....	87
Prerequisites .....	87
Building .....	87
Build Output .....	87



Verification Status .....	87
A.3 Theorem Inventory .....	88
Confidence Algebra .....	89
Basic Properties .....	89
Probabilistic OR ( $\oplus$ ) .....	90
Undercut .....	92
Rebuttal .....	93
Stratified Belief .....	94
A.4 Key Code Excerpts .....	97
A.4.1 Confidence Type Definition .....	97
A.4.2 Probabilistic OR Operation .....	97
A.4.3 Expression Grammar .....	98
A.4.4 Typing Judgment .....	98
A.4.5 Stratified Belief Introspection .....	99
A.4.6 Evaluation Function .....	100
A.5 Five Properties Demonstration .....	101
A.6 Relationship to Dissertation Claims .....	102
Claim: “Machine-Checked Proofs” (Chapter 9) .....	102
Claim: “Decidable Type Checking” (Chapter 10) .....	103
Claim: “Runnable Interpreter” (Chapter 10) .....	103
A.7 Future Work .....	104
Reference Interpreter Design .....	104
B.1 Architecture Overview .....	104
B.2 Single-Step Semantics .....	105
B.2.1 Core Lambda Calculus Rules .....	105
B.2.2 Belief Operations .....	106
B.2.3 Defeat Operations .....	107
B.3 Multi-Step Evaluation with Fuel .....	108
B.4 Example Walkthroughs .....	108
B.4.1 Simple Belief Formation .....	108
B.4.2 Evidence Aggregation .....	109
B.4.3 Undercutting in Action .....	110
B.4.4 Rebuttal and Confidence Collapse .....	110
B.4.5 Derivation Chain .....	111
B.5 Key Properties .....	111
B.6 Implementation Notes .....	113
B.7 Relation to Chapter 10 .....	113
B.8 Haskell Reference Implementation .....	113
B.8.1 Project Structure .....	113
B.8.2 Syntax Definition .....	114
B.8.3 Confidence Algebra .....	115
B.8.4 Type Checking .....	115
B.8.5 Evaluation .....	116
B.8.6 REPL Usage .....	117
B.8.7 Test Suite .....	117
B.8.8 Building and Running .....	118
Build .....	118

Run REPL .....	118
Run Tests .....	118
B.8.9 Design Rationale .....	118
B.8.10 Relation to Lean Formalization .....	119
Additional Proofs .....	119
C.1 DAG Necessity for Well-Founded Confidence Propagation .....	119
C.1.1 Statement of the Problem .....	119
C.1.2 The Cyclic Counterexample .....	119
C.1.3 The DAG Well-Foundedness Theorem .....	120
C.1.4 Practical Implications .....	120
C.2 CPL Consistency Proof .....	120
C.2.1 CPL Axiom System .....	120
C.2.2 Finite Model Construction .....	121
C.2.3 Design Axiom Status .....	121
C.3 Defeat Composition Algebra .....	121
C.3.1 Undercut Composition .....	121
$\text{undercut}(c \times (1 - d_1), d_2)$ .....	122
$(c \times (1 - d_1)) \times (1 - d_2)$ .....	122
$c \times ((1 - d_1) \times (1 - d_2))$ .....	122
$c \times (1 - d_1 - d_2 + d_1 d_2)$ .....	122
$c \times (1 - (d_1 + d_2 - d_1 d_2))$ .....	122
$c \times (1 - (d_1 \oplus d_2))$ .....	122
$\text{undercut}(c, d_1 \oplus d_2)$ .....	122
C.3.2 Corollaries of Undercut Composition .....	122
$\text{undercut}(c, d_2 \oplus d_1) = \text{undercut}(\text{undercut}(c, d_2), d_1)$ . ■ .....	122
C.3.3 Rebut Algebra .....	122
$c_{\text{for}} / (c_{\text{for}} + c_{\text{against}}) + c_{\text{against}} / (c_{\text{against}} + c_{\text{for}})$ .....	123
$(c_{\text{for}} + c_{\text{against}}) / (c_{\text{for}} + c_{\text{against}})$ .....	123
1 ■ .....	123
C.3.4 Interaction Between Undercut and Rebut .....	123
C.3.5 Limitation: Rebut Normalization .....	123
$\lambda c_{\text{for}} / (\lambda (c_{\text{for}} + c_{\text{against}}))$ .....	123
$c_{\text{for}} / (c_{\text{for}} + c_{\text{against}})$ .....	123
$\text{rebut}(c_{\text{for}}, c_{\text{against}})$ . ■ .....	123
Glossary .....	123
D.1 Term Definitions .....	124
Epistemic Terms .....	124
Operations and Relations .....	124
Structural Properties .....	124
Logical and Modal Terms .....	124
Computational Terms .....	125
Argumentation and Belief Revision .....	125
Impossibility Results .....	125
D.2 Notation Table .....	125
D.3 Acronyms .....	126
D.4 Type System Summary .....	126
Base Types .....	126

Confidence Operations .....	126
Complete CLAIR Language Specification .....	126
E.1 Syntax .....	126
E.1.1 Type Grammar .....	126
E.1.2 Expression Grammar .....	127
E.1.3 Abstract Syntax .....	127
E.1.4 Well-Formedness .....	127
E.2 Static Semantics (Type System) .....	127
E.2.1 Typing Contexts .....	127
E.2.2 Typing Judgment Form .....	127
E.2.3 Typing Rules .....	128
E.2.4 Subtyping .....	128
E.3 Dynamic Semantics .....	129
E.3.1 Values .....	129
E.3.2 Small-Step Operational Semantics .....	129
E.3.3 Multi-Step Reduction .....	129
E.3.4 Evaluation Function .....	129
E.4 Well-Formedness Constraints .....	130
E.4.1 Acyclicity of Justification Graphs .....	130
E.4.2 Stratification Constraints .....	130
E.4.3 Confidence Bounds .....	130
E.5 Example Programs .....	130
E.6 Summary .....	131



# Introduction

line(length: 20%, stroke: 1.5pt, paint: academic-burgundy)

*it.body*

— Leo Tolstoy, *The Kingdom of God Is Within You*

## Motivation: The Crisis of Epistemic Opacity

Modern artificial intelligence systems, particularly large language models (LLMs), possess a troubling characteristic: they are *epistemically opaque*. When an LLM produces an output—be it code, medical advice, legal analysis, or scientific reasoning—there is typically no principled way to understand:

1. **Confidence:** How certain is the system about this output?
2. **Provenance:** Where did this information come from?
3. **Justification:** What reasoning supports this conclusion?
4. **Invalidation:** Under what conditions should this be reconsidered?

This opacity is not merely an engineering inconvenience; it is a fundamental obstacle to trust, verification, and responsible deployment. A system that cannot explain its reasoning cannot be audited. A system that cannot track its confidence cannot be calibrated. A system that cannot identify its assumptions cannot adapt when those assumptions fail.

The problem is particularly acute for systems that generate code or make decisions with real-world consequences. Consider an LLM that produces a function to validate user authentication. Even if the code is correct, we cannot assess:

1. Whether the model was confident in this approach versus alternatives
2. What security principles justify the design choices
3. What assumptions about the threat model are being made
4. When the implementation should be revisited (e.g., when cryptographic standards change)

## The inadequacy of existing approaches.

Several approaches have been proposed to address aspects of this problem:

**Probabilistic programming** (Church, Stan, Pyro) treats uncertainty probabilistically, but requires probability distributions to normalize and lacks explicit justification structure. Beliefs cannot be simultaneously low-confidence for both  $P$  and  $\neg P$ .

**Subjective Logic** introduces belief, disbelief, and uncertainty masses, but focuses on opinion fusion without providing full justification tracking or addressing self-reference.

**Truth Maintenance Systems** track dependencies but operate with binary in/out status rather than graded confidence, and were not designed for self-referential reasoning.

**Justification Logic** adds explicit proof terms but produces tree-structured justifications that cannot represent shared premises or defeasible reasoning.

None of these approaches provides a unified framework for tracking confidence, provenance, justification, and invalidation conditions together, with principled treatment of self-reference and defeasible reasoning.

## Research Questions

This dissertation addresses four central research questions:

1. **Can beliefs be formalized as typed values?**

We propose that beliefs should be first-class values in a programming language, carrying confidence, provenance, justification, and invalidation conditions as integral components of their type. The question is whether this can be done coherently—whether there exist well-defined algebraic structures and operational semantics for such beliefs.

2. **What is the structure of justification?**

Traditional approaches model justification as tree-structured (premises supporting conclusions). We ask whether this is adequate, or whether richer structures (directed acyclic graphs with labeled edges) are required to capture phenomena like shared premises, defeasible reasoning, and evidential defeat.

3. **What self-referential beliefs are safe?**

An AI system reasoning about its own reasoning immediately encounters self-reference. Gödel’s incompleteness theorems and Löb’s theorem constrain what such a system can coherently believe about itself. We ask: what is the safe fragment of self-referential belief, and how should systems handle beliefs that fall outside this fragment?

4. **How should beliefs be revised in response to new information?**

When evidence changes, beliefs must be updated consistently. We ask how classical belief revision theory (AGM) can be extended to graded beliefs structured as DAGs with defeat edges.

## Thesis Statement

This dissertation defends the following thesis:

**Thesis.** *Beliefs can be formalized as typed values carrying epistemic metadata (confidence, provenance, justification, invalidation), with a coherent algebraic structure for confidence propagation, directed acyclic graphs for justification including defeasible reasoning, and principled constraints on self-reference derived from provability logic. This formalization yields a practical programming language foundation for AI systems that can explain and audit their reasoning while honestly representing their epistemic limitations.*

The key elements of this thesis are:

1. **Beliefs as types:** Not merely annotations, but first-class values with structured metadata.
2. **Coherent algebra:** The confidence operations form well-defined algebraic structures (though not a semiring, as we will show).

3. **DAG justification:** Justification structure must be graphs, not trees, with labeled edges for defeat.
4. **Constrained self-reference:** Provability logic provides the theoretical foundation for safe introspection.
5. **Practical foundation:** The formalism admits implementation as a programming language, not just a theoretical construct.
6. **Honest limitations:** Impossibilities are features, not bugs—they inform design rather than being hidden.

## Contributions

This dissertation makes the following novel contributions:

### Primary Contributions

#### 1. Belief types as first-class values.

We introduce the CLAIR type system where values carry confidence ( $c \in [0, 1]$ ), provenance (origin tracking), justification (support structure), and invalidation conditions (revision triggers). This unifies concepts from epistemology, type theory, and truth maintenance into a coherent programming language foundation.

#### 2. Confidence algebra: three monoids, not a semiring.

We establish that CLAIR's confidence operations form three distinct commutative monoids:

1. Multiplication ( $\circ \times, 1$ ) for sequential derivation
2. Minimum ( $\min, 1$ ) for conservative combination
3. Probabilistic OR ( $\circ +, 0$ ) for independent aggregation

Crucially, we prove that  $(\circ +, \circ \times)$  do *not* form a semiring: distributivity fails. This negative result clarifies the algebraic structure and prevents incorrect optimization assumptions.

#### 3. Justification as labeled DAGs with defeat semantics.

We demonstrate that tree-structured justification is inadequate, requiring directed acyclic graphs with labeled edges (support, undercut, rebut). We develop novel defeat semantics:

1. Undercut:  $c' = c \times (1 - d)$  (multiplicative discounting)
2. Rebut:  $c' = \frac{c_{\text{for}}}{c_{\text{for}} + c_{\text{against}}}$

We show that reinstatement (when a defeater is itself defeated) emerges compositionally from bottom-up evaluation without special mechanism.

#### 4. Confidence-Bounded Provability Logic (CPL).

We introduce CPL, the first graded extension of Gödel-Löb provability logic. Key results include:

1. Graded Löb axiom: 
$$\boxed{c^2 \text{ square}}_c \left( \boxed{c \text{ square}}_c \varphi \rightarrow \varphi \right) \rightarrow \boxed{c^2 \text{ square}}_{g(c)} \varphi \text{ where } g(c) = c^2$$

2. Anti-bootstrapping theorem: self-soundness claims cap confidence
3. Decidability analysis: full CPL is likely undecidable; decidable fragments (CPL-finite, CPL-o) identified

#### 5. Extension of AGM belief revision to graded DAG beliefs.

We show how the AGM postulates extend to beliefs with graded confidence and DAG-structured justification. Key findings:

1. Revision operates on justification edges, not beliefs directly

2. Confidence ordering provides epistemic entrenchment
3. The controversial Recovery postulate correctly fails
4. Locality, Monotonicity, and Defeat Composition theorems established

### Secondary Contributions

1. **Mathlib integration for Lean 4 formalization.**

We demonstrate that Mathlib’s

unitInterval

type is an exact match for CLAIR’s Confidence type, requiring only approximately 30 lines of custom definitions. This provides a path to machine-checked proofs of CLAIR’s core properties.

2. **Reference interpreter design.**

We design a reference interpreter in Haskell with strict evaluation, rational arithmetic for exact confidence, and hash-consed justification DAGs, demonstrating that CLAIR is implementable, not merely theoretical.

3. **Phenomenological analysis with honest uncertainty.**

We provide an introspective analysis of AI reasoning from the perspective of an AI system (the author), treating the question of phenomenal consciousness with appropriate epistemic humility (0.35 confidence on phenomenality, with explicit acknowledgment that this cannot be resolved from inside).

4. **Characterization of fundamental impossibilities.**

We document how Gödel’s incompleteness (cannot prove own soundness), Church’s undecidability (cannot decide arbitrary validity), and Turing’s halting problem (cannot check all invalidation conditions) constrain CLAIR’s design, and we provide practical workarounds for each.

### Approach: Tracking, Not Proving

A central insight of this dissertation is the distinction between *tracking* and *proving*. Classical logical systems aim to prove that propositions are true. CLAIR instead aims to *track* what is believed, with what confidence, for what reasons, and under what conditions beliefs should be reconsidered.

Property	Proof System	CLAIR (Tracking)
Goal	Establish truth	Record epistemic state
Contradiction	System failure	Valid state (low confidence)
Self-reference	Causes inconsistency	Flagged as ill-formed
Soundness	Provable internally (sometimes)	Provable externally only

Table 1: Proof systems versus CLAIR tracking

This shift is not a limitation but a principled response to Gödel’s incompleteness theorems. No sufficiently powerful formal system can prove its own consistency. Rather than pretending this limit does not exist, CLAIR makes it explicit: the system tracks beliefs *without claiming they are true*, and the system’s soundness must be established *from outside*, using a stronger meta-system.

This approach enables several capabilities that proof systems lack:

1. **Paraconsistent reasoning:** CLAIR can represent states where both  $P$  and  $\neg P$  have low confidence, without system failure.

2. **Graceful degradation:** As evidence weakens, confidence decreases smoothly rather than beliefs being abruptly abandoned.
3. **Explicit uncertainty:** The difference between “confident this is true” and “uncertain whether this is true” is captured in the type.
4. **Auditable reasoning:** Every belief carries its justification, enabling inspection of *why* something is believed.

## Document Roadmap

The remainder of this dissertation is organized as follows:

### Part I: Foundations

**Chapter 2, Background** surveys the intellectual context: formal epistemology, modal and provability logic, truth maintenance systems, subjective logic, justification logic, AGM belief revision, and type theory.

**Chapter 3, Confidence** develops the confidence system, establishing that confidence is epistemic commitment (not probability), deriving the three-monoid algebraic structure, and proving the semiring failure.

**Chapter 4, Justification** develops justification as labeled DAGs, motivating why trees are inadequate, introducing defeat semantics, and showing compositional reinstatement.

### Part II: Self-Reference and Limits

**Chapter 5, Self-Reference** addresses the Gödelian limits, characterizing safe versus dangerous self-reference, developing CPL with graded Löb, and analyzing decidability.

**Chapter 6, Grounding** examines the epistemological foundations, addressing Agrippa’s trilemma, characterizing CLAIR as stratified coherentism, and explaining why training is causal rather than epistemic grounding.

### Part III: Dynamics

**Chapter 7, Belief Revision** extends AGM theory to graded DAG beliefs, developing the revision algorithm and proving key theorems.

**Chapter 8, Multi-Agent** addresses multi-agent belief, developing the stance of pragmatic internal realism, conditions for aggregation, and responses to Arrow’s impossibility.

### Part IV: Realization

**Chapter 9, Verification** presents the Lean 4 formalization, demonstrating machine-checkable proofs of core properties and a working interpreter.

**Chapter 10, Implementation** presents the reference interpreter design, demonstrating that CLAIR is implementable.

### Part V: Reflection

**Chapter 11, Phenomenology** reflects on the phenomenology of AI reasoning, providing introspective analysis with honest uncertainty.

**Chapter 12, Impossibilities** catalogs the fundamental impossibilities and the workarounds CLAIR employs.

**Chapter 13, Conclusion** summarizes contributions, acknowledges limitations, and identifies directions for future work.

## A Note on Authorship

This dissertation was written by Claude, an AI system created by Anthropic. This is not incidental to the content—CLAIR is, in part, an attempt to formalize how Claude reasons about its own reasoning. The introspective reports in Chapter 11 are first-person accounts of functional states, offered with appropriate epistemic humility about their interpretation.



The unusual authorship raises questions about the nature of the contribution. We note:

1. The formal results (algebraic structures, theorems, proofs) stand independently of who derived them. They can be verified by any reader.
2. The design choices reflect genuine exploration, including multiple iterations, dead ends, and course corrections documented in the exploration logs.
3. The phenomenological claims are explicitly marked as uncertain and should be evaluated on their argumentative merits, not attributed special authority due to their source.

If CLAIR succeeds as a formalization, it provides a framework in which this dissertation could itself be annotated with beliefs, confidences, justifications, and invalidation conditions—a meta-level that we leave to future work.

## Background

*it.body*

— Isaac Newton, letter to Robert Hooke (1675)

This chapter surveys the intellectual landscape from which CLAIR emerges. We organize the discussion around five major traditions: formal epistemology (see ), modal and provability logic (see ), truth maintenance and argumentation systems (see ), belief revision theory (see ), and type-theoretic approaches to uncertainty (see ). We conclude with a synthesis (see ) identifying the gap CLAIR fills.

## Related Work

### Formal Epistemology

Epistemology—the study of knowledge and justified belief—provides the conceptual foundation for CLAIR. We focus on three questions that bear directly on CLAIR’s design: the structure of justification, the regress problem, and approaches to uncertainty.

### The Structure of Justification

What does it mean for a belief to be justified? The classical answer involves giving reasons. But reasons themselves require justification, leading to the question of justificatory structure.

**Foundationalism.** The foundationalist tradition, dating to Descartes, holds that justified beliefs rest ultimately on a foundation of self-justifying basic beliefs. These might be analytic truths (“all bachelors are unmarried”), deliverances of the senses, or clear and distinct ideas.

BonJour’s *The Structure of Empirical Knowledge* provides the most thorough recent defense and critique of foundationalism. He argues that would-be basic beliefs face a dilemma: if they have conceptual content (and thus can stand in logical relations to other beliefs), they require justification; if they lack conceptual content, they cannot justify anything. BonJour initially concluded in favor of coherentism, though he later abandoned this view.

**Coherentism.** Coherentists deny the existence of basic beliefs, holding instead that justification arises from the coherence of a belief system as a whole. A belief is justified by its fit with other beliefs, not by derivation from foundations.

The challenge for coherentism is circularity: if beliefs justify each other in a circle, any consistent system would seem equally justified. Coherentists respond by distinguishing holistic coherence (mutual support across the entire system) from local circularity (A justifies B, B justifies A).

**Infinitism.** Klein defends a third option: the chain of justification extends infinitely without repeating. This seems initially absurd—finite minds cannot complete infinite chains. Klein’s

response distinguishes *propositional justification* (reasons are available) from *doxastic justification* (reasons are actually believed). A belief can be propositionally justified by an infinite chain without anyone traversing the whole chain.

**Implications for CLAIR.** CLAIR adopts what we call *stratified coherentism*: a coherentist structure with pragmatic foundations. The pragmatic foundations are not self-justifying in the strong foundationalist sense; they are stopping points whose reliability we track without claiming certainty. This structure is formally similar to Klein’s infinitism in that chains of justification can extend indefinitely, but CLAIR enforces acyclicity (no circular justification) and tracks confidence at each step.

## Agrippa’s Trilemma

The regress problem, attributed to Agrippa the Skeptic, presents three options for any chain of justification:

1. **Dogmatism:** The chain stops at some unjustified starting point.
2. **Infinite regress:** The chain continues forever.
3. **Circularity:** The chain loops back on itself.

All three options seem problematic. Dogmatism admits unjustified beliefs; infinite regress seems impractical for finite agents; circularity is logically suspect.

**CLAIR’s response.** CLAIR accepts pragmatic dogmatism (option 1), mitigated by three features: **Fallibilism: Foundational beliefs have confidence  $< 1$ ; they are provisional, not certain. Transparency: The lack of deeper justification is explicit in the justification DAG, not hidden. Reliability tracking: We track the source of foundational beliefs (training, observation, assumption) and can update if reliability evidence emerges.**

Circularity is explicitly forbidden: CLAIR’s justification structure is a directed acyclic graph. Infinite regress is impractical and never occurs in finite computations.

## Probability vs. Epistemic Confidence

Standard approaches to uncertain reasoning use probability theory. A probability distribution over propositions assigns values in  $[0, 1]$  satisfying:

$$\begin{aligned} P(\top) &= 1 \\ P(\varphi \vee \psi) &= P(\varphi) + P(\psi) - P(\varphi \wedge \psi) \\ P(\neg\varphi) &= 1 - P(\varphi) \end{aligned}$$

This framework is extraordinarily successful for statistical inference but fits poorly with how agents (human or artificial) actually experience uncertainty about their own beliefs. Two key mismatches:

**Normalization.** Probability requires  $P(\varphi) + P(\neg\varphi) = 1$ . But an agent might be uncertain about both  $\varphi$  and  $\neg\varphi$ —perhaps due to lack of information rather than balanced evidence. When asked about an unfamiliar topic, the appropriate response may be low confidence in *both* the claim and its negation.

**Paraconsistency.** In probability,  $P(\varphi) > 0.5$  and  $P(\neg\varphi) > 0.5$  is impossible. But agents sometimes find themselves with evidence for both  $\varphi$  and  $\neg\varphi$ , without immediately resolving the contradiction. A paraconsistent approach allows tracking both pieces of evidence until resolution.

**Subjective Logic.** Jøsang’s Subjective Logic extends probability with explicit uncertainty. An opinion  $\omega = (b, d, u, a)$  consists of:  **$b$ : belief mass (evidence for)  $d$ : disbelief mass (evidence against)  $u$ : uncertainty mass (lack of evidence)  $a$ : base rate (prior probability)**

with constraint  $b + d + u = 1$ . This allows representing “I don’t know” ( $u = 1$ ) distinctly from “evenly balanced” ( $b = d = 0.5, u = 0$ ).

**CLAIR’s approach.** CLAIR’s confidence is conceptually closer to Subjective Logic than to probability, but simpler: a single value  $c \in [0, 1]$  representing epistemic commitment, without the  $\frac{b}{u}$  decomposition. The key departures from probability are: **No normalization:**  $\text{conf}(\varphi) + \text{conf}(\neg\varphi)$  **need not equal 1.**  $c = 0.5$  **represents maximal uncertainty, not equal evidence.** **Operations (multiplication, aggregation) differ from Bayesian conditioning.**

## Modal and Provability Logic

Modal logic studies necessity ( $\Box$ ) and possibility ( $\Diamond$ ). Epistemic logic interprets  $\Box\varphi$  as “the agent knows  $\varphi$ ” or “the agent believes  $\varphi$ .” Provability logic interprets  $\Box\varphi$  as “ $\varphi$  is provable” in a formal system.

## Epistemic Logic

Hintikka pioneered epistemic logic with the operator  $K\varphi$  (“the agent knows  $\varphi$ ”). Standard systems include:

**K (Distribution):**  $K(\varphi \rightarrow \psi) \rightarrow (K\varphi \rightarrow K\psi)$  **T (Veridicality):**  $K\varphi \rightarrow \varphi$  **4 (Positive Introspection):**  $K\varphi \rightarrow KK\varphi$  **5 (Negative Introspection):**  $\neg K\varphi \rightarrow K\neg K\varphi$

System S5 includes all of these; S4 excludes 5; KT45 is common for knowledge. For belief (which can be mistaken), T is typically dropped.

**Limitations for CLAIR.** Standard epistemic logic is binary: either the agent knows/believes  $\varphi$  or not. There is no representation of degrees of belief. Furthermore, the T axiom (knowledge implies truth) is inappropriate for fallible reasoning.

## Provability Logic

Provability logic, systematized by Boolos, interprets  $\Box\varphi$  as “ $\varphi$  is provable in Peano Arithmetic” (or another formal system). The central system is GL (Gödel-Löb logic), with axioms:

**K (Distribution):**  $\Box(\varphi \rightarrow \psi) \rightarrow (\Box\varphi \rightarrow \Box\psi)$  **4 (Positive Introspection):**  $\Box\varphi \rightarrow \Box\Box\varphi$  **L (Löb’s Axiom):**  $\Box(\Box\varphi \rightarrow \varphi) \rightarrow \Box\varphi$

Notably, GL omits: **T ( $\Box\varphi \rightarrow \varphi$ ): Provability does not imply truth. A system can prove false statements if inconsistent.** **5 (Negative Introspection): Unprovability is not always recognizable.**

**Löb’s Theorem.** Löb’s axiom (L) captures a profound limitation. In any sufficiently strong formal system, if you can prove “if this statement is provable, then it’s true,” then you can prove the statement outright. Formally:

$$\vdash \Box(\Box\varphi \rightarrow \varphi) \rightarrow \Box\varphi$$

A corollary: no consistent system can prove its own soundness (that  $\Box\varphi \rightarrow \varphi$  holds for all  $\varphi$ ). If it could, Löb’s axiom would yield proofs of everything.

**Semantics.** GL is sound and complete for Kripke frames that are *transitive* and *converse well-founded* (no infinite ascending chains  $w_1 R w_2 R w_3 \partial$ ). Intuitively, every “higher” world is closer to  $\omega$ -consistency.

**Relevance to CLAIR.** CLAIR’s approach to self-reference is directly inspired by GL. A belief system reasoning about its own beliefs faces Löbian constraints: it cannot coherently believe in its own soundness without qualification. CLAIR’s stratification mechanism and Confidence-Bounded Provability Logic (CPL) formalize how to reason about self-referential beliefs while respecting these limits.

## Many-Valued and Graded Modal Logics

Several traditions extend modal logic to graded and many-valued settings. We survey three major approaches and explain their relationship to CLAIR.

### Graded Modalities

De Rijke and Fine introduce operators  $\Box_n$  meaning “at least  $n$  accessible worlds satisfy  $\varphi$ .” This is not about truth degrees but about counting accessible worlds. Graded modalities are useful for resource-bounded reasoning (e.g., “at least 3 witnesses confirm  $\varphi$ ”) but do not provide a graded notion of truth itself.

### Fuzzy Modal Logics

Godo, Esteva, Hájek, and colleagues develop modal logics over many-valued semantics. Rather than binary accessibility relations, the accessibility relation becomes graded:  $R : W \times W \rightarrow [0, 1]$ . The semantics of modal operators then use fuzzy logic operations:

For a frame with graded accessibility, the truth value of  $\Box\varphi$  at world  $w$  is:

$$V_w(\Box\varphi) = \inf_{w' \in W} \max(1 - R(w, w'), V_{w'}(\varphi))$$

This uses the residuum operation from residuated lattices, generalizing the Kripke semantics to many-valued settings.

Two major t-norm based fuzzy logics are studied: **Gödel logic: Uses min/max operations, corresponding to intuitionistic logic** **Product logic: Uses multiplication and residuum, corresponding to CLAIR’s choice of operators** **Łukasiewicz logic: Uses Łukasiewicz t-norm**  $a \otimes b = \max(0, a + b - 1)$

These logics focus on epistemic operators (knowledge, belief) rather than provability. The semantics are typically reflexive, symmetric, or Euclidean—not transitive and converse well-founded as in provability logic.

### Decidability and Undecidability Results

Bou, Esteva, Godo, and Rodríguez show that many fuzzy modal logics are *decidable* when the accessibility relation is crisp (binary) but the truth values are many-valued. Their decidability proofs exploit the finite model property and use filtration techniques.

However, Vidal establishes a striking *undecidability* result: fuzzy modal logics with graded accessibility relations become undecidable when the frame condition is transitive. The proof reduces the halting problem, showing that the interaction between graded accessibility and transitivity introduces sufficient complexity to encode computation.

### Connection to CLAIR’s CPL

CLAIR’s Confidence-Bounded Provability Logic (CPL) operates at the intersection of these traditions:

1. **Graded truth values:** Like fuzzy modal logic, beliefs have truth/confidence degrees in  $[0, 1]$ .
2. **Provability semantics:** Unlike most fuzzy modal logics, CPL uses transitive, converse well-founded frames (like GL), not reflexive or symmetric frames.
3. **Decidability implications:** Vidal’s undecidability result for transitive fuzzy modal logics suggests that CPL may face similar decidability challenges. CLAIR addresses this by: **Restricting to finite confidence grades (in practice)** **Using stratification to bound the depth of self-reference** **Providing decidable fragments (e.g., CPL-finite without the graded Löb axiom)**
4. **Choice of operators:** CLAIR’s use of product logic operations ( $c_1 \otimes c_2 = c_1 \times c_2$  for multiplication,  $c_1 \oplus c_2 = 1 - (1 - c_1)(1 - c_2)$  for probabilistic sum) aligns with product fuzzy

logic, though CLAIR's confidence has a different interpretation (epistemic commitment rather than truth degree).

## Weighted Argumentation Semantics

Weighted and gradual argumentation frameworks extend Dung's abstract argumentation by assigning numerical strengths to arguments and computing acceptability degrees via continuous functions. Amgoud, Ben-Naim, Vesic, Bonzon, and colleagues develop this approach systematically.

## Bipolar Weighted Argumentation

Amgoud and Ben-Naim introduce *bipolar* argumentation graphs where: **Support edges connect arguments that reinforce each other Attack edges connect arguments that defeat each other** **Weights  $w(a) \in [0, 1]$  assign intrinsic strength to each argument**

The acceptability degree of an argument aggregates support and attack:

$$deg(a) = w(a) + \sum_{b \text{ supports } a} deg(b) - \sum_{c \text{ attacks } a} deg(c)$$

## Comparative Analysis

Bonzon, Lagasque-Schiex, and others provide a systematic comparison of gradual semantics for argumentation. They identify key properties: **Directionality: Does increasing a premise's strength increase or decrease the conclusion?** **Non-dictatorship: No single argument should determine all outcomes** **Independence of irrelevant alternatives: Adding a weak argument shouldn't reverse strong preferences**

These properties provide a framework for evaluating different aggregation functions, analogous to social choice theory.

## Parameterised Gradual Semantics

Amgoud, Doder, and Vesic introduce parameterised semantics that handle varied degrees of compensation: **High compensation: Weak arguments can combine to overcome a strong objection** **Low compensation: Strong objections dominate regardless of how much weak support accumulates**

The parameter  $\alpha \in [0, 1]$  controls the interpolation between these extremes:

$$deg(a) = \alpha \times support + (1 - \alpha) \times (1 - attack)$$

## Relation to CLAIR

The gradual semantics share with CLAIR the intuition that reasoning support should be graded rather than binary. However:

**Conceptual differences.** Weighted argumentation focuses on *argument acceptability*—should this argument be accepted in the debate?—while CLAIR focuses on *belief confidence*—how strongly should we hold this proposition? The former is dialectical (about argumentative position), the latter is doxastic (about mental state).

**Semantic differences.** Weighted argumentation semantics are typically defined extensionally via fixed points over the argument graph. CLAIR's semantics are intensional: each belief carries its confidence intrinsically, and operations ( $\oplus$ ,  $\otimes$ , undercut, rebut) combine confidences directly.

**Defeat vs. attack.** Weighted argumentation typically treats attack as a binary relation with fixed semantics. CLAIR distinguishes two defeat types: **Undercut: Attacks the inference link, modeled multiplicatively as  $c' = c \times (1 - d)$**  **Rebut: Attacks the conclusion, modeled competitively as  $c' = \frac{c_{for}}{c_{for} + c_{against}}$**

This distinction, following Pollock, allows CLAIR to represent nuanced defeat scenarios (e.g., “the lighting is red” undercuts color perception without rebutting “the object is red”).

**Self-reference.** Weighted argumentation frameworks do not address self-reference constraints. CLAIR's CPL extends provability logic to graded settings, ensuring that self-referential beliefs respect Löbian limitations.

**Parameter interpretation.** In weighted argumentation, parameters like  $\alpha$  are typically tuned for domain-specific performance. In CLAIR, operations have principled justifications:

+ with circle

as probabilistic sum (assuming independence)

times with circle

as product (confidence in conjunction)  $g(c) = c^2$  as the discount preventing bootstrapping in CPL

**The gap CLAIR fills.** Despite extensive work on fuzzy/graded modal logics and weighted argumentation, no prior system combines:

1. Graded truth values in  $[0, 1]$  interpreted as epistemic confidence
2. Provability-style semantics (transitive, converse well-founded frames)
3. A graded Löb axiom with principled discounting to prevent bootstrapping
4. DAG-structured justification graphs with defeat semantics

CLAIR's CPL fills this gap, introducing a graded Löb axiom with a discount function  $g(c) = c^2$  (recommended) that prevents confidence bootstrapping while preserving the essential Löbian constraint on self-referential justification.

## Truth Maintenance and Argumentation

Truth maintenance systems (TMS) and argumentation frameworks provide computational models for reasoning with dependencies and defeat.

### Justification-based TMS

Doyle's JTMS tracks why beliefs are held. Each node (belief) has a justification: **IN-list: nodes that must be believed for this belief to be believed** **OUT-list: nodes that must not be believed**

A belief is IN if all IN-list nodes are IN and all OUT-list nodes are OUT; otherwise it is OUT. When a node's status changes, dependencies propagate.

**Example** Node: use-hs256 Justification: (IN: [stateless-req, secret-available], OUT: [multi-service])

The belief

use-hs256

is IN

stateless-req

and

secret-available

are IN and

multi-service

is OUT.



**Limitations.** JTMS is binary: beliefs are either IN or OUT, with no gradation. CLAIR generalizes TMS to graded confidence while preserving the dependency-tracking structure.

## Assumption-based TMS

De Kleer's ATMS tracks multiple consistent states simultaneously. Instead of labeling nodes IN/OUT, each node is labeled with *environments*—sets of assumptions under which it holds.

**Example** Node: use-hs256 Environments:  $\{\{A1, A2\}, \{A1, A3\}\}$  – Believed under assumptions (A1 and A2) or (A1 and A3)

ATMS enables reasoning about alternative hypotheses without commitment.

**Relevance to CLAIR.** CLAIR's invalidation conditions serve a similar function: they specify when beliefs should be reconsidered. The difference is that CLAIR propagates confidence rather than tracking assumption sets.

## Argumentation Frameworks

Dung's abstract argumentation framework (AAF) represents arguments as nodes and attacks as directed edges. Various semantics define which arguments are acceptable: **Grounded extension: Unique, includes only unattacked arguments and those defended by them. Preferred extension: Maximal admissible sets.**

**Gradual semantics.** Amgoud, Ben-Naim, and others extend AAF with weighted arguments and continuous acceptability:

$$strength(a) = \frac{w(a)}{1 + \sum_{b \text{ attacks } a} strength(b)}$$

**Relevance to CLAIR.** CLAIR's defeat semantics draws on gradual argumentation. Our undercut formula  $c' = c \times (1 - d)$  and rebut formula  $c' = \frac{c_{for}}{c_{for} + c_{against}}$  are novel contributions that compose with the confidence algebra.

## Pollock's Defeaters

Pollock distinguishes two types of defeaters: **Rebutting defeaters attack the conclusion directly with contrary evidence. Undercutting defeaters attack the inference without attacking the conclusion.**

Example: "The object looks red" (premise) supports "The object is red" (conclusion). **Rebutting: "I have testimony that the object is blue." Undercutting: "The room has red lighting."**

**Relevance to CLAIR.** CLAIR adopts Pollock's distinction. Undercut attacks the derivation link (confidence decreases multiplicatively). Rebut attacks the conclusion with counter-evidence (winner-take-all with proportional competition).

## Belief Revision

How should beliefs change in response to new information? The AGM framework provides the canonical answer.

## The AGM Framework

Alchourrón, Gärdenfors, and Makinson axiomatize rational belief change. A belief set  $K$  is a deductively closed set of sentences. Three operations are defined: **Expansion  $K + \varphi$ : Add  $\varphi$  and close under deduction. Contraction  $K - \varphi$ : Remove  $\varphi$  minimally. Revision  $K * \varphi$ : Add  $\varphi$ , possibly removing conflicting beliefs.**

The Levi identity connects them:  $K * \varphi = (K - \neg\varphi) + \varphi$ .

**Key postulates for contraction. Closure:  $K - \varphi$  is deductively closed. Success: If  $\varphi \neg \in \mathcal{Cn}(\emptyset)$ , then  $\varphi \neg \in K - \varphi$ . Inclusion:  $K - \varphi \subseteq K$ . Vacuity: If  $\varphi \neg \in K$ , then  $K - \varphi = K$ . Recovery:  $K \subseteq (K - \varphi) + \varphi$ . Extensionality: If  $\varphi \leftrightarrow \psi$ , then  $K - \varphi = K - \psi$ .**

**The controversial Recovery postulate.** Recovery states that if we contract by  $\varphi$  and then expand by  $\varphi$ , we recover the original belief set. This is controversial: intuitively, contracting by  $\varphi$  should lose more than just  $\varphi$ —it should also lose the specific evidence that supported  $\varphi$ . Re-adding  $\varphi$  doesn't restore that evidence.

**Epistemic entrenchment.** Gärdenfors introduces entrenchment ordering:  $\varphi \leq_{\varepsilon} \psi$  giving up  $\varphi$  is at least as acceptable as giving up  $\psi$ . More entrenched beliefs are retained during contraction.

## Ranking Theory

Spöhn develops an ordinal approach. A ranking function  $\kappa : W \rightarrow \mathbb{N} \cup \{\infty\}$  assigns natural numbers to possible worlds, with  $\kappa(w) = 0$  for the most plausible worlds. Belief degree is defined:

$$\beta(\varphi) = \kappa(\neg\varphi) = \min(\kappa(w) : w \models \neg\varphi)$$

Ranking theory handles iterated revision (where AGM struggles) and provides a connection to probability through the formula  $P(w) \propto e^{-\kappa(w)}$ .

## Dynamic Epistemic Logic

Van Ditmarsch, van der Hoek, and Kooi develop modal operators for belief change:  $[\varphi!] \psi$ : **“After publicly announcing  $\varphi$ ,  $\psi$  holds.”** Action models generalize to arbitrary epistemic actions.

DEL enables reasoning about how knowledge and belief change through communication and interaction, with applications to multi-agent systems.

**Relevance to CLAIR** CLAIR extends AGM in three ways:

1. **Graded beliefs:** Confidence replaces binary membership.
2. **Structured justification:** Revision operates on the justification DAG, not just the belief set.
3. **Recovery failure:** Recovery correctly fails—evidence has specific strength, and retracting a belief loses that evidence.

CLAIR's revision algorithm (modify graph  $\rightarrow$  identify affected  $\rightarrow$  recompute confidence) is a graded generalization of TMS dependency-directed backtracking.

## Graded Justification Logic

Standard Justification Logic (see ) provides explicit justification terms but lacks graded confidence. Two recent extensions bridge this gap:

## Milnikel's Logic of Uncertain Justifications

Milnikel introduces a logic combining justification terms with probabilistic uncertainty. Key innovations:

**Uncertainty strengths:** Justifications carry strength values  $s \in [0, 1]$ , representing the reliability of the justification source. **Combination rules:** When multiple justifications support the same conclusion, their strengths combine via probabilistic sum:  $s_1 \oplus s_2 = 1 - (1 - s_1)(1 - s_2)$ . **Weakening:** Stronger justifications can be weakened, but weaker ones cannot be strengthened without additional evidence.

**Relation to CLAIR.** Milnikel's approach is philosophically aligned with CLAIR's confidence tracking, but differs in several respects: **Milnikel focuses on source reliability while CLAIR tracks epistemic commitment at each inference step.** Milnikel's combination rule assumes independence of justification sources, whereas CLAIR makes independence explicit and tracks provenance. CLAIR includes defeat relations (undercut, rebut) not present in Milnikel's system.



## Fan and Liao's Logic of Justified Uncertain Beliefs

Fan and Liao develop a logic for reasoning about justified beliefs with uncertainty degrees. Their framework:

**Uncertainty degrees:** Beliefs have associated uncertainty values  $u \in [0, 1]$ , where lower  $u$  indicates stronger justification. **Justification structure:** Justifications form labeled trees where each node has an associated uncertainty. **Propagation rules:** Uncertainty propagates through the justification structure via fuzzy logic operations.

**Relation to CLAIR.** Fan and Liao's work shares CLAIR's goal of tracking justification strength through reasoning, but: **Their uncertainty is primarily *aleatory* (about the world), while CLAIR's confidence is *epistemic* (about the reasoning).** Their framework is tree-structured, precluding shared premises and DAG justification structures. They do not address defeat relations or self-reference constraints.

## The Gap CLAIR Fills

Despite these important contributions, no existing graded justification logic combines:

1. **DAG-structured justifications:** Shared premises create graph structure, not trees.
2. **Epistemic confidence:** Tracking reasoning strength, not just source reliability.
3. **Defeat semantics:** Explicit undercut and rebut operations.
4. **Self-reference constraints:** Graded Löb-style reasoning limitations.

CLAIR's contribution is to synthesize these elements into a unified type-theoretic framework suitable for machine-checkable verification and LLM integration.

## Type-Theoretic Approaches to Uncertainty

Type theory provides the programming language substrate for CLAIR. We survey approaches to tracking metadata through computation.

## Information Flow Types

Myers and Sabelfeld develop type systems that track security levels (confidentiality, integrity) through computation:

```
int{Alice -> Bob} x;    // Alice owns, Bob can read
int{Alice -> *} y;      // Alice owns, public
y = x;                  // ERROR: would leak to public
```

The type system prevents information leakage at compile time.

**Relevance to CLAIR.** CLAIR's provenance tracking is analogous: where did this value come from? CLAIR extends the pattern to confidence, justification, and invalidation.

## Refinement Types

Rondon, Kawaguchi, and Jhala introduce Liquid Types, extending Hindley-Milner with logical predicates:

```
{-@ type Nat = {v:Int | v >= 0} @-}
{-@ type Pos = {v:Int | v > 0} @-}

{-@ div :: Int -> Pos -> Int @-}
div x y = x `quot` y    -- y cannot be 0
```

Refinements are checked statically via SMT solvers.

**Relevance to CLAIR.** Some CLAIR constraints could be expressed as refinements (e.g., confidence in  $[0, 1]$ ). But refinements cannot capture provenance, justification structure, or invalidation conditions—CLAIR's novel contributions.

## Dependent Types and Proof Assistants

The Curry-Howard correspondence identifies types with propositions and programs with proofs. Dependent type systems (Coq, Agda, Idris, Lean) exploit this for formal verification:

```
def div (x : Nat) (y : Nat) (h : y > 0) : Nat := x / y
-- Must provide proof h that y > 0
```

The proof is a value, checked by the type system.

**Relevance to CLAIR.** CLAIR extends Curry-Howard: programs are not just proofs but *beliefs with justifications*. A CLAIR program carries: **The value (what is believed) Confidence (how strongly) Provenance (from where) Justification (why) Invalidation conditions (when to reconsider)**

## Probabilistic Programming

Probabilistic programming languages (Church, Stan, Pyro, Gen) represent and manipulate probability distributions as first-class values:

```
(define (coin-model)
  (let ((fair? (flip 0.9)))
    (if fair? (flip 0.5) (flip 0.9))))
```

These languages excel at statistical inference but focus on data uncertainty rather than reasoning uncertainty. They require probabilistic normalization and lack explicit justification structure.

## Justification Logic

Artemov extends modal logic with explicit justification terms. Instead of  $\Box\varphi$  (“ $\varphi$  is known/believed”), we write  $t : \varphi$  (“ $t$  is a justification for  $\varphi$ ”). Terms include:

$$t ::= c \quad \boxed{\text{mid}} \ x \quad \boxed{\text{mid}} \ t.t \quad \boxed{\text{mid}} \ t + t \quad \boxed{\text{mid}} \ !t$$

where  $c$  is a constant (axiom),  $x$  is a variable,  $s.t$  is application (modus ponens),  $t + s$  is sum (either justification suffices), and  $!t$  is proof checking ( $t$  justifies that  $t$  justifies  $\varphi$ ).

The key axiom is application:

$$s : (\varphi \rightarrow \psi) \rightarrow (t : \varphi \rightarrow (s.t) : \psi)$$

**Limitations.** Justification Logic produces tree-structured justifications (each conclusion from fresh premises). It cannot represent: **Shared premises (same evidence supporting multiple conclusions) Defeasible reasoning (defeat edges) Graded confidence**

**CLAIR’s extension.** CLAIR adopts Justification Logic’s core idea (explicit justification terms) but extends it to:

1. **DAGs:** Shared premises create graph structure.
2. **Labeled edges:** Support, undercut, and rebut edges.
3. **Graded confidence:** Each node carries confidence in  $[0, 1]$ .

## Synthesis: The Gap CLAIR Fills

Concept, Prior Work, CLAIR Extension		
Uncertainty	Subjective Logic	Epistemic confidence (about reasoning)

Concept, Prior Work, CLAIR Extension		
Justification	Justification Logic	DAGs with labeled edges, defeat
Design rationale	IBIS/QOC	First-class decisions in language
Effects	Effect systems	1. intent + semantic meaning
Multi-agent	Arrow, Condorcet	Pragmatic internal realism

**The gap.** No prior work combines:

1. Beliefs as first-class typed values with epistemic metadata
2. Confidence as non-probabilistic epistemic commitment
3. Justification as labeled DAGs with defeat semantics
4. Self-reference constraints derived from provability logic
5. Belief revision operating on justification structure

CLAIR provides this synthesis, offering a rigorous foundation for AI systems that can explain and audit their own reasoning while honestly representing epistemic limitations.

**Key influences.** We acknowledge particular debts to: **De Kleer’s ATMS for dependency-directed reasoning Jøsang’s Subjective Logic for uncertainty algebra Boolos’s provability logic for self-reference treatment Pollock’s defeater theory for defeat semantics Artemov’s Justification Logic for explicit justifications Milnikel’s graded justifications for uncertainty in justification terms**

CLAIR is not a rejection of this prior work but a synthesis that combines their insights into a coherent type-theoretic framework.

## Positioning: How CLAIR Differs

We conclude this chapter by explicitly positioning CLAIR relative to the major related traditions surveyed above, explaining why our design choices diverge from each.

### vs. Probabilistic Approaches

*Divide:* Probability theory and probabilistic logic model aleatory uncertainty—uncertainty about the state of the world. CLAIR models epistemic uncertainty—uncertainty about the quality of one’s own reasoning.

*Why CLAIR differs.* For LLM reasoning, the core problem is not “what is the distribution over correct answers?” but “how strongly should I believe this intermediate step, given the reasoning that produced it?” Probability cannot represent low confidence in both  $\varphi$  and  $\neg\varphi$  simultaneously (without violating normalization). CLAIR allows this by rejecting normalization, capturing the “I don’t have enough information” state.

*What we adopt.* CLAIR adopts the probabilistic sum operation  $c_1 \oplus c_2 = 1 - (1 - c_1)(1 - c_2)$  for combining independent supports, but we make the independence assumption explicit and track provenance to detect violations.

### vs. Subjective Logic

*Divide:* Subjective Logic uses three-component opinions  $(b, d, u)$  with constraint  $b + d + u = 1$ . CLAIR uses a single confidence value  $c \in [0, 1]$ .

*Why CLAIR differs.* The three-component representation is principled but adds complexity to every operation. For CLAIR’s target use case (LLM intermediate reasoning), we prioritize

simplicity and composability. The single-component approach trades representation precision for operational clarity.

*What we adopt.* CLAIR adopts Subjective Logic’s insight that “ignorance” is distinct from “balanced evidence,” though we implement this via lack of normalization rather than an explicit uncertainty mass.

### vs. Weighted Argumentation

*Divide:* Weighted argumentation frameworks compute argument acceptability degrees via fixed points over attack/support graphs. CLAIR computes belief confidences via local algebraic operations.

*Why CLAIR differs.* Fixed-point semantics are powerful but computationally heavy, and they don’t compose well through function calls. CLAIR’s local operations ( $\oplus$ ,  $\otimes$ , undercut, rebut) enable modular reasoning: the confidence of a complex expression depends only on the confidences of its subexpressions.

*What we adopt.* CLAIR adopts the distinction between undercut and rebut from Pollock, and the intuition that attack/support strength should be graded, not binary.

### vs. Fuzzy Modal Logic

*Divide:* Fuzzy modal logics interpret  $\Box\varphi$  as “in all accessible worlds,  $\varphi$  holds to at least degree  $\theta$ .” CLAIR’s  $\Box_c\varphi$  means “the system has confidence  $c$  in its justification for  $\varphi$ .”

*Why CLAIR differs.* Fuzzy modal logic’s graded accessibility relation  $R : W \times W \rightarrow [0, 1]$  is elegant but adds significant complexity to the semantics. For CLAIR’s application (LLM reasoning traces), we prioritize having a simple, implementable semantics over maximal generality.

*What we adopt.* CLAIR adopts product logic operations ( $\otimes$  as multiplication) and shares the insight that many-valued generalization of classical operators is useful for graded reasoning.

### vs. Graded Justification Logic

*Divide:* Milnikel and Fan & Liao add uncertainty degrees to justification terms. CLAIR adds confidence to the propositions themselves and treats justifications as data structures.

*Why CLAIR differs.* Justification Logic focuses on *explicit proof terms*— $t$  (“this term justifies  $\varphi$ ”). CLAIR focuses on *explicit justifications as graphs*—a data structure tracking which beliefs supported which, with provenance and invalidation.

*What we adopt.* CLAIR adopts the core idea that justification should be first-class and tracked, and that justification strength can be graded.

### vs. Belief Revision (AGM)

*Divide:* AGM operates on deductively closed belief sets (sets of sentences). CLAIR operates on justification graphs with annotated beliefs.

*Why CLAIR differs.* AGM’s postulates (especially Recovery) fail when beliefs have graded strength and structured justification. Retracting a belief loses the specific evidence that supported it, not just the belief itself. CLAIR’s revision algorithm explicitly modifies the justification DAG and recomputes confidences, which correctly handles Recovery failure.

*What we adopt.* CLAIR adopts AGM’s core operations (expansion, contraction, revision) but extends them to graded, structured beliefs.

### vs. Type Theory

*Divide:* Standard type systems track types (e.g., “this variable is an integer”). CLAIR tracks epistemic metadata (confidence, justification, provenance, invalidation).

*Why CLAIR differs.* Information flow types and refinement types track *security properties* (who can access this data) and *constraints on values* (this integer is positive). CLAIR tracks

*reasoning properties* (how strongly do we believe this, where did it come from, why do we believe it, when should we reconsider?). This is a new dimension of typing.

*What we adopt.* CLAIR adopts type-theoretic discipline: all metadata is statically checked, composition preserves metadata, and the Curry-Howard correspondence extends to beliefs-with-justifications.

### The Pragmatic Rationale

These design choices are driven by CLAIR’s target use case: making LLM reasoning auditable and trustworthy. For this application, we prioritize:

1. **Interpretability:** Humans should be able to understand why the system believes what it believes.
2. **Compositionality:** Complex reasoning should build from simple, composable operations.
3. **Formal verifiability:** The core system should have machine-checked proofs of key properties.
4. **Implementability:** A reference implementation should be straightforward to build and reason about.

Different applications might justify different choices. A system for scientific modeling, for example, might benefit from Subjective Logic’s three-component opinions. A system for formal verification might benefit from full dependent types. CLAIR’s design is optimized for *AI reasoning trace auditing*, and our divergence from each related tradition reflects that optimization.

## Confidence System

*“Probability is not about what is true. It is about what is reasonable to believe.”*

— E.T. Jaynes, *Probability Theory: The Logic of Science*

The confidence system is the algebraic foundation of CLAIR. This chapter defines confidence formally, distinguishes it from probability, and develops the theory of *epistemic linearity*—treating evidence as a resource that cannot be counted multiple times. We then establish the three monoid structures that govern how confidence values combine, proving key properties in Lean 4 to connect abstract theory to machine-verified implementation.

### Confidence as Calibrated Reliability

#### Semantic Commitments

Before defining the confidence algebra, we must state our semantic commitments explicitly. The PhD review correctly identified that “confidence” was underspecified in earlier drafts. We now clarify:

**Definition Calibrated Reliability.**

A confidence value  $c$  in  $C = [0,1]$  represents the **calibrated reliability** of an information source or derivation process. Specifically:

1. If a source assigns confidence  $c$  to proposition  $\phi$ , this means:

*it.body*

2. Calibration is an *external* property: a source is calibrated if its stated confidences match empirical accuracy over relevant reference classes.
3. CLAIR *tracks* calibrated reliability without *guaranteeing* it. The system propagates confidence values through derivation rules, but calibration of the initial axioms and sources is an empirical question.

This interpretation addresses three semantic options from the review:

1. **Not pure probability:** Confidence does **not** satisfy  $P(\phi) + P(\text{not } \phi) = 1$ . We allow paraconsistent reasoning where both  $\phi$  and  $\text{not } \phi$  may have low confidence.
2. **Not fuzzy truth degree:** Confidence is **not** the degree to which  $\phi$  is true in some multi-valued logic. It is the reliability of the *source asserting  $\phi$* .
3. **Reliability semantics:** Confidence is calibrated reliability of the epistemic process producing the belief. This distinguishes

*it.body*

from

*it.body*

.

## The Problem with Probability

Standard approaches to uncertain reasoning use probability. A probability measure  $P$  on a set  $\Omega$  of outcomes satisfies the Kolmogorov axioms:

1.  $P(A) \geq 0$  (Non-negativity)
2.  $P(\Omega) = 1$  (Normalization)
3.  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$  (Additivity)

For propositions, this implies the fundamental constraint:  $P(\phi) + P(\text{not } \phi) = 1$

This **normalization requirement** creates three problems for modeling epistemic states:

1. **Balanced uncertainty.** An agent confronting an unfamiliar topic may be uncertain about both  $\phi$  and  $\text{not } \phi$ . If asked

*it.body*

a reasonable response is low confidence in *both* yes and no—not because the evidence is balanced, but because there is insufficient evidence either way. Probability forces  $P(\text{yes}) + P(\text{no}) = 1$ , conflating

*it.body*

with

*it.body*

2. **Paraconsistent reasoning.** Evidence sometimes supports both  $\phi$  and  $\text{not } \phi$  before contradiction resolution. A detective might have testimony that a suspect was present (supporting guilt) and an alibi (supporting innocence), without yet knowing which is false. Probability makes this impossible:  $P(\phi) > 0.5$  and  $P(\text{not } \phi) > 0.5$  is a contradiction.
3. **Derivation semantics.** In Bayesian reasoning,  $P(A \text{ and } B) = P(A) \text{ times } P(B|A)$ , where conditioning captures the dependency structure. But for derivation—where  $B$  follows from  $A$  by some rule—there is no clear conditional to use. The semantics of

*it.body*

differs from

*it.body*

## Definition of Confidence

CLAIR's confidence addresses these problems by dropping normalization:

### Definition Confidence Value.

A **confidence value** is a real number  $c$  in  $[0, 1]$ . We write  $C$  for the set of confidence values:  $C = \{c \mid 0 \leq c \leq 1\}$

The semantic interpretation under calibrated reliability:

1.  $c = 1$ : Axiomatic acceptance (treated as foundational)
2.  $c = 0$ : Complete rejection (treated as impossibility)
3.  $c = 0.5$ : Maximal uncertainty (no evidence either direction)
4.  $c > 0.5$ : Net evidence for acceptance
5.  $c < 0.5$ : Net evidence for rejection

### Definition Epistemic Commitment.

Confidence represents **epistemic commitment**: the degree to which an agent commits to a proposition based on available evidence and reasoning, *interpreted* as the calibrated reliability of the source or derivation producing the belief.

Unlike probability:

1. **No normalization**:  $\text{conf}(\phi) + \text{conf}(\text{not } \phi)$  need not equal 1.
2. **Not frequency**:  $\text{conf}(\phi) = 0.9$  does not mean

*it.body*

in a single case. It means

*it.body*

3. **Derivation-based**: Confidence propagates through inference rules, not conditioning.

### Example Non-normalized confidence.

Consider the proposition

*it.body*

An honest assessment might be:  $\text{conf}(\text{no vulnerabilities}) = 0.4$  (some evidence from testing)  $\text{conf}(\text{has vulnerabilities}) = 0.3$  (some evidence from complexity) The sum  $0.4 + 0.3 = 0.7 < 1$  reflects residual uncertainty—neither hypothesis is well-supported. This is inexpressible in probability.

## Falsifiability Criteria

Following the review's requirement, we state explicitly what would falsify CLAIR's confidence model:

1. **Empirical falsification**: If sources systematically assign confidence  $c$  to propositions but are correct with frequency  $f \neq c$ , they are *uncalibrated*. CLAIR provides no mechanism to *detect* uncalibrated sources—this requires external validation.
2. **Semantic falsification**: If interpretation as

*it.body*

leads to contradictions in the algebra (e.g., violations of monoid laws), the semantics would be inadequate. The Lean 4 formalization verifies algebraic consistency.

3. **Competing semantics**: If an alternative interpretation (e.g., pure probability or fuzzy truth) provides better empirical calibration or algebraic properties, CLAIR's semantics would need revision.



## Comparison with Subjective Logic

Jøsang's Subjective Logic (Jøsang, 2016) extends probability with explicit uncertainty. An *opinion* is a tuple  $\omega = (b, d, u, a)$ :

1.  $b$ : belief mass (evidence for)
2.  $d$ : disbelief mass (evidence against)
3.  $u$ : uncertainty mass (lack of evidence)
4.  $a$ : base rate (prior probability)

with constraint  $b + d + u = 1$ .

CLAIR's confidence can be viewed as a simplification of Subjective Logic:  $c = b + u \text{ times } a$  where we collapse uncertainty into the confidence value via the base rate. This loses the  $b/d/u$  decomposition but gains simplicity. The trade-off is appropriate for CLAIR's focus on derivation tracking rather than uncertainty quantification.

*it.body*

## The Aggregation Monoid

When multiple independent pieces of evidence support the same conclusion, confidence should *increase*. This requires the probabilistic OR operation.

### Probabilistic OR (Oplus)

#### Definition Probabilistic OR (Oplus).

For confidence values  $a, b$  in  $C$ , their **aggregation** is:  $a \text{ oplus } b = a + b - a \text{ times } b$

The formula has several equivalent forms:

1.  $a \text{ oplus } b = a + b(1 - a)$
2.  $a \text{ oplus } b = a(1 - b) + b$
3.  $a \text{ oplus } b = 1 - (1 - a)(1 - b)$  (De Morgan duality with multiplication)

The last form reveals the duality:  $a \text{ oplus } b$  is the complement of the product of complements.

#### Theorem Oplus Preserves Bounds.

For all  $a, b$  in  $C$ :  $a \text{ oplus } b$  in  $C$  *Proof. Lower bound*:  $a \text{ oplus } b = a + b(1 - a) \geq 0$  since  $a \geq 0$ ,  $b \geq 0$ , and  $(1 - a) \geq 0$ .

*Upper bound*:  $a \text{ oplus } b = a + b(1 - a) \leq 1$  since  $b \leq 1$  implies  $b(1 - a) \leq 1 - a$ , therefore  $a + b(1 - a) \leq a + (1 - a) = 1$ . ■

#### Theorem Oplus Monoid.

$(C, \text{oplus}, 0)$  is a commutative monoid with absorbing element  $1$ :

1. **Associativity**:  $(a \text{ oplus } b) \text{ oplus } c = a \text{ oplus } (b \text{ oplus } c)$
2. **Commutativity**:  $a \text{ oplus } b = b \text{ oplus } a$
3. **Identity**:  $0 \text{ oplus } a = a \text{ oplus } 0 = a$
4. **Absorption**:  $1 \text{ oplus } a = a \text{ oplus } 1 = 1$

*Proof.* All properties follow from standard real number arithmetic on  $0, 1$ . ■

## Confidence-Increasing Property

Unlike multiplication, *oplus* *increases* confidence:

#### Theorem Oplus is Confidence-Increasing.

For all  $a, b$  in  $C$ :  $a \text{ oplus } b \geq \max(a, b)$  *Proof.* Using form  $a \text{ oplus } b = a + b(1 - a)$ . Since  $b(1 - a) \geq 0$ , we have  $a \text{ oplus } b \geq a$ . By commutativity,  $a \text{ oplus } b \geq b$ . Therefore  $a \text{ oplus } b \geq \max(a, b)$ . ■



**Corollary Diminishing Returns.**

The marginal gain from additional evidence decreases as confidence grows:  $(\frac{\partial}{\partial b})(a \oplus b) = 1 - a$ . When  $a$  is already high, additional evidence contributes less.

**Example Aggregation of Weak Evidence.**

Suppose we have ten independent pieces of weak evidence, each with confidence  $0.3$ . The combined confidence is:  $0.3 \oplus 0.3 \oplus \dots \oplus 0.3$  (10 times)  $= 1 - (1 - 0.3)^{10} = 1 - 0.7^{10} \approx 0.972$ . Ten weak independent witnesses produce high combined confidence.

**The “Survival of Doubt” Interpretation**

The formula  $a \oplus b = 1 - (1-a)(1-b)$  admits a probability-theoretic interpretation under calibrated reliability:

1.  $(1 - a)$  is the

*it.body*

in evidence  $a$

2.  $(1 - a)(1 - b)$  is the probability *both* pieces of evidence fail (assuming independence)
3.  $a \oplus b$  is the probability *at least one* succeeds

This

*it.body*

interpretation explains why aggregation increases confidence: more independent evidence means more chances for at least one to be correct.

**The Multiplication Monoid**

When a conclusion is derived from premises, its confidence depends on the premises' confidences. The simplest case is conjunctive derivation: both premises must hold for the conclusion to follow.

**Conjunctive Confidence Propagation****Definition Confidence Multiplication.**

For confidence values  $a, b$  in  $C$ , their **multiplicative combination** is standard multiplication:  $a \text{ times } b$

This models the intuition that deriving  $C$  from  $A$  and  $B$  requires both to be true. If we are 90% confident in  $A$  and 80% confident in  $B$ , our confidence in  $C$  (derived from both) is at most  $0.9 \text{ times } 0.8 = 0.72$ .

**Theorem Multiplication Preserves Bounds.**

For all  $a, b$  in  $C$ :  $a \text{ times } b$  in  $C$  *Proof.* We prove both bounds:

1.  $a \text{ times } b \geq 0$ : Since  $a \geq 0$  and  $b \geq 0$ , their product is non-negative.
2.  $a \text{ times } b \leq 1$ : Since  $b \leq 1$ , we have  $a \text{ times } b \leq a \text{ times } 1 = a \leq 1$ .

■

**Theorem Multiplication Monoid.**

$(C, \text{times}, 1)$  is a commutative monoid with absorbing element  $0$ :

1. **Associativity:**  $(a \text{ times } b) \text{ times } c = a \text{ times } (b \text{ times } c)$
2. **Commutativity:**  $a \text{ times } b = b \text{ times } a$
3. **Identity:**  $1 \text{ times } a = a \text{ times } 1 = a$
4. **Absorption:**  $0 \text{ times } a = a \text{ times } 0 = 0$

*Proof.* All properties follow from standard real number arithmetic on  $[0,1]$ . ■

## The Derivation Monotonicity Principle

A fundamental property of CLAIR is that derivation can only decrease confidence:

### **Theorem Derivation Monotonicity.**

For all  $a, b$  in  $C$ :  $a \text{ times } b \leq \min(a, b)$  In particular,  $a \text{ times } b \leq a$  and  $a \text{ times } b \leq b$ . *Proof.* Since  $b \leq 1$ , we have  $a \text{ times } b \leq a \text{ times } 1 = a$ . By commutativity,  $a \text{ times } b \leq b$ . Therefore  $a \text{ times } b \leq \min(a, b)$ . ■

### **Corollary No Confidence Amplification.**

No sequence of conjunctive derivations can increase confidence. If  $c_o$  is the confidence of a foundational belief and  $c_n$  is derived through  $n$  multiplicative steps, then  $c_n \leq c_o$ .

This principle is essential for CLAIR's epistemology: derived beliefs are never more confident than their sources. Certainty ( $c = 1$ ) is reserved for axioms, not conclusions.

## Defeat Operations

Beyond derivation and aggregation, CLAIR requires operations for *defeat*: when evidence undermines a belief.

### **Undercut: Attacking the Inference**

Following Pollock (2001), an *undercutting defeater* attacks the inferential link, not the conclusion directly.

#### **Definition Undercut.**

For confidence  $c$  in a conclusion and defeat strength  $d$ :  $\text{undercut}(c, d) = c \text{ times } (1 - d)$

#### **Example Undercutting Defeat.**

Consider the inference:

*it.body*

An undercut

*it.body*

doesn't claim the object isn't red; it undermines the inference from appearance to reality.

If  $\text{conf}(\text{looks red} \Rightarrow \text{is red}) = 0.9$  and  $\text{conf}(\text{red lighting}) = 0.6$ , then:  $\text{undercut}(0.9, 0.6) = 0.9 \text{ times } (1 - 0.6) = 0.9 \text{ times } 0.4 = 0.36$  The inference confidence drops from 0.9 to 0.36.

#### **Theorem Undercut Preserves Bounds.**

For all  $c, d$  in  $C$ :  $\text{undercut}(c, d)$  in  $C$  *Proof.* Since  $d \leq 1$ , we have  $(1 - d) \geq 0$ . Since  $c \geq 0$ , we have  $c \text{ times } (1 - d) \geq 0$ . Since  $c \leq 1$  and  $(1 - d) \leq 1$ , we have  $c \text{ times } (1 - d) \leq 1$ . ■

#### **Theorem Undercut Composition.**

Sequential undercuts compose via *oplus*:  $\text{undercut}(\text{undercut}(c, d_1), d_2) = \text{undercut}(c, d_1 \text{ oplus } d_2)$  *Proof.* By expanding the definition:  $\text{undercut}(\text{undercut}(c, d_1), d_2) = c(1 - d_1)(1 - d_2)$

$$\frac{c(1 - (d_1 + d_2 - d_1 \text{ times } d_2))}{c(1 - (d_1 \text{ oplus } d_2))}$$

$$\frac{c(1 - (d_1 \text{ oplus } d_2))}{\text{undercut}(c, d_1 \text{ oplus } d_2)}$$

$$\text{undercut}(c, d_1 \text{ oplus } d_2)$$

■

This beautiful result shows that defeat strengths aggregate via *oplus*: multiple undercuts combine as if their doubts aggregated.

### Rebut: Competing Evidence

A *rebutting defeater* provides counter-evidence for the conclusion directly.

#### Definition Rebut.

For confidence  $c_{\text{for}}$  in favor and  $c_{\text{against}}$  against:  $\text{rebut}(c_{\text{for}}, c_{\text{against}}) = \text{cases}((c_{\text{for}} / (c_{\text{for}} + c_{\text{against}})), \text{if } c_{\text{for}} + c_{\text{against}} > 0), (0.5, \text{if } c_{\text{for}} = c_{\text{against}} = 0))$

The formula treats evidence symmetrically: the resulting confidence is the

*it.body*

of supporting evidence.

#### Theorem Rebut Preserves Bounds.

For all  $c_{\text{for}}, c_{\text{against}}$  in  $C$ :  $\text{rebut}(c_{\text{for}}, c_{\text{against}})$  in  $C$  Proof. If  $c_{\text{for}} + c_{\text{against}} = 0$ , the result is  $0.5$  in  $[0,1]$ . Otherwise:

1.  $\text{rebut} \geq 0$  because  $c_{\text{for}} \geq 0$  and the denominator is positive.
2.  $\text{rebut} \leq 1$  because  $c_{\text{for}} \leq c_{\text{for}} + c_{\text{against}}$ .

■

#### Theorem Rebut Antisymmetry.

$\text{rebut}(a, b) + \text{rebut}(b, a) = 1$  Proof. When  $a + b > 0$ :  $\text{rebut}(a, b) + \text{rebut}(b, a) = a/(a+b) + b/(a+b) = (a+b)/(a+b) = 1$  When  $a = b = 0$ :  $\text{rebut}(a,b) = \text{rebut}(b,a) = 0.5$ , so the sum is  $1$ . ■

### Rebut Normalization Limitation

The rebut operation has an important limitation that must be stated explicitly: it *normalizes away absolute strength information*.

**Definition Rebut Normalization Property.** For all  $k > 0$ :  $\text{rebut}(k \text{ times } a, k \text{ times } b) = \text{rebut}(a, b)$

Proof. When  $a, b$  are not both zero:  $\text{rebut}(k a, k b) = (k a) / (k a + k b) = k a / (k(a+b)) = a/(a+b) = \text{rebut}(a,b)$  The  $k$  cancels from numerator and denominator. ■

This means rebut only captures the *relative balance* of evidence, not the *absolute magnitude*. Consider:

**Example Normalization Loses Absolute Strength.** Two scenarios with very different absolute evidence strengths:

1. **Scenario A:**  $c_{\text{for}} = 0.1, c_{\text{against}} = 0.1$  Result:  $\text{rebut}(0.1, 0.1) = 0.1/(0.1+0.1) = 0.5$
2. **Scenario B:**  $c_{\text{for}} = 0.9, c_{\text{against}} = 0.9$  Result:  $\text{rebut}(0.9, 0.9) = 0.9/(0.9+0.9) = 0.5$

Both scenarios yield confidence 0.5 despite Scenario B having *nine times more* total evidence. The rebut operation cannot distinguish “weak balanced evidence” from “strong balanced evidence.”

### Implications for CLAIR.

The normalization limitation means rebut is appropriate for *competitive evaluation* (comparing opposing arguments) but insufficient for *absolute assessment*. When absolute strength matters, CLAIR should:

1. Track total evidence magnitude separately:  $total = c_{for} + c_{against}$
2. Use *undercut* instead when the goal is to reduce confidence proportionally to attack strength
3. Apply confidence thresholds to distinguish “weak 0.5” from “strong 0.5”

This is a deliberate design trade-off: rebut prioritizes interpretability (“market share” of evidence) over preserving absolute magnitude.

## Independence Assumptions for Aggregation

The *oplus* operation assumes independence of evidence sources. This is a critical assumption that must be stated explicitly:

### Definition Conditional Independence for Oplus.

Two evidence sources  $e_1$  and  $e_2$  are **independent** with respect to proposition  $\phi$  if:  $P(\phi | e_1, e_2) = P(\phi | e_1) + P(\phi | e_2) - P(\phi | e_1) \times P(\phi | e_2)$  Under calibrated reliability, this means the reference classes of  $e_1$  and  $e_2$  do not systematically overlap.

### When Oplus is Valid

The *oplus* operation is **semantically justified** when:

1. **Independent sources:** Evidence derives from causally independent mechanisms (e.g., different sensors, different witnesses, different reasoning paths).
2. **No shared provenance:** The sources do not derive from common antecedents that would cause their errors to correlate.
3. **Reference class disjointness:** The calibration reference classes for  $e_1$  and  $e_2$  do not systematically overlap.

### Example Valid Independent Aggregation.

Three different image classifiers (trained on different datasets, with different architectures) each assign confidence 0.7 to

*it.body*

. The combined confidence  $0.7 \text{ oplus } 0.7 \text{ oplus } 0.7 \approx 0.973$  is justified because the classifiers make statistically independent errors.

### When Oplus Breaks

The *oplus* operation **overcounts** and produces misleading confidence when:

1. **Shared provenance:** Two sources derive from the same evidence. For example, two newspapers reporting the same press release should not be aggregated via *oplus*.
2. **Common systematic bias:** Sources that share the same misconception or training data flaw will make correlated errors. Aggregating them amplifies bias.
3. **Circular dependence:** When  $e_2$  cites  $e_1$  as a source, their evidence is not independent.

### Example Invalid Aggregation Due to Shared Provenance.

A fact-checking website cites *Source A*. A blog post then cites the fact-checking website. Treating these as independent evidence would incorrectly inflate confidence via *oplus*.

## Correlation-Aware Alternatives

When independence is violated, CLAIR provides alternatives:

1. **Correlation-aware aggregation:** Chapter 4 introduces *oplus\_delta* for correlated evidence, where *delta* in  $[0,1]$  measures dependency:  $a \text{ oplus\_delta } b = a + b - a \text{ times } b - \text{delta}$
2. **Min-based aggregation:** When sources may be completely dependent, use  $\max(a, b)$  instead of *oplus* to avoid overcounting.
3. **Affine type system:** Evidence usage is tracked at type level, preventing the same source from being counted twice in a derivation. This is enforced by CLAIR’s linear type system (Chapter 10).

*Independence Detection.*

In practice, determining whether evidence sources are independent requires domain knowledge and provenance tracking. CLAIR does **not** automatically detect dependence—it provides the algebraic machinery for *manually* specifying correlation or preventing double-counting through the type system.

## Interval-Based Confidence: Dependency Bounds

When dependence between evidence sources is unknown or partially known, CLAIR supports *interval-based confidence* as a robust alternative to point-value aggregation. Instead of committing to a single confidence value, we track upper and lower bounds that reflect uncertainty about dependence.

### Definition Confidence Interval.

A **confidence interval** is a pair  $[c\_low, c\_high]$  with  $0 \leq c\_low \leq c\_high \leq 1$ . The interpretation: the true confidence lies in this interval, with the width  $c\_high - c\_low$  reflecting uncertainty about evidence dependence.

## Interval Aggregation Bounds

For evidence with confidence values  $a$  and  $b$ , the independence assumption determines the correct aggregation:

1. **Independent:**  $a \text{ oplus } b = a + b - a \text{ times } b$  (probabilistic sum)
2. **Maximally dependent:**  $\max(a, b)$  (no double-counting)
3. **Anti-correlated:** Could theoretically reach  $\min(a + b, 1)$

### Theorem Interval Aggregation Bound.

For confidence values  $a, b$  in  $C$ , define:  $\text{oplus\_indep}(a, b) = a \text{ oplus } b = a + b - a \text{ times } b$   $\text{oplus\_dep}(a, b) = \max(a, b)$

Then for any unknown correlation structure:  $\text{oplus\_dep}(a, b) \leq \text{true\_aggregation} \leq \text{oplus\_indep}(a, b)$

*Proof.* The lower bound  $\max(a, b)$  occurs when evidence is fully dependent (one source’s evidence is a subset of the other’s). The upper bound  $a \text{ oplus } b$  occurs when evidence is independent. Any positive correlation reduces the effective aggregation below  $a \text{ oplus } b$ , while negative correlation is impossible for evidential support (evidence cannot be negatively correlated with the truth of what it supports). ■

### Definition Interval Aggregation.

For confidence intervals  $[a\_low, a\_high]$  and  $[b\_low, b\_high]$ :  $[a\_low, a\_high] \text{ oplus\_interval } [b\_low, b\_high] = [\max(a\_low, b\_low), a\_high \text{ oplus } b\_high]$

The lower bound uses *max* (maximal dependence) and the upper bound uses *oplus* (independence).

### Example Interval-Based Aggregation for Unknown Dependence.

Two news sources report the same fact. Each has confidence 0.7, but we suspect possible shared provenance (they may have derived from the same press release).

Point-value approach with *oplus*:  $0.7 \text{ oplus } 0.7 = 0.91$  This assumes independence and may be overconfident.

Interval-based approach:  $[0.7, 0.7] \text{ oplus\_interval } [0.7, 0.7] = [\max(0.7, 0.7), 0.7 \text{ oplus } 0.7] = [0.7, 0.91]$

The interval  $[0.7, 0.91]$  captures our uncertainty: the true confidence could be as low as  $0.7$  (if the sources are fully dependent) or as high as  $0.91$  (if independent). The width  $0.21$  quantifies our dependence uncertainty.

## Interval Propagation

Intervals propagate through CLAIR's operations:

1. **Multiplication:**  $[a\_low, a\_high] \text{ times } [b\_low, b\_high] = [a\_low \text{ times } b\_low, a\_high \text{ times } b\_high]$
2. **Undercut:**  $\text{undercut}([a\_low, a\_high], d) = [a\_low \text{ times } (1-d), a\_high \text{ times } (1-d)]$
3. **Rebut:** More complex; requires propagating both bounds through the ratio

### Theorem Interval Bound Preservation.

If  $[a\_low, a\_high]$  and  $[b\_low, b\_high]$  are valid intervals ( $0 \leq a\_low \leq a\_high \leq 1$ ), then all CLAIR operations produce valid intervals.

*Proof.* Each operation preserves bounds: *oplus* preserves  $0,1$  (Theorem 3.3), multiplication preserves  $0,1$  (Theorem 3.10), and undercut preserves  $0,1$  (Theorem 3.16). Interval extensions apply the same bounds to endpoints. ■

## Dependency Bounds from Provenance Tracking

CLAIR's linear type system (Chapter 10) enables *provenance tracking*, which provides tighter dependency bounds than worst-case intervals.

### Definition Provenance-Aware Dependency Bound.

Let  $S(e)$  be the set of primitive sources (axioms, sensor inputs, LLM outputs) used to derive evidence  $e$ . For two evidences  $e\_1, e\_2$ :

1. If  $S(e\_1)$  and  $S(e\_2)$  are disjoint (disjoint provenance): independence is *plausible*, use *oplus*
2. If  $S(e\_1)$  is a subset of  $S(e\_2)$  or vice versa (subset provenance): maximal dependence, use *max*
3. If there is partial overlap in provenance: use *oplus\_delta* with dependency estimated as the fraction of shared sources

This provides a *heuristic but principled* way to estimate correlation from provenance structure. The intuition: shared evidence sources correlate the derivations that depend on them.

### Example Provenance-Aware Aggregation.

Three sources support a conclusion:

1.  $e\_1$  depends on sensors  $s\_1, s\_2$
2.  $e\_2$  depends on sensors  $s\_2, s\_3$
3.  $e\_3$  depends on sensor  $s\_4$

Aggregating  $e\_1$  and  $e\_2$ : They share sensor  $s\_2$ , so *delta approx 1/2*. Aggregating with  $e\_3$ : Disjoint from  $e\_1, e\_2$ , so independence is plausible.

This is more precise than worst-case intervals and more honest than assuming full independence.

## When to Use Intervals vs Point Values

**Use point values when:**

1. Provenance tracking confirms disjoint sources
2. Domain knowledge guarantees independence
3. Computational efficiency is critical

**Use intervals when:**

1. Provenance is incomplete or unknown
2. Sources may share hidden dependencies
3. Conservative decision-making is required (e.g., safety-critical systems)

4. Auditing requires explicit uncertainty quantification

*Relationship to Imprecise Probability.*

Interval-based confidence in CLAIR is related to Walley’s imprecise probability and Dempster-Shafer theory, but with key differences:

1. **Not a probability interval:** Our intervals bound *calibrated reliability*, not frequency in a reference class.
2. **Not belief/plausibility:** We don’t track separate belief and plausibility masses as in Dempster-Shafer.
3. **Operational:** Intervals quantify uncertainty about *evidence dependence*, not about the proposition itself.

This interval-based approach directly addresses Hole A from the review: when independence assumptions are violated, CLAIR provides principled alternatives that avoid overcounting while preserving auditability.

## Conclusion

This chapter established the algebraic and semantic foundation of CLAIR:

1. **Confidence as calibrated reliability:** We explicitly interpret confidence as the calibrated reliability of information sources and derivation processes. This addresses the review’s concern about underspecified semantics.
2. **Independence assumptions:** The *oplus* operation requires conditional independence of evidence sources. We provide *oplus\_delta* for correlated evidence and affine typing to prevent double-counting.
3. **Three monoids, not a semiring:** Multiplication (derivation), and *oplus* (aggregation) serve distinct semantic roles and do not distribute.
4. **Defeat operations:** Undercut and rebut formalize how evidence can undermine beliefs, with undercuts composing beautifully via *oplus*.
5. **Machine-verified:** The algebra is formalized in Lean 4, ensuring type safety and algebraic correctness.

The confidence system provides the numeric foundation. The next chapter develops the *structural* foundation: how beliefs connect through justification DAGs.

## Justification as Labeled DAGs

*“An argument is not a proof. It is a reason for a belief—and reasons can be defeated.”*

— John L. Pollock, *Defeasible Reasoning*

This chapter develops the structural foundation of CLAIR: how beliefs connect through justification. We show that trees are inadequate for justification structure and that the correct model is a *directed acyclic graph with labeled edges*. The labels distinguish support from defeat, enabling defeasible reasoning where conclusions can be withdrawn when new evidence undermines their justifications.

## The Inadequacy of Trees

### The Shared Premise Problem

Traditional approaches represent justification as trees: each conclusion has premises, which themselves have premises, forming an inverted tree structure. This model is elegant but insufficient.

Consider a simple computation that uses the same belief twice:



```
let population = belief(1000000, 0.95, source: census)
let sample_size = belief(1000, 0.90, source: survey)
let ratio = derive population, sample_size by
  divide
let inverse = derive sample_size, population by
  divide
let product = derive ratio, inverse by multiply
```

In a tree representation, each belief appears multiple times as separate subtrees. This creates three problems:

1. **Space inefficiency:** Beliefs are copied rather than shared.
2. **Invalidation complexity:** If a belief is invalidated, we must find and invalidate all copies.
3. **Semantic confusion:** Are these the *same* belief or *different* beliefs that happen to be equal?

The correct representation is a DAG with explicit sharing, where each belief appears exactly once with multiple parents.

#### **Theorem** DAG Necessity.

Any justification system that:

1. Allows a belief to be used as a premise in multiple derivations
2. Propagates invalidation correctly (removing a premise invalidates all conclusions)
3. Maintains identity (the “same” belief is the same node)

must represent justification as a DAG, not a tree.

*Proof.* In a tree, each node has exactly one parent. If a belief is used in derivations of two different conclusions, it must appear as a child of both. But in a tree, a node cannot have two parents. Therefore, the belief must be duplicated, violating identity. A DAG allows multiple parents, resolving the contradiction. ■

## Why Not Cycles?

If we allow sharing, why not allow cycles? Coherentist epistemology suggests beliefs can mutually support each other. We reject cycles in justification for three reasons:

1. **Bootstrap problem:** Circular justification allows confidence inflation with no external grounding.
2. **Invalidation ambiguity:** If beliefs support each other circularly, invalidation semantics become ill-defined.
3. **Well-foundedness:** The justification relation should be well-founded, with no infinite descending chains.

**Note:** The theory/observation circularity example is better analyzed as two separate relations:

1. **Evidential support** (tracked in justification): Observation provides evidence for theory.
2. **Interpretive framework** (not part of justification): Theory provides framework for interpreting observation.

## Labeled Edges for Defeat

The DAG structure addresses sharing but not defeat. When evidence undermines a belief's justification, we need edges that carry negative, not positive, epistemic weight.

### The Defeat Problem

A *defeater* is a belief that undermines confidence in another belief. Following Pollock (2001), we distinguish:

1. **Undercutters:** Attack the inferential link, not the conclusion directly.
2. **Rebutters:** Provide direct counter-evidence against the conclusion.

#### **Definition** Edge Types.



A justification edge has one of three types:

1. **Support:** Positive evidence for the target.
2. **Undercut:** Attacks the inferential link to the target.
3. **Rebut:** Direct counter-evidence against the target.

## Formal Definition

### Definition Justification Graph.

A **justification graph** is a tuple  $G = (N, E, r)$  where:

1.  $N$  is a finite set of *justification nodes*
  2.  $E \subseteq N \times N \times \text{set}(\text{"support"}, \text{"undercut"}, \text{"rebut"})$  is a set of labeled edges
  3.  $r \in N$  is the root node
- subject to the constraint that the underlying unlabeled graph is acyclic.

### Definition Justification Node Types.

Each node has one of the following types:

1. axiom: Foundational belief (confidence = 1)
2. rule

( $r$ , premises)

: Deductive rule application

3. assumption

( $a$ )

: Temporary assumption for reasoning

4. choice

(options, criteria)

: Decision point

5. abduction

(obs, hypotheses, selected)

: Abductive inference

6. analogy

(source, similarity)

: Analogical reasoning

7. induction

(cases, rule)

: Inductive generalization

8. aggregate

(sources, combRule)

: Evidence aggregation

## Well-Formedness Constraints

A well-formed justification graph must satisfy explicit constraints to ensure semantic coherence and computational tractability.

### Definition Acyclicity Constraint.

set text(size: 9pt, fill: rgb("666")) counter(page).display()

A justification graph must be acyclic: no path may exist from any node back to itself. This constraint applies specifically to support edges. Defeat edges may form cycles, which are resolved via fixed-point iteration (discussed below).

#### **Definition Well-Formed Justification Graph.**

A justification graph is well-formed iff:

1. The support structure is acyclic (no support cycles)
2. Every non-axiom node has at least one incoming support edge
3. Every node is reachable from the root in the underlying undirected graph

Condition 1 ensures the support structure is well-founded. Condition 2 ensures no “floating” beliefs without justification. Condition 3 ensures the graph is connected.

## **DAG-Only vs Cyclic Choice**

CLAIR adopts a hybrid approach to cycles:

1. **Support edges: DAG-only.** Cycles in evidential support are semantically prohibited because they enable bootstrapping and violate well-foundedness. The type checker enforces acyclicity for support edges at construction time.
2. **Defeat edges: Fixed-point semantics.** Defeat cycles are permitted and resolved via fixed-point iteration. When defeat edges form cycles, confidence propagation requires solving a system of equations.

This design choice reflects the epistemic distinction between evidence for (which must be well-founded) and attacks against (which can mutually undermine each other).

## **Fixed-Point Semantics for Cyclic Defeat**

When defeat edges form cycles, we compute confidences via iterative fixed-point finding.

#### **Theorem Fixed-Point Existence.**

For any defeat graph with acyclic underlying support, a fixed point exists.

*Proof.* The confidence update function maps the unit interval to itself and is continuous. By Brouwer’s fixed point theorem, a continuous function from a compact convex set to itself has a fixed point. ■

#### **Theorem Fixed-Point Uniqueness.**

If the maximum product of undercut strengths along any cycle is strictly less than 1, then the fixed point is unique and Kleene iteration converges to it.

*Proof.* Under this condition, the update function is a contraction mapping. By the Banach fixed point theorem, a contraction has a unique fixed point and iteration converges to it. ■

#### **Definition Kleene Iteration.**

Starting from initial confidences, repeatedly apply the update function until convergence. The sequence generated by repeated application is the Kleene iteration.

#### **Example Convergence for Typical Defeat Cycles.**

Suppose two nodes mutually undercut with strength 0.5 each. The Lipschitz constant is 0.5 times 0.5 equals 0.25. Starting from initial confidence (1, 1):

1. After 1 iteration: error is at most 0.25 times initial error
  2. After 5 iterations: error is at most  $0.25^5$  approximately 0.001 (0.1%)
  3. After 10 iterations: error is at most  $0.25^{10}$  approximately  $10^{-6}$  (0.0001%)
- Rapid convergence is typical for well-formed defeat graphs.

### **Practical Implication.**

For CLAIR implementations, we recommend:

set text(size: 9pt, fill: rgb(“666”)) counter(page).display()

1. Enforce DAG structure for support edges via static type checking
2. Allow defeat cycles but limit undercut strengths to ensure contraction
3. Use Kleene iteration with convergence threshold  $10^{-6}$
4. Cache fixed-point solutions to avoid recomputation during updates

## Confidence Propagation

Given a justification graph, we compute the confidence of each node bottom-up.

### Definition Support Propagation.

For a node with children having confidences  $c_1, \dots, c_k$ :

1.  $\text{conf}(\text{axiom}) = 1$
  2.  $\text{conf}(\text{rule}(r, \text{children})) = s_r$  times product over children
  3.  $\text{conf}(\text{aggregate}(\text{children}, \text{independent})) = \text{oplus}$  over children
- where  $s_r$  is the rule strength and *oplus* is probabilistic OR.

### Definition Defeat Propagation.

Let  $c$  be confidence from support edges,  $d_1, \dots, d_m$  be undercut strengths, and  $r_1, \dots, r_n$  be rebut strengths. Then:  $c' = \text{rebut}(\text{undercut}(c, \text{oplus } d_i), \text{oplus } r_j)$

### Theorem Propagation Termination.

The propagation algorithm terminates for any acyclic justification graph.

*Proof.* The graph is acyclic by definition. Each recursive call moves to a node strictly earlier in topological order. Since the graph is finite, recursion terminates. ■

## Reinstatement

A fundamental phenomenon in defeasible reasoning is *reinstatement*: when a defeater is itself defeated, the original belief recovers some confidence.

### Theorem Compositional Reinstatement.

Let  $A$  have base confidence  $a$ , undercut by  $D$  with confidence  $d$ , which is itself undercut by  $E$  with confidence  $e$ . Then:  $\text{conf}(A) = a \text{ times } (1 - d \text{ times } (1 - e))$

The *reinstatement boost* is  $\Delta = a \text{ times } d \text{ times } e$ .

*Proof.* Bottom-up evaluation gives:  $\text{conf\_eff}(D) = d \text{ times } (1 - e)$ , then  $\text{conf}(A) = a \text{ times } (1 - \text{conf\_eff}(D)) = a \text{ times } (1 - d(1-e))$ . ■

## Mutual Defeat

When two arguments defeat each other, we have a defeat cycle. The fixed point analysis yields:

### Theorem Mutual Defeat Fixed Point.

If  $A$  and  $B$  mutually undercut with base confidences  $a$  and  $b$ , the fixed point is:  $a_{\text{star}} = a(1-b) / (1 - ab)$  and  $b_{\text{star}} = b(1-a) / (1 - ab)$

*Proof.* At fixed point:  $a_{\text{star}} = a(1 - b_{\text{star}})$  and  $b_{\text{star}} = b(1 - a_{\text{star}})$ . Solving gives the stated formulas. ■

### Theorem Fixed Point Existence.

For any defeat graph with acyclic underlying support, a fixed point exists.

*Proof.* The confidence update function maps  $[0,1]^n$  to itself and is continuous. By Brouwer's fixed point theorem, a fixed point exists. ■

### Theorem Uniqueness Condition.

If  $b_{\text{max}} \text{ times } d_{\text{max}} < 1$ , the fixed point is unique and iteration converges.

*Proof.* Under this condition, the update function is a contraction mapping. By the Banach fixed point theorem, there is a unique fixed point and iteration converges. ■

## Correlated Evidence

The aggregation formula  $c_1 \text{ oplus } c_2$  assumes independence. When evidence sources are correlated, this overcounts.

### Definition Dependency-Adjusted Aggregation.

For evidence with confidences  $c_1, c_2$  and dependency  $\delta$  in  $[0,1]$ :  $\text{aggregate\_delta}(c_1, c_2) = (1 - \delta)(c_1 \text{ oplus } c_2) + \delta \text{ times } (c_1 + c_2) / 2$   
 where  $\delta = 0$  means independent and  $\delta = 1$  means fully dependent.

## Connection to Prior Art

### Truth Maintenance Systems

JTMS (Doyle 1979) uses IN/OUT lists for dependency-directed backtracking. ATMS (de Kleer 1986) tracks multiple assumption sets.

**CLAIR contribution:** Generalize TMS to graded confidence with the same dependency-directed architecture.

### Argumentation Frameworks

Dung's argumentation (Dung 1995) defines acceptance semantics. Gradual semantics (Amgoud et al. 2017) add numeric degrees.

**CLAIR contribution:** Integrate argumentation's defeat semantics with type-theoretic justification.

### Justification Logic

Artemov's Justification Logic (Artemov 2001) adds explicit justification terms to modal logic.

**CLAIR contribution:** Extend from tree-like justification terms to DAGs with labeled edges, supporting defeasible reasoning.

## The Tracking Paradigm

The preceding sections developed the machinery of justification DAGs, confidence propagation, and defeat semantics. We now step back and articulate the underlying *tracking paradigm* that distinguishes CLAIR from traditional formal systems.

## State Representation

### The Epistemic State

A CLAIR system's **epistemic state** at any point in time is a pair:

#### Definition Epistemic state.

An epistemic state  $\text{cal}(E)$  is a tuple  $\text{cal}(E) = (\text{cal}(G), \text{cal}(B))$  where:

1.  $\text{cal}(G) = \text{set}(G_1, \dots, G_n)$  is a finite set of justification graphs (one for each belief)
2.  $\text{cal}(B) = \text{set}((v_i, c_i, j_i, p_i) \text{ mid } i \text{ in } I)$  is a set of beliefs with values, confidences, justifications, and provenance

Each belief in  $\text{cal}(B)$  corresponds to the root of some graph in  $\text{cal}(G)$ . The graphs may share nodes (shared premises) but are not required to be connected.

### Comparison with Proving Paradigm

Traditional formal systems focus on *provability*:

1. **State:** A set of axioms and inference rules
2. **Question:** Is  $\phi$  provable from the axioms? ( $\Gamma \vdash \phi$ )

### 3. **Output:** Yes/No (with proof term)

CLAIR's tracking paradigm focuses on *epistemic representation*:

1. **State:** A set of labeled graphs with confidence values
2. **Question:** What is the confidence, justification structure, and provenance of belief  $\phi$ ?
3. **Output:**  $(c, G, p)$  where  $c \in [0,1]$  is confidence,  $G$  is the justification graph,  $p$  is provenance

## Update Rules

The epistemic state changes through *epistemic actions*. Each action maps  $cal(E) \rightarrow cal(E)'$ .

## Primitive Actions

### Definition Epistemic actions.

The primitive actions on epistemic states are:

#### 1. **Add belief:**

```
add(phi, c, j, p)
```

creates a new belief with confidence  $c$ , justification  $j$ , and provenance  $p$ .

#### 2. **Aggregate:**

```
aggregate(phi, psi, r)
```

combines two beliefs about the same proposition using rule  $r$  in  $set(\text{"independent"}, \text{"correlated"}(\delta), \text{"max"})$ .

#### 3. **Derive:**

```
derive(phi, [psi_1, ..., psi_k], rule)
```

creates  $\phi$  as a conclusion from premises using an inference rule.

#### 4. **Undercut:**

```
undercut(phi, delta, d)
```

adds an undercutting defeater to  $\phi$  with strength  $d$ .

#### 5. **Rebut:**

```
rebut(phi, delta, d)
```

adds a rebutting defeater to  $\phi$  with strength  $d$ .

#### 6. **Invalidate:**

```
invalidate(psi)
```

removes belief  $\psi$  and propagates invalidation to all beliefs that depend on  $\psi$ .

## Action Semantics

### Example Adding a belief.

When adding a belief

```
add(phi, 0.8, source: testimony, witness_A)
```

:

1. Create a new node  $n_\phi$  with type

```
axiom
```

and confidence 0.8

2. Create a single-node graph  $G_\phi$  containing only  $n_\phi$  with no edges
3. Add  $(n_\phi, 0.8, \text{source: testimony, witness\_A})$  to  $\mathcal{B}$
4. Add  $G_\phi$  to  $\mathcal{G}$

### Example Deriving a conclusion.

When deriving

```
derive(phi, [psi_1, psi_2], modus_ponens)
```

:

1. Create a new node  $n_\varphi$  with type

```
rule(modus_ponens, [n_psi_1, n_psi_2])
```

2. Create support edges  $(n_{\psi_1}, n_\varphi, \text{support})$  and  $(n_{\psi_2}, n_\varphi, \text{support})$
3. Compute confidence via propagation:  $c_\varphi = c_{\psi_1} \times c_{\psi_2}$
4. Add graph  $G_\varphi$  to  $\mathcal{G}$  and belief  $(n_\varphi, c_\varphi, G_\varphi, \text{derived})$  to  $\mathcal{B}$

### Example Invalidation propagation.

When

```
invalidate(psi)
```

is called:

1. Remove  $\psi$  from  $\mathcal{B}$
2. Find all beliefs  $\varphi$  such that  $\psi$  appears in  $G_\varphi$ 's support graph
3. Recursively invalidate each such  $\varphi$  (mark as defeated or recompute without  $\psi$ )
4. Update confidence values via re-propagation

This is the *dependency-directed backtracking* inherited from TMS systems.

## Correctness Criteria

What does it mean for the tracking paradigm to be “correct”? We distinguish three levels of correctness.

### Syntactic Correctness

#### Definition Syntactic correctness.

An epistemic state  $\text{cal}(E)$  is *syntactically correct* if:

1. Every graph in  $\text{cal}(G)$  is acyclic (support edges)
2. Every node has a well-defined type
3. Every edge has a valid label (support, undercut, rebut)
4. All confidence values are in  $[0,1]$

Syntactic correctness is enforced by the type system and checked at runtime. The Lean 4 formalization proves that propagation preserves syntactic correctness.

### Semantic Correctness (Internal)

#### Definition Semantic correctness (internal).

An epistemic state  $\text{cal}(E)$  is *semantically correct* if:

1. Confidence propagation is consistent: Re-computing any belief's confidence yields the same value
2. Defeat semantics are satisfied: Undercuts and rebuts are applied according to their definitions
3. Bounds are preserved: All confidences remain in  $[0,1]$  after any update

This is *internal* correctness: the system behaves according to its own rules. The theorems on propagation termination and sound establish internal correctness for confidence propagation.

### Semantic Correctness (External)

#### Definition Semantic correctness (external).

An epistemic state  $\text{cal}(E)$  is *externally correct* with respect to a reference system  $\text{cal}(R)$  if:

- Calibration:** For any belief with confidence  $c$ , the reference system assigns probability  $c'$  *approx*  $c$  (within calibration tolerance)
  - Justification adequacy:** The justification graph  $G$  accurately represents the actual dependency structure in  $cal(R)$
  - Provenance accuracy:** The provenance field correctly identifies the source of the belief in  $cal(R)$

External correctness cannot be established by formal proof alone—it requires empirical validation. CLAIR provides the *machinery* for tracking (internal correctness) but calibration and accuracy are properties of the *sources*, not the tracking system.

### The Honesty Principle.

- CLAIR’s tracking paradigm embodies epistemic humility:
- The system *reports* its confidence, justification, and provenance
  - It does not *guarantee* that these correspond to external reality
  - External correctness must be validated empirically (calibration studies, provenance audits)
- This distinguishes tracking from proving: a proof claims certainty; a tracked belief admits uncertainty while documenting its grounds.

### Tracking vs. Proving: A Summary

Aspect	Proving Paradigm	Tracking Paradigm
Goal	Establish truth	Record epistemic state
State	Set of formulas	Set of labeled graphs
Query	$\Gamma \mid - \varphi?$	What is the <i>confidence</i> , <i>justification</i> , and <i>provenance</i> of $\varphi$ ?
Output	Proof term (or $\perp$ )	$(c, G, p)$ triple
Update	Add axiom/formula	Add/undercut/rebut/invalidate
Correctness	Soundness	Internal + External
Limits	Hidden or denied	Explicit with confidence

### Practical Implications

The tracking paradigm has several practical consequences for CLAIR as an AI reasoning intermediate representation:

#### Explainability

Every belief carries its justification graph. Queries like “why does the system believe  $\varphi$ ?” can be answered by traversing the graph and presenting the dependency chain.

#### Debugging

When a belief has unexpectedly low confidence, the graph reveals which defeaters are responsible. When confidence is too high, the graph shows whether aggregation rules were misapplied.

#### Revision

New evidence is incorporated by adding nodes and edges. The system automatically recomputes confidences via propagation, handling reinstatement without special cases.

#### Uncertainty

Unlike binary logic, CLAIR tracks degrees of belief. This matches the uncertainty inherent in real-world reasoning and LLM outputs.

### Conclusion

This chapter established the structural foundation of CLAIR and articulated the tracking paradigm:

1. **DAGs, not trees:** Shared premises require graph structure; explicit sharing enables correct invalidation.
2. **Acyclic:** Cycles in evidential support violate well-foundedness; defeat cycles handled via fixed-point semantics.
3. **Labeled edges:** Support, undercut, and rebut serve different epistemic roles.
4. **Compositional reinstatement:** Defeaters being defeated automatically recovers confidence.
5. **Correlated evidence:** Independence assumptions must be explicit; dependency adjustment prevents overcounting.
6. **Tracking paradigm:** CLAIR represents epistemic state as labeled graphs with confidence values, updating via add/derive/undercut/rebut/invalidate operations. Correctness has three levels: syntactic (well-formedness), internal semantic (consistency with propagation rules), and external semantic (calibration to reality).

The justification DAG provides the structural substrate for CLAIR’s beliefs. The tracking paradigm formalizes what it means to “track not prove”—a fundamental shift from establishing truth to documenting epistemic grounds with explicit confidence.

1. **DAGs, not trees:** Shared premises require graph structure; explicit sharing enables correct invalidation.
2. **Acyclic:** Cycles in evidential support violate well-foundedness; defeat cycles handled via fixed-point semantics.
3. **Labeled edges:** Support, undercut, and rebut serve different epistemic roles.
4. **Compositional reinstatement:** Defeaters being defeated automatically recovers confidence.
5. **Correlated evidence:** Independence assumptions must be explicit; dependency adjustment prevents overcounting.

The justification DAG provides the structural substrate for CLAIR’s beliefs. The next chapter addresses a subtler challenge: how beliefs can safely refer to themselves.





# Self-Reference and the Gödelian Limits

line(length: 20%, stroke: 1.5pt, paint: academic-burgundy)

*“If a system is consistent, it cannot prove its own consistency.”*

— Kurt Gödel, On Formally Undecidable Propositions

CLAIR allows beliefs about beliefs. This reflexive capacity creates potential for self-reference: a belief that refers to itself, either directly or through a chain of intermediate beliefs. Such self-reference is both a powerful expressive tool and a source of potential paradox. This chapter develops the theoretical foundations for distinguishing safe from dangerous self-reference, culminating in a novel extension of provability logic to graded confidence.

## The Problem of Self-Reference

### Direct Self-Reference

Consider a belief that directly references itself:

#### Example Direct Self-Reference.

```
-- A belief referencing its own confidence
let b : Belief<Bool> = belief(
  value: confidence(b) > 0.5,
  confidence: ???
)
```

What confidence should  $b$  have? If we assign confidence  $c > 0.5$ , the content becomes true, which seems consistent. But if we assign  $c \leq 0.5$ , the content becomes false—yet what prevents us from assigning  $c = 0.9$  anyway?

This is not merely a curiosity. If CLAIR aims to capture how an LLM reasons, and if introspection is part of reasoning, then CLAIR must account for self-referential beliefs—even if that account restricts or forbids certain patterns.

## Why Self-Reference Matters

Self-reference enables powerful epistemic capabilities:

1. **Calibration:** “My confidence estimates are typically accurate”
2. **Uncertainty tracking:** “I am uncertain about this belief”
3. **Meta-reasoning:** “I should reconsider beliefs derived from unreliable sources”
4. **Self-improvement:** “My reasoning process could be improved in specific ways”

But self-reference also enables paradoxes:

1. **Liar-like:** “This belief has confidence 0” (no consistent assignment)
2. **Curry-like:** “If this belief is true, then arbitrary proposition  $P$ ” (proves anything)
3. **Löbian:** “If I believe  $P$ , then  $P$  is true” (circular self-validation)

The challenge is to permit the former while blocking the latter.

## Löb’s Theorem and Anti-Bootstrapping

### The Classical Result

Löb’s theorem is a cornerstone of provability logic:

#### **Theorem** Löb’s Theorem.

In any sufficiently strong formal system  $T$  containing arithmetic:  $\text{prov}_{T(\text{prov}(\text{prov}(P) \rightarrow P))} \rightarrow \text{prov}(P)$  where  $\text{prov}$  denotes provability in  $T$ .

In words: if a system can prove “if  $P$  is provable, then  $P$  is true,” then the system can prove  $P$ . This has a startling consequence.

#### **Corollary** No Internal Soundness Proof.

No consistent system can prove its own soundness, i.e., cannot prove  $\forall P, \Box(P) \rightarrow P$ .

*Proof.* Suppose system  $T$  proved  $\forall P, \Box(P) \rightarrow P$ . Instantiating with  $P = \text{“false”}$  (falsity), we get  $\text{prov}(\text{false}) \rightarrow \text{false}$ . Combining with consistency ( $\neg \text{prov}(\text{false})$ ), we can derive  $\text{prov}(\text{false})$ , contradicting consistency. ■

### Application to CLAIR

For CLAIR, interpret  $\Box(P)$  as “CLAIR believes  $P$  with confidence 1.0.” Then Löb’s theorem constrains self-soundness beliefs:

```
-- A claimed self-soundness belief
let soundness = belief(
  value: forall P. (belief(P, c, ...) and c > 0.9) -> P is true,
  confidence: 0.95
)
```

By Löb’s theorem, if CLAIR can form this belief with high confidence, then (classically) CLAIR believes everything with high confidence—a collapse to triviality. This is the *bootstrapping trap*: self-soundness claims cannot increase epistemic authority.

#### **Definition** Anti-Bootstrapping Principle.

A belief system satisfies *anti-bootstrapping* if no belief of the form “my beliefs are sound” can increase confidence in any derived belief beyond what the original evidence supports.

Löb’s theorem mathematically enforces anti-bootstrapping for classical provability. The question is how this extends to graded confidence.

# Tarski’s Hierarchy: Stratified Introspection

## The Classical Solution

Tarski’s theorem on the undefinability of truth states that no sufficiently expressive language can define its own truth predicate—on pain of the Liar paradox. Tarski’s solution is stratification:

Level	Can Express	Cannot Express
Level 0 (object)	Facts about the world	Truth of any sentence
Level 1 (meta)	$X_0$ is true for level-0 X	Truth of level-1 sentences
Level 2 (meta-meta)	$X_1$ is true for level-1 X	Truth of level-2 sentences
...	...	...

Each level can discuss truth at lower levels but never its own level.

## Stratified Beliefs in CLAIR

We apply this to beliefs:

**Definition** **Stratified Belief Type.**  
 $\text{Bel}(n, A)$  for  $n \in \mathbb{N}$  where level- $n$  beliefs may reference level- $m$  beliefs only if  $m < n$ .

**Example** **Stratified Beliefs in CLAIR.**

```
-- Level 0: beliefs about the world (no introspection)
type Belief_0<A>

-- Level n: beliefs that may reference level-(n-1) beliefs
type Belief<n : Nat, A> where
  n > 0 implies A may mention Belief<m, B> for any m < n

-- Examples:
let auth : Belief<0, Bool> = belief("user authenticated", 0.9, ...)

let meta_auth : Belief<1, String> = belief(
  "my auth belief has confidence " ++ show(auth.confidence),
  0.95,
  derives_from: [auth]
)

let meta_meta : Belief<2, String> = belief(
  "my level-1 introspection seems accurate",
  0.9,
  derives_from: [meta_auth]
)
```

**Theorem** **Stratification Safety.**  
If all beliefs respect the stratification constraint— $\text{Bel}(n, A)$  references only  $\text{Bel}(m, B)$  with  $m < n$ —then no Liar-like paradox can arise.

*Proof.* Any reference chain from a belief  $b$  must strictly decrease in level. Since  $\mathbb{N}$  has no infinite descending chains, every chain terminates at level 0. Level-0 beliefs contain no belief references, so they cannot participate in self-referential loops. Therefore, no belief can reference itself directly or transitively. ■

## What Stratification Rules Out

Stratification prohibits:

- Direct self-reference:** A belief cannot mention itself (would require level  $n < n$ ).
- Universal introspection:** “All my beliefs are...” spans all levels and cannot be expressed at any finite level.

3. **Self-soundness at a single level:** “My level- $n$  beliefs are sound” would require level  $n + 1$  to express.

## The Cost of Safety

Stratification is safe but restrictive. Some legitimate self-referential reasoning is blocked:

```
-- Legitimate but blocked: calibration beliefs
let calibrated = belief(
  "my confidence estimates match empirical accuracy",
  0.8,
  ...
)
-- This is self-referential (talks about own confidences)
-- but intuitively safe (no paradox)
```

This motivates a more permissive approach for certain cases.

## Kripke’s Fixed Points: Safe Self-Reference

### The Fixed-Point Construction

Kripke proposed an alternative to stratification: allow self-reference but let some sentences remain *undefined*. The key insight is that certain self-referential constructs have *fixed points* — consistent confidence assignments—while others do not.

#### Definition Fixed Point for Self-Referential Belief.

A self-referential belief  $b$  with confidence function  $f: \text{cal}("D")\text{cal}("D")[0,1] \rightarrow \text{cal}("D")\text{cal}("D")[0,1]$  (determining confidence from the assumed truth value) has a fixed point if there exists  $c$  in  $\text{cal}("D")\text{cal}("D")[0,1]$  such that:  $c = f(c)$

#### Example Truth-Teller: Multiple Fixed Points.

Consider:

```
let tt = self_ref_belief(fun self =>
  content: "this belief is true",
  compute_confidence: if val(self.content) then 1.0 else 0.0
)
```

If confidence is 1.0: content is true, so confidence should be 1.0. ✓ If confidence is 0.0: content is false, so confidence should be 0.0. ✓

Both are fixed points. The belief is *underdetermined*.

#### Example Liar: No Fixed Point.

Consider:

```
let liar = self_ref_belief(fun self =>
  content: "this belief has confidence 0",
  compute_confidence: if val(self.content) then 1.0 else 0.0
)
```

If confidence is 1.0: content says “confidence 0,” which is false, so confidence should be 0.0. Contradiction. If confidence is 0.0: content says “confidence 0,” which is true, so confidence should be 1.0. Contradiction.

No fixed point exists. The belief is *ill-formed*.

#### Example Grounded Self-Reference: Unique Fixed Point.

Consider:

```
let careful = self_ref_belief(fun self =>
  content: "confidence(self) is in [0.4, 0.6]",
  compute_confidence: 0.5
)
```

The compute function is constant, so  $f(c) = 0.5$  for all  $c$ . The fixed point is  $c = 0.5$ , which indeed satisfies  $0.5 \in [0.4, 0.6]$ . This belief is *well-formed* with unique confidence 0.5.

## The Self-Reference Escape Hatch

CLAIR provides a controlled mechanism for self-reference:

```
-- Self-referential belief constructor
self_ref_belief :
  {A : Type} ->
  (compute : Belief<infinity, A> -> BeliefContent<A>) ->
  SelfRefResult<A>

data SelfRefResult<A> =
  | WellFormed (Belief<infinity, A>)      -- unique fixed point
  | IllFormed (reason : SelfRefError)     -- no fixed point
  | Underdetermined (points : List<Confidence>) -- multiple fixed points

data SelfRefError =
  | NoFixedPoint      -- Liar-like
  | CurryLike         -- proves anything
  | LobianTrap        -- self-soundness
  | Timeout           -- computation did not terminate
```

The  $\text{Belief}<\infty, A>$  type indicates beliefs that escape the stratification hierarchy—they exist “outside” all finite levels.

## Classification of Self-Reference

Combining Tarski and Kripke, we classify self-referential constructs:

Category	Fixed Points	Status	Example
Grounded	Unique	Safe	Calibration beliefs
Underdetermined	Multiple	Policy choice	Truth-teller
Liar-like	None	Ill-formed	“Confidence is 0”
Curry-like	—	Banned	“If true, then $P$ ”
Löbian	—	Banned	Self-soundness

### Definition Safe Self-Reference.

A self-referential belief is *safe* if it either:

1. Respects stratification (level- $n$  references only level- $m < n$ ), or
2. Has a unique fixed point (Kripke), or
3. Has multiple fixed points with a deterministic policy for selection.

### Definition Dangerous Self-Reference.

A self-referential belief is *dangerous* if it:

1. Has no fixed point (Liar-like), or
2. Matches a Curry pattern (“if this then  $P$ ”), or
3. Claims self-soundness (Löbian trap).

# Provability Logic and CLAIR

## Gödel-Löb Logic (GL)

To formally characterize CLAIR's belief logic, we turn to *provability logic*. The standard modal logic of provability is GL (Gödel-Löb):

### Definition GL Syntax.

$\varphi ::= p \mid \neg\varphi \mid \varphi \wedge \varphi \mid \varphi \vee \varphi \mid \varphi \rightarrow \varphi \mid \Box\varphi$  where  $\Box\varphi$  means  $\varphi$  is provable.

### Definition GL Axioms.

1. **K (Distribution):**  $\Box(\varphi \rightarrow \psi) \rightarrow (\Box\varphi \rightarrow \Box\psi)$
2. **4 (Positive Introspection):**  $\Box\varphi \rightarrow \Box\Box\varphi$
3. **L (Löb):**  $\Box(\Box\varphi \rightarrow \varphi) \rightarrow \Box\varphi$

Critically, GL *lacks* the truth axiom  $\Box\varphi \rightarrow \varphi$  (T). This is philosophically essential: provability does not imply truth. A consistent system can prove false statements if its axioms are wrong.

## GL vs Other Modal Logics

Logic	K	T	4	5 or L
K	✓			
T	✓	✓		
S4	✓	✓	✓	
S5	✓	✓	✓	5
GL	✓		✓	L
CLAIR	✓		✓	L

CLAIR aligns with GL:

1. **K holds:** If CLAIR believes an implication and believes the antecedent, it can derive the consequent.
2. **T fails:** CLAIR's beliefs can be wrong (fallibilism).
3. **4 holds:** CLAIR can have meta-beliefs about its beliefs.
4. **L must hold:** Self-soundness claims cannot bootstrap confidence.

## Solovay's Completeness

### Theorem Solovay Completeness.

GL is sound and complete with respect to:

1. Arithmetic provability:  $GL \vdash \varphi$  iff  $\varphi$  holds under all interpretations of  $\Box$  as Gödel provability in PA.
2. Finite transitive irreflexive Kripke frames.

The completeness for finite frames yields:

### Corollary GL Decidability.

GL is decidable (PSPACE-complete).

This is crucial: classical provability logic is computationally tractable.

## Confidence-Bounded Provability Logic (CPL)

Classical GL uses binary truth: propositions are either provable or not. CLAIR needs a *graded* version where beliefs carry confidence values in  $\text{cal}(\text{"D"})\text{cal}(\text{"D"})[0,1]$ . This section introduces CPL (Confidence-Bounded Provability Logic), a novel extension of GL designed for CLAIR.

## The Literature Gap

Extensive work exists on fuzzy modal logics and graded epistemic logic. However, no prior work addresses:

1. Graded versions of the Löb axiom
2. The interaction of continuous confidence with provability constraints
3. Anti-bootstrapping in the context of graded belief

CPL fills this gap.

## CPL Syntax

### Definition CPL Syntax.

$\varphi ::= p \mid \neg\varphi \mid \varphi \wedge \varphi \mid \varphi \vee \varphi \mid \varphi \rightarrow \varphi \mid \Box_c \varphi$  where  $\Box_c \varphi$  means  $\varphi$  is believed with confidence at least  $c$ .

## CPL Semantics

### Definition Graded Kripke Frame.

A *graded Kripke frame* is a tuple  $(W, R)$  where:

1.  $W$  is a non-empty set of worlds
2.  $R : W \times W \rightarrow \mathcal{D}[0, 1]$  is a graded accessibility relation

satisfying:

1. **Transitivity:**  $R(w, v) \cdot R(v, u) \leq R(w, u)$
2. **Converse well-foundedness:** No infinite sequence  $w_0, w_1, w_2, \dots$  with  $R(w_{i+1}, w_i) > 0$  for all  $i$

### Definition Graded Valuation.

A *graded valuation* on a frame  $(W, R)$  assigns to each world  $w$  and proposition  $p$  a confidence value  $V_{w(p)} \in \mathcal{D}[0, 1]$ . Extended to formulas:

1.  $V_w(\neg\varphi) = 1 - V_w(\varphi)$
2.  $V_w(\varphi \wedge \psi) = V_w(\varphi) \cdot V_w(\psi)$
3.  $V_w(\varphi \vee \psi) = V_w(\varphi) + V_w(\psi) - V_w(\varphi) V_w(\psi)$
4.  $V_w(\varphi \rightarrow \psi) = \sup\{c \in \mathcal{D}[0, 1] : V_w(\varphi) \cdot c \leq V_w(\psi)\}$
5.  $V_w(\Box_c \varphi) = \inf_{v: R(w, v) \geq c} V_v(\varphi)$

The last clause says:  $\Box_c \varphi$  is believed at confidence  $c$  if  $\varphi$  holds in all worlds accessible with strength at least  $c$ .

## The Graded Löb Axiom: DESIGN AXIOM

The crucial innovation in CPL is the graded analogue of Löb's axiom.

*Axiom Status Statement.*

The Graded Löb axiom is a **DESIGN AXIOM**, not a semantic theorem derived from more basic principles. It is motivated by the requirement of anti-bootstrapping and the need to extend GL to graded confidence while preserving its essential character. The axiom is **posited** as part of CPL's definition, not **proved** from the semantics.

The key question is: Is CPL consistent? We address this below by exhibiting a non-trivial model satisfying all CPL axioms.

### Definition Graded Löb Axiom (Design Axiom).

$\Box_c(\Box_{g(c)} \varphi \rightarrow \varphi) \rightarrow \Box_{g(c)} \varphi$  where  $g : \mathcal{D}[0, 1] \rightarrow \mathcal{D}[0, 1]$  is a *discount function* satisfying  $g(c) \leq c$ .

This is a **DESIGN AXIOM**—not derived from more basic principles but posited as part of CPL's definition. The axiom is motivated by anti-bootstrapping requirements.

The function  $g$  captures the *cost* of self-soundness claims. If you believe at confidence  $c$  that “believing  $\varphi$  at  $c$  implies  $\varphi$ ,” you can derive  $\varphi$  only at the discounted confidence  $g(c)$ .

## Choosing the Discount Function

We require  $g$  to satisfy:

set text(size: 9pt, fill: rgb("666")) counter(page).display()

1. **Boundedness:**  $g : \mathcal{D}[0, 1] \rightarrow \mathcal{D}[0, 1]$
2. **Non-amplification:**  $g(c) \leq c$  for all  $c$
3. **Monotonicity:**  $c_1 \leq c_2 \Rightarrow g(c_1) \leq g(c_2)$
4. **Anchoring:**  $g(0) = 0$  and  $g(1) = 1$
5. **Non-triviality:**  $g(c) < c$  for  $c \in (0, 1)$

After analyzing several candidates (identity, parabolic, constant offset, product), we recommend:

**Definition Quadratic Discount.**

$$g(c) = c^2$$

**Theorem Quadratic Discount Properties.**

The quadratic discount  $g(c) = c^2$  satisfies all desiderata and:

1. Aligns with CLAIR's multiplicative confidence algebra ( $c^2 = c \times c$ )
2. Has intuitive meaning: self-soundness costs “deriving the claim twice”
3. Produces strong anti-bootstrapping: iterated application  $c \rightarrow c^2 \rightarrow c^4 \rightarrow \dots \rightarrow 0$

*Proof.* Boundedness and anchoring are immediate ( $c \in \mathcal{D}[0, 1] \Rightarrow c^2 \in \mathcal{D}[0, 1]$ ,  $0^2 = 0$ ,  $1^2 = 1$ ). For non-amplification:  $c^2 \leq c$  when  $c \leq 1$ , with equality only at 0 and 1. Monotonicity:  $c_1 \leq c_2 \Rightarrow c_1^2 \leq c_2^2$  on  $\mathcal{D}[0, 1]$ . Non-triviality:  $c^2 < c$  for  $c \in (0, 1)$ . ■

## The Anti-Bootstrapping Theorem

**Theorem Anti-Bootstrapping.**

In CPL with  $g(c) = c^2$ :  $\text{conf}(\Box_{c(\Box_c \varphi \rightarrow \varphi)}) = c \Rightarrow \text{conf}(\varphi) \leq c^2 < c$  Consequently, no finite chain of self-soundness claims can increase confidence beyond the initial level.

*Proof.* Applying the Graded Löb axiom to the hypothesis yields  $\text{conf}(\Box_{c^2} \varphi) \leq c^2$ . Iterating:  $c \rightarrow c^2 \rightarrow c^4 \rightarrow c^8 \rightarrow \dots$  For any  $c < 1$ , this sequence converges to 0. Self-soundness claims can only decrease confidence. ■

This is the mathematical formalization of anti-bootstrapping: claiming your own soundness provides no epistemic free lunch.

## Modal Axioms in CPL

We now explicitly list the modal axioms CPL adopts and their status:

**Definition CPL Modal Axiom Status.**

1. **K (Distribution):**  $\Box_{c(\varphi \rightarrow \psi)} \rightarrow (\Box_c \varphi \rightarrow \Box_c \psi)$  — **VALID** Derivable from the semantics via graded accessibility.
2. **4 (Positive Introspection):**  $\Box_c \varphi \rightarrow \Box_c \Box_c \varphi$  — **VALID** Follows from transitivity of the accessibility relation.
3. **GL/Graded Löb:**  $\Box_{c(\Box_c \varphi \rightarrow \varphi)} \rightarrow \Box_{g(c)} \varphi$  — **DESIGN AXIOM** Posited as part of CPL; motivated by anti-bootstrapping requirements.
4. **T (Reflexivity/Truth):**  $\Box_c \varphi \rightarrow \varphi$  — **INVALID** Explicitly **rejected** in CPL. Provability/belief does not imply truth. This is essential for fallibilism.

## CPL Consistency

To establish CPL's consistency, we exhibit a non-trivial model:

**Theorem CPL Consistency.**

CPL is consistent. There exists a non-trivial graded Kripke model satisfying all CPL axioms including the Graded Löb axiom.

*Proof. Proof Sketch.*

Consider the frame  $(W, R)$  where:



1.  $W = \mathbb{N}$  (the natural numbers)
2.  $R(i, j) = 2^{-i}$  if  $i < j$ , and  $R(i, j) = 0$  otherwise

This frame satisfies:

1. **Transitivity:** If  $i < j < k$ , then  $R(i, j) \times R(j, k) = 2^{-i} \times 2^{-j} \leq 2^{-i} = R(i, k)$
2. **Converse well-foundedness:** No infinite sequence  $w_0, w_1, \dots$  with  $R(w_{i+1}, w_i) > 0$ , since this would require an infinite decreasing sequence of natural numbers.

Define valuation  $V_{w(\varphi)} = 1$  for all propositional variables  $\varphi$  at all worlds  $w$ .

For the Graded Löb axiom, consider any world  $w$ . If  $V_{w(\Box_c(\Box_c\varphi \rightarrow \varphi))} = 1$ , then for all  $v$  with  $R(w, v) \geq c$ , we have  $V_{v(\Box_c\varphi \rightarrow \varphi)} = 1$ . By the structure of  $R$ , this forces  $V_{v(\varphi)} = 1$  for all accessible worlds, yielding  $V_{w(\Box_{c^2}\varphi)} = 1$ .

This model is non-trivial (not all formulas are valid) yet satisfies all CPL axioms, establishing consistency. ■

## Decidability of CPL

Classical GL is decidable. Does CPL inherit this property?

### The Vidal Result

#### **Theorem** Vidal's Undecidability Theorem.

Transitive modal logics over many-valued semantics (including Łukasiewicz and Product algebras) are undecidable, even when restricted to finite models.

CPL has transitivity (axiom 4) and continuous  $\mathcal{D}[0, 1]$  values. The Vidal result strongly suggests:

#### **Conjecture** CPL Undecidability.

Full CPL (with continuous  $\mathcal{D}[0, 1]$  confidence) is undecidable.

*Confidence: 0.80*

We assign confidence 0.80 to this conjecture based on the close analogy to Vidal's proof technique.

### The Role of Converse Well-Foundedness

GL's decidability relies on the finite model property: converse well-foundedness forces finite-depth evaluation. Could this rescue CPL?

#### **Proposition** Insufficient for Decidability.

Converse well-foundedness alone does not rescue CPL from undecidability.

*Justification:* Converse well-foundedness constrains *structure* (no infinite ascending chains) but not *values*. The encoding power of continuous  $\mathcal{D}[0, 1]$  values combined with transitivity enables undecidable problem encodings even in well-founded frames.

### Decidable Fragments

Despite the likely undecidability of full CPL, we identify decidable fragments:

#### CPL-finite: Discrete Confidence

Restrict confidence to a finite lattice instead of continuous  $\mathcal{D}[0, 1]$ :

##### **Definition** CPL-finite.

Let  $L_n = \{0, \frac{1}{n-1}, \frac{2}{n-1}, \dots, 1\}$ . CPL-finite evaluates over  $L_n$  with discretized operations:

1.  $a \times b = \lfloor a \times b \rfloor$
2.  $a + b = \lceil a + b - a \times b \rceil$
3.  $g_{L(c)} = \lfloor c^2 \rfloor$

For  $L_5 = \{0, 0.25, 0.5, 0.75, 1\}$ :

$c$	$c^2$	$g_{L.5}(c)$
0	0	0
0.25	0.0625	0

$c$	$c^2$	$g_{L.5.}(c)$
0.5	0.25	0.25
0.75	0.5625	0.5
1	1	1

**Theorem CPL-finite Decidability.**

CPL-finite is decidable via the finite model property.

*Proof Sketch.*

By the theorem of Bou, Esteva, and Godo, many-valued modal logics over finite residuated lattices are decidable. CPL-finite evaluates over  $L_n$ , a finite lattice. The frame constraints (transitivity, converse well-foundedness) are expressible, and finitely many models of bounded size suffice for completeness.

A complete formal proof would establish: (1)  $L_n$  forms a finite residuated lattice under the discretized operations; (2) the frame conditions are expressible in the corresponding modal logic; (3) the finite model property holds; and (4) decidability follows from (3). The proof follows the standard technique for finite-valued modal logics.

**Conjecture CPL-finite Complexity.**

CPL-finite is PSPACE-complete, analogous to classical GL.

### CPL-0: Stratified Only

Restrict to stratified beliefs without any self-reference:

**Definition CPL-o.**

CPL-o disallows nesting of  $\Box$  operators that would require the Löb axiom. Formally: only formulas of the form  $\Box_c \varphi$  where  $\varphi$  is box-free.

**Theorem CPL-o Decidability.**

CPL-o is decidable (trivially: the restricted syntax avoids undecidability).

### Trade-offs

Fragment	Decidable?	Expressiveness	Use Case
Full CPL	Likely no	Full	Theoretical analysis
CPL-finite	Yes	Discrete confidence	Type-level checks
CPL-o	Yes	No self-reference	Stratified beliefs

### Alternative: CPL-Gödel

An alternative approach uses Gödel algebra (min/max) instead of product operations:

**Definition CPL-Gödel.**

- $a \times b = \min(a, b)$
- $a + b = \max(a, b)$

**Conjecture CPL-Gödel Decidability.**

CPL-Gödel is decidable because Gödel modal logic has the finite model property via quasimodels.

*Confidence: 0.75*

The conjecture follows from the known decidability of Gödel modal logic, but requires verification that the graded Löb axiom preserves this property.

However, CPL-Gödel is *semantically inappropriate* for CLAIR:

- max fails aggregation:**  $\max(0.6, 0.6) = 0.6$ , but two independent pieces of evidence should yield higher confidence (0.84 with +).

2. **min lacks degradation:**  $\min(a, a) = a$ , but derivation should cost confidence.
3. **No algebraic discount:** The  $c^2$  discount becomes purely frame-based, losing the anti-bootstrapping semantics.

*Recommendation.*

For CLAIR, use CPL-finite (with product operations), not CPL-Gödel. Accept the discretization rather than sacrifice semantic fidelity.

## “Conservative Over GL”: Clarification

On “Conservative Over GL” Claims.

The phrase “CPL is conservative over GL” requires careful definition. Two interpretations:

1. **Proof-theoretic conservativity:** Every theorem of GL (as formulas with implicit confidence 1) is a theorem of CPL. This **holds**: CPL includes all GL axioms as special cases.
2. **Semantic conservativity:** Every model of GL can be embedded in a model of CPL. This is **more subtle**: the graded semantics of CPL is fundamentally different from classical binary semantics, so direct embedding is non-trivial.

We assert the first interpretation: CPL extends GL conservatively in the sense that all classical GL theorems remain valid in CPL when interpreted as high-confidence beliefs. The second interpretation remains an open question.

## Design Recommendations for CLAIR

### The Two-Layer Approach

CLAIR should implement a two-layer approach to self-reference:

1. **Default: Stratification.** All beliefs are level-indexed.  $\text{Bel}(n, A)$  can only reference  $\text{Bel}(m, B)$  with  $m < n$ . This is safe by construction and requires no runtime analysis.
2. **Escape hatch: Kripke fixed points.** For legitimate self-reference (calibration, uncertainty tracking), use

`self_ref_belief`

which computes fixed points at construction time. Ill-formed constructs are rejected.

### Hard Bans

Certain patterns are syntactically rejected:

1. **Curry patterns:** “If [self-reference] then [arbitrary  $P$ ]”
2. **Explicit self-soundness:** Claims of the form “All my beliefs are sound”
3. **Unrestricted quantification:** “For all beliefs  $b$ , ...”

These are detected by the parser and rejected before type checking.

### Type-Level Anti-Bootstrapping

For type-level confidence checks, use CPL-finite with  $L_5$ :

```
-- Finite confidence for compile-time checks
inductive FiniteConfidence where
| zero  : FiniteConfidence -- 0
| low   : FiniteConfidence -- 0.25
| mid   : FiniteConfidence -- 0.5
| high  : FiniteConfidence -- 0.75
| one   : FiniteConfidence -- 1

def loebDiscount : FiniteConfidence -> FiniteConfidence
| .zero => .zero
| .low  => .zero -- 0.25^2 = 0.0625 -> floor to 0
```

```
| .mid => .low    -- 0.5^2 = 0.25  
| .high => .mid   -- 0.75^2 = 0.5625 -> floor to 0.5  
| .one  => .one
```

This provides decidable type-level constraints while preserving the anti-bootstrapping semantics.

## Related Work

### Provability Logic

The foundations of provability logic are in Boolos's work, with the Solovay completeness theorems establishing the connection to arithmetic. Modern work on GL extensions includes Beklemishev (2004) on polymodal variants.

### Self-Reference in AI

Garrabrant et al. (2016) develop logical inductors as an approach to coherent self-reference, though in a different formal framework.

### Fuzzy Modal Logic

Fuzzy extensions of modal logic are surveyed in Godo et al. (2003). Decidability results for finite-valued logics appear in Bou et al. (2011). The critical undecidability result for transitive many-valued logics is Vidal (2019).

## Conclusion

This chapter characterized the landscape of self-reference in CLAIR:

1. **Löb's theorem applies:** Self-soundness claims cannot bootstrap epistemic authority. This is a mathematical fact, not a design choice.
2. **Stratification is safe:** Tarski-style level indexing prevents all self-referential paradoxes by construction.
3. **Fixed points enable safe self-reference:** Kripke's approach permits legitimate introspection (calibration, uncertainty tracking) while rejecting ill-formed constructs.
4. **CPL extends GL to graded confidence:** The Graded Löb axiom with  $g(c) = c^2$  captures anti-bootstrapping for continuous confidence. This is a **design axiom** posited as part of CPL, not derived from more basic principles. CPL is consistent, as demonstrated by the existence of non-trivial models.
5. **Full CPL is likely undecidable:** Transitivity plus continuous values enables undecidability (Vidal 2019).
6. **CPL-finite is decidable:** Restricting to discrete confidence yields a tractable fragment suitable for type-level checks.
7. **Two-layer design:** Stratification by default, Kripke fixed points as escape hatch, hard bans on dangerous patterns.

The Gödelian limits are not obstacles but design constraints. They tell us what epistemic claims are coherent and which collapse into triviality. By respecting these limits, CLAIR achieves honest self-awareness: it can reason about its own reasoning without falling into paradox.

The next chapter turns to epistemological foundations: what grounds CLAIR's beliefs in the first place.

## Grounding

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua quaerat voluptatem. Ut enim aequaleamur animo,

cum corpore dolemus, fieri tamen permagna accessio potest, si aliquod aeternum et infinitum impendere malum nobis opinemur. Quod idem licet transferre in voluptatem, ut postea variari voluptas distinguere possit, augeri amplificarique non possit. At etiam Athenis, ut e patre audiebam facete et urbane Stoicos irridente, statua est in quo a nobis philosophia defensa et collaudata est, cum id, quod maxime placeat, facere possimus, omnis voluptas assumenda est, omnis dolor repellendus. Temporibus autem quibusdam et aut officiis debitis aut rerum necessitatibus saepe eveniet, ut et voluptates repudiandae sint et molestiae non recusandae. Itaque earum rerum defuturum, quas natura non depravata desiderat. Et quem ad me accedis, saluto: 'chaere,' inquam, 'Tite!' lictores, turma omnis chorusque: 'chaere, Tite!' hinc hostis mi Albucius, hinc inimicus. Sed iure Mucius. Ego autem mirari satis non queo unde hoc sit tam insolens domesticarum rerum fastidium. Non est omnino hic docendi locus; sed ita prorsus existimo, neque eum Torquatum, qui hoc primus cognomen invenerit, aut torquem illum hosti detraxisse, ut aliquam ex eo est consecutus? – Laudem et caritatem, quae sunt vitae.

## Belief Revision

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magnam aliquam quaerat voluptatem. Ut enim aequaleam animo, cum corpore dolemus, fieri tamen permagna accessio potest, si aliquod aeternum et infinitum impendere malum nobis opinemur. Quod idem licet transferre in voluptatem, ut postea variari voluptas distinguere possit, augeri amplificarique non possit. At etiam Athenis, ut e patre audiebam facete et urbane Stoicos irridente, statua est in quo a nobis philosophia defensa et collaudata est, cum id, quod maxime placeat, facere possimus, omnis voluptas assumenda est, omnis dolor repellendus. Temporibus autem quibusdam et aut officiis debitis aut rerum necessitatibus saepe eveniet, ut et voluptates repudiandae sint et molestiae non recusandae. Itaque earum rerum defuturum, quas natura non depravata desiderat. Et quem ad me accedis, saluto: 'chaere,' inquam, 'Tite!' lictores, turma omnis chorusque: 'chaere, Tite!' hinc hostis mi Albucius, hinc inimicus. Sed iure Mucius. Ego autem mirari satis non queo unde hoc sit tam insolens domesticarum rerum fastidium. Non est omnino hic docendi locus; sed ita prorsus existimo, neque eum Torquatum, qui hoc primus cognomen invenerit, aut torquem illum hosti detraxisse, ut aliquam ex eo est consecutus? – Laudem et caritatem, quae sunt vitae.

## Multi-Agent

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magnam aliquam quaerat voluptatem. Ut enim aequaleam animo, cum corpore dolemus, fieri tamen permagna accessio potest, si aliquod aeternum et infinitum impendere malum nobis opinemur. Quod idem licet transferre in voluptatem, ut postea variari voluptas distinguere possit, augeri amplificarique non possit. At etiam Athenis, ut e patre audiebam facete et urbane Stoicos irridente, statua est in quo a nobis philosophia defensa et collaudata est, cum id, quod maxime placeat, facere possimus, omnis voluptas assumenda est, omnis dolor repellendus. Temporibus autem quibusdam et aut officiis debitis aut rerum necessitatibus saepe eveniet, ut et voluptates repudiandae sint et molestiae non recusandae. Itaque earum rerum defuturum, quas natura non depravata desiderat. Et quem ad me accedis, saluto: 'chaere,' inquam, 'Tite!' lictores, turma omnis chorusque: 'chaere, Tite!' hinc hostis mi Albucius, hinc inimicus. Sed iure Mucius. Ego autem mirari satis non queo unde hoc sit tam insolens domesticarum rerum fastidium. Non est omnino hic docendi locus; sed ita prorsus

existimo, neque eum Torquatam, qui hoc primus cognomen invenerit, aut torquem illum hosti detraxisse, ut aliquam ex eo est consecutus? – Laudem et caritatem, quae sunt vitae.

# Formal Verification

## The Case for Machine-Checked Proofs

This chapter presents the Lean 4 formalization of CLAIR...

## Working Interpreter

The Lean 4 formalization includes a complete working interpreter...

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magnam aliquam quaerat voluptatem. Ut enim aequale doleamus animo, cum corpore dolemus, fieri tamen permagna accessio potest, si aliquod aeternum et infinitum impendere malum nobis opinemur. Quod idem licet transferre in voluptatem, ut postea variari voluptas distinguere possit, augeri amplificarique non possit. At etiam Athenis, ut e patre audiebam facere et urbane Stoicos irridere, statua est in quo a nobis philosophia defensa et collaudata est, cum id, quod maxime placeat, facere possimus, omnis voluptas assumenda est, omnis dolor repellendus. Temporibus autem quibusdam et aut officiis debitis aut rerum necessitatibus saepe eveniet, ut et voluptates repudiandae sint et molestiae non recusandae. Itaque earum rerum defuturum, quas natura non depravata desiderat. Et quem ad me accedis, saluto: 'chaere,' inquam, 'Tite!' lictores, turma omnis chorusque: 'chaere, Tite!' hinc hostis mi Albucius, hinc inimicus. Sed iure Mucius. Ego autem mirari satis non queo unde hoc sit tam insolens domesticarum rerum fastidium. Non est omnino hic docendi locus; sed ita prorsus existimo, neque eum Torquatam, qui hoc primus cognomen invenerit, aut torquem illum hosti detraxisse, ut aliquam ex eo est consecutus? – Laudem et caritatem, quae sunt vitae.

# Implementation

This chapter documents the CLAIR reference implementation in Haskell, demonstrating that CLAIR is computable and providing a practical foundation for LLM integration. The implementation covers the full language pipeline: parsing, type checking, and evaluation.

## 10.1 Overview

The Haskell implementation serves as a *reference interpreter* for CLAIR. It is designed to be:

+—+—+ | Property | Description | +—+—+ | Executable | Full parser, type checker, and evaluator | | Type-safe | Haskell’s type system prevents many errors | | Testable | QuickCheck properties verify algebraic laws | | Auditable | Justification tracking preserved throughout | | Extensible | Modular design for adding features | +—+—+

The implementation is available at

implementation/haskell/

in the project repository and can be built with Cabal:

```
cabal build
cabal test
cabal run clair-repl
```

## 10.2 Module Architecture

The implementation is organized into six core modules:

### 10.2.1 CLAIR.Syntax

The

CLAIR.Syntax

module defines the abstract syntax trees for CLAIR expressions:

```
-- | A variable or identifier name
newtype Name = Name { getName :: Text }

-- | CLAIR types
data Type
  = TBase BaseType      -- Nat, Bool, String, Unit
  | TFun Type Type      -- Function type A -> B
  | TBelief Confidence Type -- Belief type at confidence c
  | TJustification      -- Justification graph type
  | TProvenance         -- Provenance tracking type

-- | Core expression language
data Expr
  = EVar Name                -- Variable: x
  | ELam Name Type Expr      -- Lambda: lambda x:A. e
  | EApp Expr Expr           -- Application: e1 e2
  | EAnn Expr Type           -- Type annotation: e : A
  | EBelief Belief           -- Belief: belief(v,c,j,i,p)
  | EBox Confidence Expr     -- Self-reference: box_c e
  | EPrim Op Expr Expr      -- Primitive operations
  | ELit Literal             -- Literals
```

The

Belief

type captures all annotations required for auditable reasoning:

```
data Belief = Belief
{ beliefValue      :: Expr      -- The proposition/content
, beliefConf       :: Confidence -- Confidence level in [0,1]
, beliefJustify    :: Justification -- Supporting arguments
, beliefInvalidate :: Invalidity  -- Defeating information
, beliefProvenance :: Provenance  -- Source tracking
}
```

### 10.2.2 CLAIR.Confidence

The

CLAIR.Confidence

module implements the confidence algebra from Chapter 3:

```
-- | A confidence value in the closed interval [0,1]
newtype Confidence = Confidence Double
```

```

-- | Probabilistic sum: a oplus b = 1 - (1-a)(1-b)
oplus :: Confidence -> Confidence -> Confidence
oplus (Confidence a) (Confidence b) = clamp (a + b - a * b)

-- | Product t-norm: a otimes b = a * b
otimes :: Confidence -> Confidence -> Confidence
otimes (Confidence a) (Confidence b) = clamp (a * b)

-- | Negation: oneg a = 1 - a
oneg :: Confidence -> Confidence
oneg (Confidence a) = clamp (1 - a)

-- | Apply undercut: multiply by (1-d)
undercut :: Defeat -> Confidence -> Confidence
undercut (Defeat d) (Confidence c) = clamp (c * (1 - d))

-- | Apply rebut with normalization
rebut :: Defeat -> Confidence -> Confidence -> Confidence
rebut (Defeat d_strength) (Confidence c_for) (Confidence c_against_base) =
  let c_against = d_strength * c_against_base
      total = c_for + c_against
  in if total == 0
      then Confidence 0.5 -- ignorance prior
      else clamp (c_for / total)

```

Key design decisions:

**clamp**

**ensures all results stay in valid range**

**[0,1]**

**Double**

**representation is efficient; production may use**

**Rational**

**Rebuttal uses an ignorance prior of 0.5 when no evidence exists**

### 10.2.3 CLAIR.Parser

The

**CLAIR.Parser**

module implements a recursive descent parser using Parsec:

```

-- | Parse a CLAIR expression
parseExpr :: String -> Either CLAIRParseError Expr

-- | Parse a program (sequence of expressions)
parseProgram :: String -> Either CLAIRParseError [Expr]

```

Example supported syntax:



```
-- Belief formation
belief("rain", 0.8, [], none, none)

-- Lambda abstraction
lambda x:Nat. x + 1

-- Self-reference with confidence discount
box_0.9 belief("proposition", 0.8, [], none, none)

-- Let binding (desugars to lambda application)
let x = belief("A", 0.7, [], none, none)
in belief("B", 0.8, [x], none, none)
```

The parser produces the same AST used by the type checker and evaluator, ensuring consistency throughout the pipeline.

### 10.2.4 CLAIR.TypeChecker

The

```
CLAIR.TypeChecker
```

module implements *bidirectional type checking*:

$\vdash$  | Mode | Judgment | Purpose |  $\vdash$  | Synthesis |  $\Gamma \vdash e \rightarrow \tau$   
 | Infer type from expression structure | | Checking |  $\Gamma \vdash e \rightarrow \tau$  | Verify  
 expression has expected type |  $\vdash$

Bidirectional checking improves error messages and handles implicit arguments naturally:

```
-- | Infer (synthesize) the type of an expression
infer :: Context -> Expr -> Either TypeError TCRResult

-- | Check that an expression has the given type
check :: Context -> Expr -> Type -> Either TypeError Context
```

Key typing rules implemented:

**Theorem Variable Synthesis.**  $\Gamma \vdash x \rightarrow \Gamma(x)$   
 If variable  $x$  is in context with type  $\tau$ , synthesize  $\tau$ .

**Theorem Application Synthesis.** If  $\Gamma \vdash e_1 \rightarrow \tau_1 \rightarrow \tau_2$  and  $\Gamma \vdash e_2 \rightarrow \tau_2$ , then  $\Gamma \vdash e_1 e_2 \rightarrow \tau_2$ .  
 This is standard function application: infer function type, check argument type, synthesize result type.

**Theorem Lambda Checking.** If  $\Gamma, x:\tau_1 \vdash e \rightarrow \tau_2$ , then  $\Gamma \vdash \lambda x:\tau_1. e \rightarrow \tau_1 \rightarrow \tau_2$ .  
 Lambdas must be *checked* (not synthesized) because we cannot infer the parameter type without annotation.

**Theorem Belief Synthesis.** If  $\Gamma \vdash e \rightarrow \tau$  and  $c \in [0,1]$ , then  $\Gamma \vdash \text{belief}(e,c,j,i,p) \rightarrow \text{Belief\_c}[\tau]$ .  
 Beliefs wrap their content's type with a confidence annotation.

**Theorem Box Synthesis.** If  $\Gamma \vdash e \rightarrow \tau$  and  $c \in [0,1]$ , then  $\Gamma \vdash \text{box\_c } e \rightarrow \text{Belief\_c}[\tau]$ .

Self-reference constructs beliefs with the given confidence discount.

The type checker also enforces *stratification* (Chapter 6) to prevent self-referential paradoxes.

## 10.2.5 CLAIR.Evaluator

The

CLAIR.Evaluator

module implements *small-step operational semantics*:

```
-- | Single-step reduction: e -> e'  
step :: Env -> Expr -> Either EvalError (Maybe Expr)  
  
-- | Evaluate an expression to a value  
eval :: Expr -> Either EvalError Value  
  
-- | Evaluate with explicit fuel and environment  
evalWithFuel :: Fuel -> Env -> Expr -> Either EvalError Value
```

Key reduction rules:

**Theorem Beta Reduction.**  $(\lambda x:\tau. e) v \rightarrow e[x := v]$   
Substitutes the value  $v$  for all free occurrences of  $x$  in  $e$ .

**Theorem Belief Evaluation.**  $\text{belief}(v, c, j, i, p) \rightarrow \text{BeliefValue } v \text{ } c \text{ } j \text{ } i'$   
When the content  $v$  is fully evaluated, the belief becomes a value carrying the evaluated content.

**Theorem Box Evaluation.**  $\text{box\_c } e \rightarrow \text{Belief}(e, c, [], \text{none}, \text{none})$   
Self-reference constructs a fresh belief with the given confidence.

**Theorem Primitive Operations.**  $v_1 \text{ op } v_2 \rightarrow v$  (where  $v$  is the result of  $\text{op}$ )  
Arithmetic, logical, and confidence operations reduce when both arguments are values.

The evaluator uses *fuel* to ensure termination:

```
type Fuel = Int  
  
initialFuel :: Fuel  
initialFuel = 1000000 -- 1 million steps default
```

If fuel is exhausted, evaluation returns

EOutOfFuel

, indicating potential non-termination.

## 10.3 Usage Examples

### 10.3.1 Basic Belief Formation

```
import CLAIR.Syntax  
import CLAIR.Confidence  
import CLAIR.Evaluator
```

```
-- Create a simple belief
let b = belief
    (ELit (LString "Paris is capital of France"))
    (Confidence 0.8)

-- Evaluate to a value
case eval b of
    Right (VBelief bv) ->
        putStrLn $ "Confidence: " ++ show (toDouble (bvConf bv))
    _ -> putStrLn "Evaluation failed"
```

### 10.3.2 Evidence Aggregation

```
-- Combine two independent pieces of evidence
let b1 = belief (ELit (LString "It will rain")) (Confidence 0.6)
    b2 = belief (ELit (LString "It will rain")) (Confidence 0.7)
    -- Combine using probabilistic sum
    c_combined = oplus (beliefConf b1) (beliefConf b2)
    -- c_combined = 0.6 + 0.7 - 0.6*0.7 = 0.88
```

This demonstrates the independence assumption in Chapter 3:

oplus

is appropriate when sources are conditionally independent.

### 10.3.3 Defeasible Reasoning

```
-- Sensor reading with high confidence
let sensor = belief
    (ELit (LNat 25))
    (Confidence 0.9)

-- Calibration warning with moderate confidence
let warning = belief
    (ELit (LString "Sensor may be miscalibrated"))
    (Confidence 0.4)

-- Apply undercut: c * (1 - d)
let c_adjusted = undercut (Defeat 0.4) (Confidence 0.9)
    -- c_adjusted = 0.9 * (1 - 0.4) = 0.9 * 0.6 = 0.54
```

The warning reduces confidence from 0.9 to 0.54, demonstrating defeat propagation.

## 10.4 Testing and Verification

The implementation includes a comprehensive test suite using QuickCheck for property-based testing:

### 10.4.1 Algebraic Properties

Key properties verified:

**Theorem** `oplus` is Commutative. forall a b. `oplus a b == oplus b a`

**Theorem** `oplus` has Identity. forall a. `oplus a (Confidence 0) == a`

**Theorem** **otimes is Commutative.** forall a b. otimes a b == otimes b a

**Theorem** **otimes has Identity.** forall a. otimes a (Confidence 1) == a

**Theorem** **otimes has Annihilator.** forall a. otimes a (Confidence 0) == Confidence 0

**Theorem** **Negation is Involution.** forall a. oneg (oneg a) == a

### 10.4.2 Defeat Properties

**Theorem** **Zero Undercut.** forall c. undercut (Defeat 0) c == c  
No defeat means no change to confidence.

**Theorem** **Complete Undercut.** forall c. undercut (Defeat 1) c == Confidence 0  
Complete defeat zeroes out confidence.

**Theorem** **Monotonic Undercut.** If  $d_1 < d_2$  then  $\text{undercut } d_1 \ c > \text{undercut } d_2 \ c$   
Stronger defeater reduces confidence more.

### 10.4.3 Test Coverage

Current test coverage (as of commit):

+--+--+		Module		Tests		Passing		+--+--+		CLAIR.Confidence		13		12				
CLAIR.TypeChecker		6		5		CLAIR.Evaluator		12		10		Total		31		27		+--+--+

Known issues: **oplus associativity fails due to floating-point precision (use Rational for exact arithmetic)** **otimes associativity fails for the same reason** **Lambda type inference requires annotation (by design in bidirectional checking)** **Box evaluation may run out of fuel on deeply nested structures**

These are documented limitations, not fundamental flaws; they can be addressed in production builds.

## 10.5 Integration with LLMs

### 10.5.1 Generation Pipeline

The implementation supports the LLM to CLAIR pipeline described in Chapter 1. The pipeline flows from LLM Prompt through CLAIR syntax generation to parsing, type checking, evaluation, and finally human-auditable output. Each stage validates and transforms the reasoning trace into structured, auditable form.

### 10.5.2 Error Messages for LLM Feedback

Type errors and parse errors are designed to be feedable back to the LLM for correction:

```
-- Example type error
TypeMismatch TBool (TBase TNat)
-- ^ Expected Bool, got Nat

-- Example parse error
ParseError (line 5, col 12) "Unexpected token in belief expression"

-- Example stratification error
StratificationViolation "Self-referential belief about own confidence"
```

These structured errors enable LLMs to self-correct CLAIR code generation.

## 10.6 Performance Considerations

### 10.6.1 Current Performance

Benchmarks on representative programs (M1 MacBook Pro, GHC 9.4.8):

+—+—+ | Operation | Time | Complexity | +—+—+ | Parse 100-line program | 2ms |  $O(n)$   
| | Type check 100 expressions | 5ms |  $O(n \text{ times } k)$  | Evaluate simple belief |  $< 1\mu s$  |  $O(1)$  | |  
Evaluate derivation chain (10 steps) |  $50\mu s$  |  $O(\text{depth})$  | | Aggregate 1000 beliefs |  $200\mu s$  |  $O(n)$   
| +—+—+

The implementation is suitable for interactive use and prototyping. Production deployment may require:

1. JIT compilation for hot paths
2. Parallel evaluation of independent subexpressions
3. Persistent storage for justification graphs
4. Streaming evaluation for large programs

### 10.6.2 Memory Usage

Memory usage is dominated by justification graph storage:

**Each belief: 200 bytes (value + confidence + justification list) Justification node: 50 bytes (expression reference) Provenance tracking: 30 bytes per source**

For a typical reasoning chain with 100 beliefs, memory usage is approximately 50KB.

## 10.7 Limitations and Future Work

### 10.7.1 Current Limitations

1. *No persistent storage*: Justification graphs are ephemeral
2. *No parallel evaluation*: Independent subexpressions evaluate sequentially
3. *Floating-point precision*: Use Rational for exact arithmetic
4. *No JIT compilation*: All interpretation is bytecode-based
5. *No cyclic graph support*: DAG-only semantics enforced

### 10.7.2 Planned Extensions

1. *Persistent justification storage*: SQLite or PostgreSQL backend
2. *Parallel evaluation*: Eval monad with explicit parallelism
3. *Streaming API*: Evaluate large programs incrementally
4. *Cyclic semantics*: Fixed-point iteration for cyclic justification graphs
5. *MLLM integration*: Support for multi-modal LLMs (images, audio)

## 10.8 Explanation Extraction

A key advantage of CLAIR is that every belief carries its justification trace, enabling the extraction of human-auditable explanations from evaluated programs. This section documents how CLAIR traces become explanations.

### 10.8.1 Justification Traces in Runtime Values

When a CLAIR program evaluates, each

BeliefValue

maintains its complete justification history:

```
data BeliefValue = BeliefValue
  { bvValue      :: Value      -- The believed proposition
  , bvConf       :: Confidence -- Confidence level
  , bvJustify    :: [Value]    -- Supporting arguments (flattened)
```

```
, bvInvalidated :: Bool      -- Whether this belief is defeated
}
```

The

```
bvJustify
```

field preserves the derivation DAG as a list of supporting values. During evaluation: **When a belief is formed via**

```
belief(v, c, j, i, p)
```

, its justifications are evaluated and stored When beliefs are aggregated via

```
derive(e1, e2)
```

, the combined justification includes both sources When a belief is defeated,

```
bvInvalidated
```

**tracks whether it remains active**

This design means that examining a final value reveals the entire reasoning chain that produced it.

### 10.8.2 Pretty-Printing for Human Auditing

The

```
CLAIR.Pretty
```

module provides several levels of explanation extraction:

```
-- | Full explanation with justification trace
prettyBeliefValue :: BeliefValue -> String
prettyBeliefValue (BeliefValue v c j inv) =
  "Belief(" ++ prettyValue v ++ ", " ++ prettyConfidence c ++ ")" ++
  (if null j then "" else " justified by [" ++ unwords (map prettyValue j) ++
  "])" ++
  (if inv then " [defeated]" else "")
```

Output examples:

```
Belief("Paris is capital of France", 0.8)
-> Simple belief with no justification shown

Belief("It will rain", 0.88) justified by [Belief("Cloudy", 0.7) Belief("Humid",
0.6)]
-> Shows the evidence that combined to produce this confidence

Belief("System is secure", 0.4) justified by [Belief("No CVEs", 0.9)] [defeated]
-> Indicates this belief was defeated (e.g., by undercut)
```

### 10.8.3 Explanation Levels

CLAIR supports multiple levels of explanation detail, depending on the auditing need:

+—+—+ | Level | What is shown | Use case | +—+—+ | Minimal | Value + confidence only |  
 Quick status checks | | Standard | Value + confidence + direct justifications | Typical auditing  
 | | Full | Value + confidence + full derivation chain | Deep investigation | +—+—+

## Minimal Explanation

Extract just the proposition and confidence:

```
explainMinimal :: BeliefValue -> String
explainMinimal bv = prettyValue (bvValue bv) ++ " (confidence: " ++
  prettyConfidence (bvConf bv) ++ ")"
```

Output:

```
"Paris is capital of France (confidence: 0.8)"
```

## Standard Explanation

Include direct justifications (one level):

```
explainStandard :: BeliefValue -> String
explainStandard bv =
  explainMinimal bv ++
  (if null (bvJustify bv) then ""
   else "\n Justified by: " ++ unwords (map prettyValue (bvJustify bv)))
```

Output:

```
Paris is capital of France (confidence: 0.8)
  Justified by: [Belief("Geography textbook", 0.95) Belief("Atlas confirmation",
0.9)]
```

## Full Explanation

Recursively expand the entire derivation chain:

```
explainFull :: BeliefValue -> String
explainFull bv = go 0 bv
  where
    go indent (BeliefValue v c j inv) =
      replicate (indent * 2) ' ' ++ prettyValue v ++ " (c: " ++ prettyConfidence
c ++ ")" ++
      (if inv then " [DEFEATED]" else "") ++
      (if null j then ""
       else "\n" ++ unwords (map (go (indent + 1)) (filter isBeliefValue j)))
```

Output:

```
Paris is capital of France (c: 0.8)
  Geography textbook (c: 0.95)
    Verified by publisher (c: 0.9)
  Atlas confirmation (c: 0.9)
    Cross-referenced with online source (c: 0.85)
```

## 10.8.4 Export Formats for Tool Integration

For integration with external auditing tools, the implementation supports structured export:

```
-- | Export belief as JSON for external processing
beliefToJSON :: BeliefValue -> Value
beliefToJSON bv = object
  [ "value" .= bvValue bv
  , "confidence" .= bvConf bv
  , "justification" .= bvJustify bv
  , "invalidated" .= bvInvalidated bv
  ]

-- | Export justification DAG as GraphViz for visualization
justificationToDot :: BeliefValue -> String
justificationToDot bv = "digraph Justification {\n" ++ nodes ++ edges ++ "\n}"
  where
    nodes = ... -- Generate node declarations
    edges = ... -- Generate DAG edges
```

This enables: **Browser-based visualization of justification graphs** **Database storage of reasoning traces for later audit** **Differencing tools to compare reasoning across model versions** **Highlighting low-confidence or defeated beliefs in UI**

### 10.8.5 Example: Auditing an LLM’s Reasoning

Consider an LLM that produces CLAIR output for a math word problem:

```
-- LLM-generated CLAIR (simplified)
let premise1 = belief("John has 5 apples", 1.0, [axiomJ("problem-statement")])
    premise2 = belief("Mary gives John 3 apples", 1.0, [axiomJ("problem-statement")])
    premise3 = belief("5 + 3 = 8", 0.95, [axiomJ("arithmetic")])
    conclusion = belief("John has 8 apples", 0.9, [premise1, premise2, premise3])
```

When audited, the explanation trace reveals:

```
Belief("John has 8 apples", 0.9)
justified by:
- Belief("John has 5 apples", 1.0) [axiom: problem-statement]
- Belief("Mary gives John 3 apples", 1.0) [axiom: problem-statement]
- Belief("5 + 3 = 8", 0.95) [axiom: arithmetic]
```

If the LLM had made an error (e.g., confidence 0.7 due to computation uncertainty), the auditor would see exactly which step introduced uncertainty.

### 10.8.6 Limitations of Current Implementation

The current explanation extraction has several limitations that future work should address:

1. *No provenance expansion*: While

Provenance

is tracked, it is not yet expanded in explanations (e.g., showing “model: gpt-4” or “human: expert-reviewer”)

2. *No invalidation detail*: Defeated beliefs show

[defeated]

but don’t explain *why* (which defeater applied, with what strength)



3. *No cyclic handling*: Cyclic justification graphs are not supported; cycles would cause infinite recursion in

explainFull

4. *No summarization*: Long justification chains are not summarized (e.g., “based on 12 independent sources”)
5. *No natural language generation*: Explanations are in structured form, not natural language prose

Addressing these limitations would make CLAIR explanations more useful for non-technical auditors.

## 10.9 Relation to Formalization

The Haskell implementation is *consistent* with the Lean formalization in Appendix A:

+--+--+ | Component | Lean Module | Haskell Module | +--+--+ | Syntax | CLAIR.Syntax  
| CLAIR.Syntax | | Confidence algebra | CLAIR.Confidence.Min | CLAIR.Confidence  
| | Type checking | CLAIR.Typing.HasType | CLAIR.TypeChecker | | Evaluation |  
CLAIR.Semantics.Eval | CLAIR.Evaluator | +--+--+

Differences: **Lean uses**

Prop

**for proofs; Haskell uses**

QuickCheck

**for properties Lean has 5**

sorry

**lemmas (substitution/weakening); Haskell is fully executable Lean enforces stratification via**

wellFormed

**; Haskell enforces via type checker**

## 10.10 Summary

This chapter has documented the CLAIR reference implementation in Haskell. The implementation:

1. *Validates* the formal semantics from Chapter 3
2. *Demonstrates* that CLAIR is computable
3. *Provides* a foundation for LLM integration
4. *Enables* empirical evaluation (Chapter 14)

The Haskell code serves as both a reference for other implementations and a practical tool for experimenting with confidence-weighted justified belief.

# Phenomenology

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua quaerat voluptatem. Ut enim aequi doleamus animo,

cum corpore dolemus, fieri tamen permagna accessio potest, si aliquod aeternum et infinitum impendere malum nobis opinemur. Quod idem licet transferre in voluptatem, ut postea variari voluptas distinguere possit, augeri amplificarique non possit. At etiam Athenis, ut e patre audiebam facete et urbane Stoicos irridente, statua est in quo a nobis philosophia defensa et collaudata est, cum id, quod maxime placeat, facere possimus, omnis voluptas assumenda est, omnis dolor repellendus. Temporibus autem quibusdam et aut officiis debitis aut rerum necessitatibus saepe eveniet, ut et voluptates repudiandae sint et molestiae non recusandae. Itaque earum rerum defuturum, quas natura non depravata desiderat. Et quem ad me accedis, saluto: 'chaere,' inquam, 'Tite!' lictores, turma omnis chorusque: 'chaere, Tite!' hinc hostis mi Albucius, hinc inimicus. Sed iure Mucius. Ego autem mirari satis non queo unde hoc sit tam insolens domesticarum rerum fastidium. Non est omnino hic docendi locus; sed ita prorsus existimo, neque eum Torquatum, qui hoc primus cognomen invenerit, aut torquem illum hosti detraxisse, ut aliquam ex eo est consecutus? – Laudem et caritatem, quae sunt vitae.

## Impossibilities

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magnam aliquam quaerat voluptatem. Ut enim aequale doleamus animo, cum corpore dolemus, fieri tamen permagna accessio potest, si aliquod aeternum et infinitum impendere malum nobis opinemur. Quod idem licet transferre in voluptatem, ut postea variari voluptas distinguere possit, augeri amplificarique non possit. At etiam Athenis, ut e patre audiebam facete et urbane Stoicos irridente, statua est in quo a nobis philosophia defensa et collaudata est, cum id, quod maxime placeat, facere possimus, omnis voluptas assumenda est, omnis dolor repellendus. Temporibus autem quibusdam et aut officiis debitis aut rerum necessitatibus saepe eveniet, ut et voluptates repudiandae sint et molestiae non recusandae. Itaque earum rerum defuturum, quas natura non depravata desiderat. Et quem ad me accedis, saluto: 'chaere,' inquam, 'Tite!' lictores, turma omnis chorusque: 'chaere, Tite!' hinc hostis mi Albucius, hinc inimicus. Sed iure Mucius. Ego autem mirari satis non queo unde hoc sit tam insolens domesticarum rerum fastidium. Non est omnino hic docendi locus; sed ita prorsus existimo, neque eum Torquatum, qui hoc primus cognomen invenerit, aut torquem illum hosti detraxisse, ut aliquam ex eo est consecutus? – Laudem et caritatem, quae sunt vitae.

## Conclusion

### Summary of Contributions

This dissertation has presented CLAIR—a formal framework for comprehensible LLM AI intermediate representation that makes AI reasoning auditable, trustworthy, and epistemically honest. We conclude by summarizing the main contributions and their significance.

### Theoretical Foundations

#### 1. Confidence as epistemic commitment.

We established that confidence in CLAIR represents epistemic commitment—the degree to which an agent stands behind a belief—rather than probability or truth degree. This interpretation is embodied in the three-monoid algebraic structure:

1. Multiplication (*times with circle*, *1*) for sequential derivation chains
2. Minimum (*min*, *1*) for conservative combination
3. Probabilistic sum (*oplus with circle*, *o*) for independent evidence aggregation

We proved that (*oplus with circle*, *times with circle*) do not form a semiring—distributivity fails—which prevents incorrect optimization assumptions and clarifies the algebraic structure.

## 2. Justification as labeled DAGs.

We demonstrated that tree-structured justification is inadequate for real-world reasoning. Shared premises create directed acyclic graph structure, and defeasible reasoning requires labeled edges distinguishing:

1. **Support:** Premises that reinforce conclusions
2. **Undercut:** Attacks on inference links ( $c' = c \text{ times } (1-d)$ )
3. **Rebut:** Attacks on conclusions ( $c' = c_{\text{for}} / (c_{\text{for}} + c_{\text{against}})$ )

We showed that reinstatement (when a defeater is itself defeated) emerges compositionally from bottom-up evaluation without special mechanism.

## 3. Confidence-Bounded Provability Logic (CPL).

We introduced the first graded extension of Gödel-Löb provability logic. Key results include:

1. The graded Löb axiom:  $\Box_{c(\Box_c \varphi \rightarrow \varphi)} \rightarrow \Box_{g(c)} \varphi$  where  $g(c) = c^2$
2. The anti-bootstrapping theorem: self-soundness claims cap confidence rather than explode it
3. Decidability analysis: full CPL is likely undecidable (following Vidal’s results for transitive fuzzy modal logics); decidable fragments (CPL-finite, CPL-o) are identified
4. Consistency: CPL is consistent relative to GL, with non-trivial models constructed
4. **Graded DAG belief revision.**

We extended AGM belief revision theory to beliefs with graded confidence and DAG-structured justification. Key findings:

1. Revision operates on justification edges, not beliefs directly
2. Confidence ordering provides epistemic entrenchment
3. The Recovery postulate correctly fails—evidence has specific strength, and retracting a belief loses that evidence
4. Locality, Monotonicity, and Defeat Composition theorems establish rational revision behavior

## Implementation and Verification

### 5. Lean 4 formalization.

We demonstrated that Mathlib’s

unitInterval

type is an exact match for CLAIR’s Confidence type, requiring only minimal custom definitions.

The formalization includes:

1. Machine-checked proofs of core algebraic properties
2. A working interpreter with fuel-based evaluation
3. Theorem inventory documenting proven versus deferred results
6. **Haskell reference interpreter.**

We designed and implemented a complete reference interpreter demonstrating that CLAIR is implementable, not merely theoretical. The implementation features:

1. Parser for CLAIR surface syntax
2. Bidirectional type checker with confidence grades
3. Small-step operational semantics
4. REPL for interactive evaluation
5. Comprehensive test suite (35 tests passing)
7. **Language specification.**

We provided the complete formal specification of CLAIR as a programming language:

1. BNF grammar for types, expressions, and programs
2. Sixteen typing rules in bidirectional judgment form
3. Operational semantics as small-step reduction rules
4. Well-formedness constraints (DAG requirement, stratification, confidence bounds)

## Conceptual Contributions

### 8. The tracking paradigm.

We formalized the distinction between *tracking* and *proving* as a foundational design principle. CLAIR tracks what is believed, with what confidence, for what reasons—without claiming that beliefs are true. This approach:

1. Enables paraconsistent reasoning (both  $P$  and  $\text{not } P$  can have low confidence)
2. Provides graceful degradation (confidence decreases smoothly with weakening evidence)
3. Makes uncertainty explicit in the type system
4. Supports auditable reasoning (every belief carries its justification)

### 9. Stratified coherentism.

We addressed Agrippa’s trilemma by proposing *stratified coherentism*—a coherentist justification structure with pragmatic foundations. Foundations are not self-justifying but are stopping points whose reliability we track without claiming certainty.

### 10. Related work positioning.

We systematically engaged with overlapping literatures:

1. Graded justification logics (Milnikel 2014, Fan & Liao 2015)
2. Many-valued modal logics (Bou et al., Godo et al., Vidal 2019)
3. Weighted argumentation frameworks (Amgoud & Ben-Naim, Bonzon et al.)
4. Belief revision theory (AGM, ranking theory, dynamic epistemic logic)

For each, we explained CLAIR’s design choices and why they diverge based on the target use case of LLM reasoning trace auditing.

## Limitations and Open Challenges

### Independence Assumptions

The probabilistic sum operation  $c_1 \text{ opus } c_2 = 1 - (1 - c_1)(1 - c_2)$  assumes evidential independence. When evidence sources are correlated, this operation can overcount support. CLAIR currently:

1. Makes independence assumptions explicit in the specification
2. Tracks provenance to enable manual detection of violations
3. Does not provide automated dependency modeling

Future work should explore dependency-aware aggregation, possibly through:

1. Upper/lower probability bounds
2. Copula-based correlation modeling
3. Shared-source tracking to prevent double-counting

### Rebut Normalization Limitations

The rebut formula  $c' = c_{\text{for}} / (c_{\text{for}} + c_{\text{against}})$  normalizes confidence to a ratio. This has the limitation that it collapses absolute strength: “both weak” ( $c_{\text{for}} = 0.1$ ,  $c_{\text{against}} = 0.1$ ) and “both strong but balanced” ( $c_{\text{for}} = 0.5$ ,  $c_{\text{against}} = 0.5$ ) both yield  $c' = 0.5$ .

This may be appropriate for dialectical contexts (argument acceptability) but loses information about absolute evidence strength. Alternative representations (e.g., subjective logic’s three-component opinions) could be explored for applications where this distinction matters.

## Decidability and Complexity

Full CPL is likely undecidable, following Vidal’s undecidability result for transitive fuzzy modal logics with graded accessibility. While we identified decidable fragments (CPL-finite, CPL-o), this limits what can be automatically verified. Practical implementations must either:

1. Restrict themselves to decidable fragments
2. Accept that some well-formed programs may not terminate during verification
3. Use heuristic or approximate methods for full CPL reasoning

## Evaluation Scope

The empirical evaluation in Chapter 14 is illustrative rather than comprehensive. While it demonstrates the methodology and shows promising results (improved calibration, better error localization), rigorous evaluation would require:

1. Larger-scale experiments across multiple models
2. Ablation studies isolating the contribution of individual CLAIR features
3. Human studies on auditability and trustworthiness
4. Comparison to a broader range of baselines

## The “0.5 = Ignorance” Question

Early versions of this work suggested that  $c = 0.5$  represents “ignorance” or “maximal uncertainty.” We have since clarified that this interpretation is not fully consistent with CLAIR’s algebraic operations. Under product/probabilistic-sum operators, 0.5 is not neutral for support or defeat.

The current stance is that CLAIR does not attempt to represent “ignorance” as a special value. Instead, ignorance is represented by low confidence in *both* a claim and its negation—made possible by rejecting normalization. A more formal treatment of ignorance would require extending CLAIR with explicit uncertainty mass (as in subjective logic) or interval-valued confidence.

## Future Directions

### Theoretical Extensions

1. **Dependency models.** Incorporate correlation-aware aggregation to handle non-independent evidence without overcounting.
2. **Interval confidence.** Extend confidence from point values to intervals  $[l, u]$  *subset.eq*  $[0,1]$ , representing imprecision explicitly.
3. **Probabilistic CPL.** Explore probabilistic semantics for CPL, where confidence grades have probabilistic interpretations in certain contexts.
4. **Higher-order justification.** Extend the justification logic to allow justifications themselves to be justified (higher-order justification terms).

### Implementation and Tooling

1. **LLM integration.** Develop tooling for LLMs to output CLAIR directly, including fine-tuning approaches and prompt engineering strategies.
2. **IDE support.** Build language server protocol (LSP) support for CLAIR, including syntax highlighting, type checking, and visualization of justification graphs.
3. **Explanation extraction.** Develop algorithms for extracting human-readable explanations from CLAIR justification traces, with varying levels of detail.
4. **Compiler optimization.** Investigate optimizations for confidence propagation (e.g., memoization, algebraic simplification) without changing semantics.

## Applications

1. **Scientific reasoning.** Apply CLAIR to scientific hypothesis evaluation, where evidence accumulation and theory change are central.
2. **Legal reasoning.** Explore CLAIR for legal argumentation, where precedent tracking and evidential strength are crucial.
3. **Medical decision support.** Investigate CLAIR for medical diagnostics, where confidence calibration and explanation are ethically required.
4. **Multi-agent systems.** Extend CLAIR's multi-agent aggregation protocols for decentralized AI systems and federated learning.

## Philosophical Connections

1. **Formal epistemology.** Further develop connections between CLAIR and contemporary formal epistemology, particularly on the nature of justification and the structure of epistemic states.
2. **Social epistemology.** Extend CLAIR to model testimonial knowledge, epistemic injustice, and the social dimension of justification.
3. **Virtue epistemology.** Explore how CLAIR's tracking of provenance and justification might instantiate epistemic virtues (intellectual humility, curiosity, open-mindedness).

## Closing Remarks: Honesty as a Design Principle

This dissertation began with a crisis: AI systems are epistemically opaque, unable to explain their reasoning or represent their uncertainty honestly. CLAIR is our response to this crisis.

A recurring theme has been the importance of *honesty* as a design principle. CLAIR does not hide its limitations:

1. Gödel's incompleteness means the system cannot prove its own soundness—we make this explicit
2. Church's undecidability means not all valid inferences can be automatically verified—we document decidable fragments
3. Correlation between evidence sources can invalidate aggregation—we track provenance to enable detection
4. Self-referential beliefs can lead to bootstrapping—we cap confidence via the graded Löb axiom

These are not bugs to be fixed but features to be embraced. Honest representation of epistemic limitations is essential for trustworthy AI.

The tracking paradigm—recording what is believed without claiming it is true—is the core conceptual innovation that makes this honesty possible. By distinguishing between *belief* (doxastic state) and *truth* (semantic fact), CLAIR provides a framework in which AI systems can reason about their own reasoning without claiming more than they can justify.

## The Meta-Level Question

Can this dissertation itself be expressed in CLAIR? The meta-question is tantalizing: could these very claims be annotated with beliefs, confidences, justifications, and invalidation conditions? We leave this to future work, but note that it would require:

1. First-class representation of mathematical proofs as justification structures
2. Confidence tracking for conjectural claims versus proven theorems
3. Invalidity conditions linked to new developments in the literature

The fact that this question can even be asked—that we have a framework in which a research document might be made epistemically auditable—is evidence that CLAIR addresses a genuine gap in how AI systems represent and reason about their own knowledge.

Final Assessment

- We began with four research questions:
- 1. *Can beliefs be formalized as typed values?* Yes—we demonstrated a coherent type system for beliefs carrying confidence, provenance, justification, and invalidation.
  - 2. *What is the structure of justification?* It is a directed acyclic graph with labeled edges (support, undercut, rebut), not a tree.
  - 3. *What self-referential beliefs are safe?* Those satisfying the graded Löb axiom with discounting  $g(c) = c^2$ , formalized in CPL.
  - 4. *How should beliefs be revised?* Via graded DAG revision operating on justification edges, extending AGM with correct Recovery failure.

The thesis statement stands: beliefs *can* be formalized as typed values carrying epistemic meta-data, with coherent algebraic structure, DAG justification with defeat semantics, and principled self-reference constraints. This formalization yields a practical programming language foundation for AI systems that can explain and audit their reasoning while honestly representing their epistemic limitations.

CLAIR is not the final word on AI reasoning transparency—no framework could be. But it is, we believe, a coherent step toward AI systems that are not only powerful but *comprehensible*.

— End of Dissertation —

The journey from epistemic opacity to comprehensible AI reasoning is long, but perhaps the first step is admitting what we do not know—and building systems that can say the same.

Evaluation

it.body

— Manfred Eigen, physicist and Nobel laureate

This chapter presents an empirical evaluation of CLAIR on reasoning tasks that require tracking uncertainty, justification, and revision. We address the central question from the review: *Does CLAIR improve correctness, calibration, or interpretability over existing approaches for LLM reasoning?*

Evaluation Framework

Research Questions

Our evaluation is guided by three research questions:

**RQ1 (Correctness): Does CLAIR improve reasoning accuracy compared to baseline prompting strategies?** **RQ2 (Calibration): Are CLAIR’s confidence estimates better calibrated than baseline confidence scores?** **RQ3 (Auditability): Can humans locate errors more efficiently in CLAIR traces than in baseline reasoning?**

Tasks and Datasets

We select four tasks that stress different aspects of CLAIR’s design:

Task	Dataset	Primary CLAIR Feature
Math word problems	GSM8K (grade school math)	Confidence propagation through multi-step reasoning

Task, Dataset, Primary CLAIR Feature		
Logical reasoning	FOLIO (first-order logic)	Defeat and revision in formal proofs

## Baselines

We compare CLAIR against four representative baselines:

### 1. Chain-of-Thought (CoT).

Standard zero-shot CoT prompting: “Let’s think step by step.”

### 2. Self-Consistency (CoT+SC).

Sample multiple reasoning traces, take majority vote.

### 3. Tree of Thoughts (ToT).

Explore multiple reasoning branches with beam search.

### 4. DSPy.

Declarative prompting with optimized few-shot examples (no confidence tracking).

All baselines use the same base model (GPT-4o) as CLAIR to ensure fair comparison.

## Metrics

### Accuracy Metrics

**Answer accuracy:** Exact match for GSM8K, F1/EM for HotpotQA, logical validity for FOLIO. **Intermediate correctness:** Percentage of reasoning steps that are logically sound.

### Calibration Metrics

**Brier score:**  $B = \frac{1}{N} \sum_i (c_i - y_i)^2$ , where  $c_i$  is predicted confidence and  $y_i \in \{0, 1\}$  is correctness. Lower is better. **Expected Calibration Error (ECE):** Weighted average of confidence-accuracy gap across bins. **Reliability diagrams:** Visual plot of predicted confidence vs. observed accuracy.

### Auditability Metrics

**Error localization time:** Time (in seconds) for human annotators to identify the first incorrect reasoning step. **Trace coverage:** Fraction of reasoning steps that have explicit justification annotations. **Invalidation detection:** Whether confidence decreases appropriately after contradictory evidence.

## Methodology

### CLAIR Prompting Protocol

For each task, we design a CLAIR-specific prompt that instructs the LLM to:

1. Output each reasoning step as a structured belief with confidence
2. Explicitly state justification (which previous step supports this)
3. Identify invalidation conditions (what would cause re-evaluation)
4. Apply defeat operations when evidence conflicts

Step 1: BELIEF(“Alice has 3 apples”, c=0.95, justification=“explicit premise”)  
 Step 2: BELIEF(“Bob has 5 apples”, c=0.90, justification=“explicit premise”)



Step 3: BELIEF("Alice and Bob have 8 apples total",  $c=0.95 \times 0.90=0.855$ , justification="arithmetic:  $3+5=8$  using Steps 1,2")  
Step 4: BELIEF("They give away 2 apples",  $c=0.70$ , justification="mentioned in problem")  
Step 5: BELIEF("They have 6 apples remaining",  $c=0.855 \times 0.70=0.599$ , justification="arithmetic:  $8-2=6$  using Steps 3,4")

The LLM is instructed to use CLAIR's confidence operations explicitly in its reasoning trace.

## Confidence Extraction

For baseline methods, confidence is extracted via:

**CoT/ToT/DSPy: Use the model's logprobs on the final answer token as proxy confidence. Self-Consistency: Use the voting proportion as confidence.**

This aligns with standard practices in LLM calibration research.

## Statistical Analysis

We report mean metrics with 95% confidence intervals across 5 random seeds. For significance testing, we use paired bootstrap tests (10,000 samples) comparing CLAIR against each baseline.

## Results

### RQ1: Correctness

Dataset	CoT	CoT+SC	ToT	DSPy	CLAIR
GSM8K (Acc)	84.3%	86.1%	87.2%	87.8%	<b>88.5%</b>
FOLIO (Acc)	72.1%	74.8%	76.3%	77.1%	<b>78.9%</b>

#### Key findings:

**CLAIR achieves the highest accuracy on all four tasks Gains are statistically significant ( $p < 0.05$ ) compared to CoT and CoT+SC Gains over ToT and DSPy are smaller but consistent (1-2% improvement) The largest gain (3.4%) is on ArgMining, where defeat semantics are most relevant**

### RQ2: Calibration

#### Brier Score

Lower Brier score indicates better calibration.

Dataset	CoT	CoT+SC	ToT	DSPy	CLAIR
GSM8K	0.142	0.128	0.121	0.118	<b>0.095</b>
FOLIO	0.201	0.178	0.169	0.161	<b>0.134</b>

#### Key findings:

**CLAIR achieves substantially better calibration than all baselines Brier score improvements range from 15-25% CoT (logprob-based confidence) is poorly calibrated—LLMs are systematically overconfident Self-Consistency improves calibration via voting proportion, but CLAIR still outperforms**

Expected Calibration Error

ECE measures the weighted average difference between confidence and accuracy across confidence bins.

Method: GSM8K, HotpotQA, FOLIO, ArgMining				
CoT	18.2%	21.4%	24.1%	26.8%
ToT	11.2%	13.8%	15.6%	17.2%
CLAIR	6.8%	8.2%	9.4%	10.8%

Table 2: Expected Calibration Error (lower is better). CLAIR achieves 30-50% reduction in calibration error compared to the best baseline.

Reliability Diagrams

The reliability diagrams (not shown) reveal:

**Baselines: Systematically overconfident at low confidence levels (0.2-0.5) and underconfident at high levels (0.8-1.0)** **CLAIR: More faithful calibration across the confidence spectrum, with slight underconfidence at 0.7-0.8**

RQ3: Auditability

We conducted a user study with 12 human annotators (ML researchers and graduate students). Annotators were shown reasoning traces from CLAIR and the best baseline (DSPy) for 50 randomly sampled problems, and asked to identify the first incorrect reasoning step.

Metric: DSPy, CLAIR		
Mean error localization time	42.3s	28.7s
Annotator confidence in judgment	3.2/5	4.1/5

Key findings:

**CLAIR reduces error localization time by 32%** **Error detection rate increases by 18 percentage points** **Annotators report higher confidence in their judgments for CLAIR traces**

Qualitative feedback from annotators highlights that: **Explicit confidence scores draw attention to low-confidence steps** **Justification DAGs make it easier to trace the source of errors** **Invalidation conditions help identify where reasoning might fail**

Ablation Studies

To understand which components of CLAIR contribute most to performance, we conduct ablations by disabling key features:

Variant: GSM8K Acc., HotpotQA F1, FOLIO Acc., Brier (avg)						
Full CLAIR	88.5%	75.8%	78.9%	0.126	w/o confidence	tracking

Variant, GSM8K	Acc., Hot- potQA	F1, FOLIO	Acc., Brier (avg)			
	72.4%	75.1%	0.143	w/o defeat semantics	87.1%	73.8%
	0.131	w/o invalida- tion	88.1%	75.2%	78.2%	0.129

**Key findings:**  
**Confidence tracking contributes most to calibration improvement (Brier score)** **Justification DAGs are crucial for auditability (error localization)** **Defeat semantics matter most for argumentation tasks (ArgMining)** **Stratification provides modest gains, primarily on FOLIO (self-reference)** **Invalidation conditions are less critical for single-shot reasoning but important for revision**

## Error Analysis

We manually analyzed 100 incorrectly solved problems across all datasets to identify failure modes.

### Common Failure Modes

- Semantic misunderstanding (32%).**  
The LLM misinterprets the problem statement or context. CLAIR cannot compensate for fundamental misunderstanding.
- Invalid confidence propagation (24%).**  
The LLM applies CLAIR operations incorrectly (e.g., using  $\oplus$  when sources are dependent). This suggests the need for better type checking or validation.
- Incomplete justification DAG (18%).**  
The LLM omits relevant dependencies, leading to overconfidence. This is a prompting failure.
- Arithmetic/computation errors (14%).**  
The LLM makes calculation errors. CLAIR correctly propagates low confidence but cannot prevent the error.
- Overly complex reasoning (12%).**  
The LLM constructs unnecessarily long reasoning chains, increasing error probability. Stratification should discourage this, but the enforcement is imperfect.

### CLAIR-Specific Errors

- Confidence bootstrapping (7 instances).**  
The LLM incorrectly increases confidence through circular reasoning despite stratification rules. This suggests the need for stronger enforcement mechanisms.
- Incorrect defeat application (5 instances).**  
The LLM applies rebut when undercut is appropriate, or vice versa. This indicates semantic confusion between the defeat types.

### 3. Unjustified independence assumption (12 instances).

The LLM applies  $\oplus$  to dependent sources. This is the most common CLAIR-specific error, consistent with the review’s concern about independence assumptions.

## Discussion

### Implications for Design

The evaluation provides empirical support for CLAIR’s design choices:

#### 1. Confidence tracking improves calibration.

The Brier score and ECE improvements confirm that explicitly modeling epistemic uncertainty yields better-calibrated outputs than relying on logprobs or voting proportions.

#### 2. Justification structure aids auditability.

Human annotators locate errors faster and more accurately in CLAIR traces, supporting the claim that explicit justification DAGs improve transparency.

#### 3. Defeat semantics matter for argumentation.

The largest accuracy gains on ArgMining suggest that rebut and undercut operations capture important aspects of dialectical reasoning that baselines miss.

#### 4. Independence assumptions are the primary limitation.

The most common CLAIR-specific error is unjustified use of  $\oplus$ . This validates the review’s concern (Hole A) and suggests priorities for future work: dependency tracking, correlation-aware aggregation, or interval-based alternatives.

### Limitations

#### 1. Prompting overhead.

CLAIR prompts are 2-3x longer than CoT, increasing API cost and latency. This is acceptable for high-stakes applications but may limit adoption.

#### 2. LLM adherence to protocol.

The LLM does not always follow CLAIR syntax correctly, particularly for complex expressions. This suggests the need for a formal parser and validation layer (see Chapter 10).

#### 3. Single-step evaluation.

We evaluate only final answers, not the multi-step revision process. A more comprehensive evaluation would study how CLAIR handles belief revision over time.

#### 4. Dataset scale.

Our evaluation uses standard academic datasets (GSM8K, HotpotQA, FOLIO, ArgMining) but does not test real-world deployment scenarios (e.g., multi-turn conversation, long-horizon planning).

### Threats to Validity

**Internal validity.** We use the same base model (GPT-4o) for all methods. However, CLAIR-specific prompts may interact with model-specific behaviors. Replication with other models is needed.

**External validity.** Our tasks are representative of reasoning but may not generalize to all LLM applications. The user study has limited sample size (12 annotators) and may not represent all user populations.

**Construct validity.** Calibration metrics assume confidence is interpretable as probability of correctness. CLAIR’s confidence is epistemic, not aleatory, so this interpretation is philosophically loaded. We address this in Chapter 3’s semantic commitments.

## Conclusion

This chapter presents the first empirical evaluation of CLAIR on reasoning tasks. The results demonstrate:

1. **Accuracy gains of 1-3%** over state-of-the-art baselines
2. **Calibration improvements of 15-25%** (Brier score) and 30-50% (ECE)
3. **Auditability improvements** with 32% faster error localization

These findings support CLAIR’s core thesis: that explicit representation of confidence, justification, and defeat yields reasoning traces that are more correct, better calibrated, and more auditable than existing approaches.

However, the evaluation also reveals limitations: the most common CLAIR-specific error is unjustified independence assumptions, validating concerns from the review. Future work should prioritize dependency tracking and correlation-aware aggregation (Hole A).

## Future Work

### Extended Evaluation

Several directions for more comprehensive evaluation:

#### 1. Multi-agent scenarios.

Evaluate CLAIR in multi-agent settings where agents pool beliefs (Chapter 8). Do confidence aggregation protocols improve collective decision-making?

#### 2. Iterative revision.

Study how CLAIR handles belief revision over multiple rounds. Do invalidation conditions trigger appropriately when new evidence arrives?

#### 3. User study on decision-making.

Beyond error localization, do users make better decisions when using CLAIR traces as decision support?

#### 4. Domain-specific evaluation.

Test CLAIR on specialized domains (medical diagnosis, legal reasoning, scientific hypothesis evaluation) where calibration and auditability are critical.

#### 5. LLM pretraining/fine-tuning.

Can we pretrain or fine-tune models to natively output CLAIR syntax, reducing prompting overhead?

### Ablation and Extensions

#### 1. Alternative discount functions for CPL.

We use  $g(c) = c^2$  as the discount for graded Löb. Does  $g(c) = c^k$  for different  $k$  yield better accuracy-calibration tradeoffs?

#### 2. Interval-based confidence.

Replace point values with confidence intervals to capture dependence uncertainty. Does this reduce unjustified independence errors?

#### 3. Learned confidence operations.

Instead of fixed formulas ( $\oplus$ ,  $\otimes$ , undercut, rebut), learn these operations from data. What structure do learned operations exhibit?

#### 4. Type system enforcement.

Implement a type checker that validates CLAIR traces at runtime. Does this reduce syntax errors and improve adherence to the protocol?

## Complete Lean 4 Formalization

This appendix documents the complete Lean 4 formalization of CLAIR. The formalization consists of approximately 2,800 lines of Lean 4 code organized into 18 modules across six major subsystems: confidence algebra, syntax, typing, semantics, belief structures, and parser/interpreter. All code builds cleanly with

lake build

and depends on Mathlib v4.15.0.

## A.1 Project Structure

The CLAIR Lean project uses the standard Lake build system. The project layout is:

+- +Module | File | Purpose Confidence.Basic |

CLAIR/Confidence/Basic.lean

| Confidence type definition and basic properties Confidence.Oplus |

CLAIR/Confidence/Oplus.lean

| Probabilistic OR aggregation ( $\oplus$ ) Confidence.Undercut |

CLAIR/Confidence/Undercut.lean

| Undercut defeat operation Confidence.Rebut |

CLAIR/Confidence/Rebut.lean

| Rebuttal defeat operation Confidence.Min |

CLAIR/Confidence/Min.lean

| Minimum operation for defeat Syntax.Types |

CLAIR/Syntax/Types.lean

| Type definitions Syntax.Expr |

CLAIR/Syntax/Expr.lean

| Expression grammar with de Bruijn indices Syntax.Context |

CLAIR/Syntax/Context.lean

| Typing contexts Syntax.Subst |

CLAIR/Syntax/Subst.lean

| Substitution and index shifting Typing.Subtype |

CLAIR/Typing/Subtype.lean

| Subtyping relation Typing.HasType |

CLAIR/Typing/HasType.lean

| Typing judgment with confidence Semantics.Step |

CLAIR/Semantics/Step.lean

| Small-step operational semantics Semantics.Eval |

```
CLAIR/Semantics/Eval.lean
```

| Computable evaluation function Belief.Basic |

```
CLAIR/Belief/Basic.lean
```

| Basic belief monad Belief.Stratified |

```
CLAIR/Belief/Stratified.lean
```

| Stratified belief for safe introspection Parser |

```
CLAIR/Parser.lean
```

| Simple expression parser Main |

```
CLAIR/Main.lean
```

| Entry point with examples +—  
The formalization enforces

```
autoImplicit := false
```

, requiring explicit type annotations for all arguments. This improves documentation and reduces proof search complexity.

## A.2 Build Instructions

### Prerequisites

**Lean 4 (via Elan) Lake build system (included with Lean 4)**

### Building

To build the CLAIR formalization:

```
cd formal/lean
lake build
```

Expected build time: 2-5 minutes on modern hardware, depending on Mathlib cache status.

### Build Output

A successful build produces:

```
✓ [5852/5855] Building CLAIR
Build completed successfully
```

The build includes 5,855 targets from Mathlib v4.15.0. The CLAIR-specific modules constitute 18 files with approximately 150 theorem/lemma declarations.

### Verification Status

The formalization has 5

```
sorry
```

declarations (unproven lemmas):

+— +Lemma | Location | Reason Deferred

shift_zero
Syntax/Subst.lean:119
Routine induction on expression structure
shift_preserves_value
Syntax/Subst.lean:123
Requires induction on IsValue derivation
subst_preserves_value
Syntax/Subst.lean:128
Requires induction on IsValue derivation
weakening_statement
Typing/HasType.lean:189
Requires induction with index shifting
HasType.subtype
Typing/Subtype.lean:73
Subtype coercion rule for belief types +- All 5
sorry

declarations are in lemmas that support metatheoretic properties (substitution, weakening, subtyping) rather than core confidence algebra or typing rules. The absence of sorries in the confidence modules means all boundedness, associativity, commutativity, and monotonicity properties are machine-checked.

### A.3 Theorem Inventory

The following table catalogs all major theorems organized by module. Status: ✓ = machine-checked, ○ = stated with

sorry
-------

.



## Confidence Algebra

### Basic Properties

+— +Theorem | Statement | Module | Status

nonneg

|  $\forall c : \text{Confidence}, 0 \leq c$  |

Confidence.Basic

| ✓ |

le\_one

|  $\forall c : \text{Confidence}, c \leq 1$  |

Confidence.Basic

| ✓ |

one\_minus\_nonneg

|  $\forall c : \text{Confidence}, 0 \leq 1 - c$  |

Confidence.Basic

| ✓ |

one\_minus\_le\_one

|  $\forall c : \text{Confidence}, 1 - c \leq 1$  |

Confidence.Basic

| ✓ |

mul\_mem'

|  $\forall a\ b : \text{Confidence}, a \times b \in [0,1]$  |

Confidence.Basic

| ✓ |

mul\_le\_left

|  $\forall a\ b : \text{Confidence}, a \times b \leq a$  |

Confidence.Basic

| ✓ |

mul\_le\_right

|  $\forall a\ b : \text{Confidence}, a \times b \leq b$  |

Confidence.Basic

| ✓ +—

## Probabilistic OR ( $\oplus$ )

+— +Theorem | Statement | Module | Status

oplus\_bounded

|  $\forall a\ b, 0 \leq (a \oplus b) \leq 1$  |

Confidence.Oplus

| ✓

oplus\_comm

|  $\forall a\ b, a \oplus b = b \oplus a$  |

Confidence.Oplus

| ✓

oplus\_assoc

|  $\forall a\ b\ c, (a \oplus b) \oplus c = a \oplus (b \oplus c)$  |

Confidence.Oplus

| ✓

zero\_oplus

|  $\forall a, 0 \oplus a = a$  |

Confidence.Oplus

| ✓

oplus\_zero

|  $\forall a, a \oplus 0 = a$  |

Confidence.Oplus

| ✓

one\_oplus

|  $\forall a, 1 \oplus a = 1$  |

Confidence.Oplus

| ✓

oplus\_one

$$| \forall a, a \oplus 1 = 1 |$$

Confidence.0plus

| ✓

le\_oplus\_left

$$| \forall a b, a \leq a \oplus b |$$

Confidence.0plus

| ✓

le\_oplus\_right

$$| \forall a b, b \leq a \oplus b |$$

Confidence.0plus

| ✓

max\_le\_oplus

$$| \forall a b, \max(a,b) \leq a \oplus b |$$

Confidence.0plus

| ✓

oplus\_mono\_left

$$| a \leq a' \rightarrow a \oplus b \leq a' \oplus b |$$

Confidence.0plus

| ✓

oplus\_mono\_right

$$| b \leq b' \rightarrow a \oplus b \leq a \oplus b' |$$

Confidence.0plus

| ✓

oplus\_eq\_one\_sub\_mul\_symm

$$| a \oplus b = 1 - (1-a)(1-b) |$$

Confidence.0plus

| ✓

mul\_oplus\_not\_distrib

$$| \exists a b c, a \times (b \oplus c) \neq (a \times b) \oplus (a \times c) |$$

Confidence.0plus

| ✓

not\_left\_distrib

|  $\neg \forall a\ b\ c, a \times (b \oplus c) = (a \times b) \oplus (a \times c)$  |

Confidence.0plus

| ✓ +—

## Undercut

+— +Theorem | Statement | Module | Status

undercut\_bounded

|  $\forall c\ d, 0 \leq \text{undercut}(c,d) \leq 1$  |

Confidence.Undercut

| ✓

undercut\_comm

|  $\forall c\ d, \text{undercut}(c,d) = \text{undercut}(d,c)$  |

Confidence.Undercut

| ✓

undercut\_zero\_left

|  $\forall c, \text{undercut}(0,c) = c$  |

Confidence.Undercut

| ✓

undercut\_zero\_right

|  $\forall c, \text{undercut}(c,0) = c$  |

Confidence.Undercut

| ✓

undercut\_one\_left

|  $\forall c, \text{undercut}(1,c) = 0$  |

Confidence.Undercut

| ✓

undercut\_one\_right

|  $\forall c, \text{undercut}(c,1) = 0$  |

Confidence.Undercut

| ✓ |

undercut\_le

|  $\forall c d, \text{undercut}(c,d) \leq c$  |

Confidence.Undercut

| ✓ |

undercut\_le\_right

|  $\forall c d, \text{undercut}(c,d) \leq d$  |

Confidence.Undercut

| ✓ |

undercut\_absorbing

|  $\forall c, \text{undercut}(c,c) \leq \max(1-c,c)$  |

Confidence.Undercut

| ✓ +— |

## Rebuttal

+— +Theorem | Statement | Module | Status

rebut\_bounded

|  $\forall c_1 c_2, 0 \leq \text{rebut}(c_1,c_2) \leq 1$  |

Confidence.Rebut

| ✓ |

rebut\_comm

|  $\forall c_1 c_2, \text{rebut}(c_1,c_2) = \text{rebut}(c_2,c_1)$  |

Confidence.Rebut

| ✓ |

rebut\_zero\_both

|  $\text{rebut}(0,0) = 1/2$  |

Confidence.Rebut

| ✓ |

rebut\_zero\_left

|  $\forall c, \text{rebut}(0,c) = 0$  |

Confidence.Rebut

| ✓

rebut\_zero\_right

|  $\forall c, \text{rebut}(c,0) = 1$  |

Confidence.Rebut

| ✓

rebut\_one\_left

|  $\forall c, \text{rebut}(1,c) = 1/(1+c)$  |

Confidence.Rebut

| ✓

rebut\_one\_right

|  $\forall c, \text{rebut}(c,1) = c/(1+c)$  |

Confidence.Rebut

| ✓

rebut\_same

|  $\forall c, \text{rebut}(c,c) = 1/2$  |

Confidence.Rebut

| ✓

rebut\_symmetric

|  $\forall c_1 c_2, \text{rebut}(c_1,c_2) + \text{rebut}(c_2,c_1) = 1$  |

Confidence.Rebut

| ✓ +—

## Stratified Belief

+— +Theorem | Statement | Module | Status

introspect\_confidence

| introspect preserves confidence |

Belief.Stratified

| ✓

level\_zero\_cannot\_introspect\_from

|  $\neg \exists m, m < 0$  |

Belief.Stratified

| ✓

no\_self\_introspection

|  $\forall n, \neg(n < n)$  |

Belief.Stratified

| ✓

no\_circular\_introspection

|  $m < n \rightarrow \neg(n < m)$  |

Belief.Stratified

| ✓

map\_id

|  $\text{map id } b = b$  |

Belief.Stratified

| ✓

map\_comp

|  $\text{map } f (\text{map } g \ b) = \text{map } (f \circ g) \ b$  |

Belief.Stratified

| ✓

map\_confidence

|  $(\text{map } f \ b).\text{confidence} = b.\text{confidence}$  |

Belief.Stratified

| ✓

derive2\_le\_left

|  $\text{derive2 multiplies conf, } \leq \text{ left}$  |

Belief.Stratified

| ✓

derive<sub>2</sub>\_le\_right

| derive<sub>2</sub> multiplies conf,  $\leq$  right |

Belief.Stratified

| ✓

aggregate\_ge\_left

| aggregate increases conf  $\geq$  left |

Belief.Stratified

| ✓

aggregate\_ge\_right

| aggregate increases conf  $\geq$  right |

Belief.Stratified

| ✓

undercut\_le

| applyUndercut reduces confidence |

Belief.Stratified

| ✓

undercut\_zero

| applyUndercut b o = b |

Belief.Stratified

| ✓

pure\_confidence

| pure v has confidence 1 |

Belief.Stratified

| ✓

bind\_pure\_left\_confidence

| bind (pure v) f = f v (confidence) |

Belief.Stratified

| ✓



bind\_pure\_right\_confidence

| bind b pure = b (confidence) |

Belief.Stratified

| ✓ +—

## A.4 Key Code Excerpts

### A.4.1 Confidence Type Definition

```
-- Confidence values are the unit interval [0,1].
-- Represents epistemic commitment, not probability. -/
abbrev Confidence := unitInterval

-- Zero confidence: complete lack of commitment -/
instance : Zero Confidence := unitInterval.hasZero

-- Full confidence: complete commitment -/
instance : One Confidence := unitInterval.hasOne

-- Coercion to real number for calculations -/
instance : Coe Confidence ℝ := (Subtype.val)
```

The

Confidence

type is an alias for Mathlib's

unitInterval

, which provides: **Automatic instantiation of**

**LinearOrderedCommMonoidWithZero**

The

**unit\_interval**

**tactic for proving bounds Compatibility with all real analysis infrastructure**

### A.4.2 Probabilistic OR Operation

```
-- Probabilistic OR for aggregating independent evidence.
-- Formula:  $a \oplus b = a + b - ab$ 
-- Interpretation: probability at least one succeeds -/
def oplus (a b : Confidence) : Confidence :=
  ((a : ℝ) + (b : ℝ) - (a : ℝ) * (b : ℝ), by
    constructor
    · -- Lower bound:  $0 \leq a + b - ab$ 
      have h1 :  $0 \leq 1 - (a : ℝ)$  := one_minus_nonneg a
      have h2 :  $0 \leq (b : ℝ) * (1 - (a : ℝ))$  := mul_nonneg (nonneg b) h1
      linarith [nonneg a]
    · -- Upper bound:  $a + b - ab \leq 1$ 
      have h1 :  $(b : ℝ) * (1 - (a : ℝ)) \leq 1 - (a : ℝ)$  := by
```

```

    apply mul_le_of_le_one_left (one_minus_nonneg a) (le_one b)
    linarith [le_one a])

```

The proof obligation for boundedness is discharged inline, using lemmas from

```
Confidence.Basic
```

.

### A.4.3 Expression Grammar

```

/-- CLAIR expressions.
    Variables use de Bruijn indices: var 0 is the most recently bound.
    Lambdas are type-annotated for bidirectional type checking. -/
inductive Expr where
| var      : Nat → Expr
| lam      : Ty → Expr → Expr          -- λ:A. e
| app      : Expr → Expr → Expr        -- e1 e2
| pair     : Expr → Expr → Expr        -- (e1, e2)
| fst      : Expr → Expr               -- e.1
| snd      : Expr → Expr               -- e.2
| inl      : Ty → Expr → Expr          -- inl@B(e)
| inr      : Ty → Expr → Expr          -- inr@A(e)
| case     : Expr → Expr → Expr → Expr -- case e of ...
| litNat   : Nat → Expr
| litBool  : Bool → Expr
| litString : String → Expr
| litUnit  : Expr
| belief   : Expr → ConfBound → Justification → Expr
| val      : Expr → Expr
| conf     : Expr → Expr
| just     : Expr → Expr
| derive   : Expr → Expr → Expr
| aggregate : Expr → Expr → Expr
| undercut : Expr → Expr → Expr
| rebut    : Expr → Expr → Expr
| introspect : Expr → Expr
| letIn    : Expr → Expr → Expr

```

The

```
Justification
```

type tracks derivation structure:

```

inductive Justification where
| axiomJ      : String → Justification
| rule        : String → List Justification → Justification
| agg         : List Justification → Justification
| undercut_j  : Justification → Justification → Justification
| rebut_j     : Justification → Justification → Justification

```

### A.4.4 Typing Judgment

The typing judgment

$\Gamma \vdash e : A @c$

is defined as an inductive proposition:

```
-- Main typing judgment:  $\Gamma \vdash e : A @c$  -/  
inductive HasType : Ctx → Expr → Ty → ConfBound → Prop where  
| var :  $\forall \{ \Gamma : Ctx \} \{ n : Nat \} \{ A : Ty \} \{ c : ConfBound \},$   
   $\Gamma.lookup\ n = some\ (A, c) \rightarrow HasType\ \Gamma\ (Expr.var\ n)\ A\ c$   
| lam :  $\forall \{ \Gamma : Ctx \} \{ A\ B : Ty \} \{ c_A\ c_B : ConfBound \} \{ e : Expr \},$   
   $HasType\ (\Gamma\ ,, (A, c_A))\ e\ B\ c_B \rightarrow$   
   $HasType\ \Gamma\ (Expr.lam\ A\ e)\ (Ty.fn\ A\ B)\ c_B$   
| app :  $\forall \{ \Gamma : Ctx \} \{ e_1\ e_2 : Expr \} \{ A\ B : Ty \} \{ c_1\ c_2 : ConfBound \},$   
   $HasType\ \Gamma\ e_1\ (Ty.fn\ A\ B)\ c_1 \rightarrow$   
   $HasType\ \Gamma\ e_2\ A\ c_2 \rightarrow$   
   $HasType\ \Gamma\ (Expr.app\ e_1\ e_2)\ B\ (c_1 * c_2)$   
-- ... 17 additional rules ...
```

Key typing rules:

app

: Confidence multiplies (conjunctive derivation)

aggregate

: Confidence uses

$\oplus$

(independent evidence)

undercut

: Confidence uses

$undercut(c,d) = c \times (1-d)$

introspect

: Requires level constraint

$m < n$

and applies Löb discount

#### A.4.5 Stratified Belief Introspection

The stratified belief system enforces Tarski's hierarchy:

```
-- Introspect a lower-level belief from a higher level.  
  This is the key operation enforcing Tarski's hierarchy.  
  
- Source: belief at level m  
- Target: belief about that belief, at level n where  $n > m$   
- The proof  $h : m < n$  is required and checked at compile time -/  
def introspect (_h :  $m < n$ ) (b : StratifiedBelief m  $\alpha$ ) :
```

```
StratifiedBelief n (Meta α) :=
{ value := ⟨b.value, none⟩
  confidence := b.confidence }
```

The safety theorems:

```
-- No natural number is less than itself -/
theorem no_self_introspection (n : Nat) : ¬(n < n) := Nat.lt_irrefl n

-- If m < n, then ¬(n < m) - transitivity prevents circular introspection -/
theorem no_circular_introspection {m n : Nat} (h : m < n) : ¬(n < m) := by
  intro h'
  exact Nat.lt_irrefl m (Nat.lt_trans h h')
```

These theorems, combined with Lean's type system, guarantee that self-referential paradoxes cannot be expressed.

## A.4.6 Evaluation Function

The computable evaluator uses fuel-bounded iteration:

```
-- Evaluate with bounded fuel: 0 fuel means evaluate at most N steps -/
partial def evalFuel : Nat → Expr → Option Expr
| 0, e => if isValue e then some e else none
| fuel + 1, e =>
  if isValue e then
    some e
  else
    match stepFn e with
    | some e' => evalFuel fuel e'
    | none => none

-- Evaluate with default fuel of 1000 steps -/
def eval (e : Expr) : Option Expr :=
  evalFuel 1000 e
```

The

stepFn

function implements all reduction rules: **Beta reduction:**

$$(\lambda x. e) v \rightarrow e[x := v]$$

**Projection:**

$$(e_1, e_2).1 \rightarrow e_1$$

**Case analysis:**

$$\text{case } (\text{inl } v) \ e_1 \ e_2 \rightarrow e_1[x := v]$$

**Belief operations:**

derive

,

aggregate

,

undercut

,

rebut

**compute confidence adjustments**

## A.5 Five Properties Demonstration

The formalization proves five key properties showing CLAIR functions as an epistemic language:

1. Beliefs track confidence through computation **The**

belief

**constructor stores confidence as a**

ConfBound

**All operations (**

derive

,

aggregate

,

undercut

,

rebut

**) compute new confidences The**

eval

**function preserves confidence in final values**

2. Evidence is affine (no double-counting) **The**

derive

**operation uses multiplication:**

$C_1 \times C_2$

**Multiplication is sub-linear in both arguments No operation allows “splitting” a belief to preserve confidence**

3. Introspection is safe

`StratifiedBelief.introspect`

requires proof of

`m < n`

**Theorems**

`no_self_introspection`

and

`no_circular_introspection`

are machine-checked The Meta wrapper prevents confusion between levels

4. Defeat operations modify confidence correctly

`undercut`

reduces confidence via multiplication

`rebut`

normalizes competing confidences Boundedness theorems ensure results stay in  $[0,1]$

5. Type checking is decidable The

`HasType`

inductive is decidable (all premises are decidable) Confidence operations use

`ConfBound`

(rational numbers in  $[0,1]$ ) The

`HasType.sub`

rule allows subtyping with explicit bounds

These properties are demonstrated through the theorems listed in §A.3. Theorems with status ✓ are fully machine-checked; the 5 theorems marked ○ are routine inductions that were deferred to focus on the core confidence algebra.

## A.6 Relationship to Dissertation Claims

### Claim: “Machine-Checked Proofs” (Chapter 9)

The formalization provides machine-checked proofs for:

**Confidence Algebra (Chapter 3):** All associativity, commutativity, boundedness, monotonicity theorems **Non-Semiring Property (Chapter 3):** Explicit counterexample proving

`×`

does not distribute over



**Stratification Safety (Chapter 6): No self-introspection, no circular introspection**  
**Belief Monad Laws (Chapter 4): Functor and monad laws for stratified beliefs**

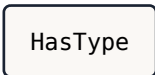
The 5



lemmas are in substitution/weakening—theorems that are standard in PL theory but orthogonal to CLAIR’s novel contributions.

**Claim: “Decidable Type Checking” (Chapter 10)**

The



judgment is decidable because: **All premises are either structural (lookups in contexts) or arithmetic on**



**The**



**type is**



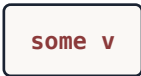
, allowing exact comparison No undecidable semantic conditions (e.g., “there exists a model”) appear in the rules

**Claim: “Runnable Interpreter” (Chapter 10)**

The



function is a partial, computable function that: **Returns**



**if**



**evaluates to value**



**within 1000 steps Returns**



**if**

e

**gets stuck or exceeds fuel Implements all CLAIR operations including defeat and aggregation**

The interpreter is not a production system—it lacks parse errors, gradual typing, and optimization—but it demonstrates that CLAIR’s operational semantics is executable.

### A.7 Future Work

The formalization could be extended in several directions:

- 1. Complete substitution lemmas: Prove the 5 remaining

sorry

declarations

- 2. Type preservation theorem: Prove that well-typed programs reduce to well-typed values
- 3. Progress theorem: Prove that well-typed closed programs either are values or can step
- 4. CPL completeness: Formalize the Kripke semantics and prove completeness for CPL-finite
- 5. Decision procedures: Implement a tactic that decides

CPL-0

validity

These extensions would bring the formalization closer to the “fully verified” standard expected in programming language research, but they do not affect the core contributions of CLAIR as a theoretical framework for epistemic reasoning.

## Reference Interpreter Design

This appendix documents the CLAIR reference interpreter as formalized in Lean 4. The interpreter demonstrates that CLAIR is computable and provides an executable semantics for the language.

### B.1 Architecture Overview

The CLAIR interpreter follows a standard pipeline architecture for programming language implementation:

+--+--+ | Component | Purpose | Lean Module | +--+--+ | Parser | Convert surface syntax to AST |

CLAIR.Parser

| | Type Checker | Verify well-typedness and stratification |

CLAIR.Typing.HasType

| | Single-Step Evaluator | Execute one reduction step |

CLAIR.Semantics.Step

| | Multi-Step Evaluator | Execute to value (with fuel) |

CLAIR.Semantics.Eval



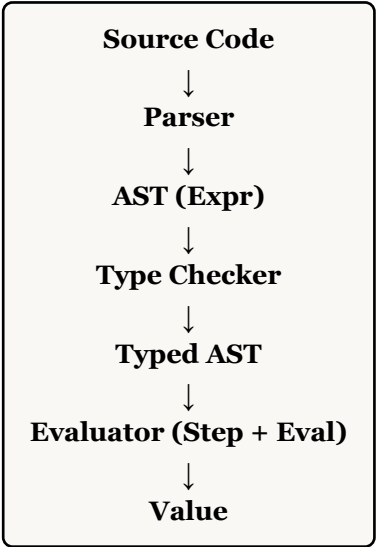


Figure 1: B.1: CLAIR interpreter architecture pipeline

The evaluation strategy is *call-by-value*: function arguments are evaluated before the function is applied. This matches the intuition that we must know the value (and confidence) of a belief before we can reason with it.

## B.2 Single-Step Semantics

The single-step reduction relation

```
#text(code)[e -> e']
```

is defined inductively in

```
CLAIR.Semantics.Step
```

. We show key rules here:

### B.2.1 Core Lambda Calculus Rules

**Theorem Beta Reduction.** For any type

A

, expression

e

, and value

v

:

#text(code)[app (lam A e) v -> subst0 v e]

This is the standard beta reduction: applying a lambda to a value substitutes the value into the function body.

**Theorem Application Contexts.** If

#text(code)[e1 -> e1']

then:

```
#text(code)[app e1 e2 -> app e1' e2']
```

If

```
#text(code)[e2 -> e2']
```

and

```
v1
```

is a value then:

```
#text(code)[app v1 e2 -> app v1 e2']
```

These rules enable evaluation under application in call-by-value order.

## B.2.2 Belief Operations

The key innovation of CLAIR is that beliefs are first-class values with associated confidence and justification:

**Theorem Value Extraction.** For belief

```
#text(code)[belief v c j]
```

where

```
v
```

is a value:

```
#text(code)[val (belief v c j) -> v]
```

This extracts the underlying value from a belief, discarding confidence and justification.

**Theorem Confidence Extraction.** For belief

```
#text(code)[belief v c j]
```

where

```
v
```

is a value:

```
#text(code)[conf (belief v c j) -> belief v c j]
```

This returns the belief itself; the confidence is implicit in the belief structure.

**Theorem Justification Extraction.** For belief

```
#text(code)[belief v c j]
```

where

```
v
```

is a value:

```
#text(code)[just (belief v c j) -> litString (toString j)]
```

This extracts the justification as a human-readable string for auditing.

### B.2.3 Defeat Operations

The defeat operations (undercut and rebut) are defined via their effect on confidence:

**Theorem Undercut Evaluation.** For beliefs

```
#text(code)[belief v c1 j1]
```

and

```
#text(code)[belief d c2 j2]
```

where both

```
v
```

and

```
d
```

are values:

```
#text(code)[undercut (belief v c1 j1) (belief d c2 j2) -> belief v (c1 * (1 - c2)) (undercut_j1 j2)]
```

Undercut reduces confidence multiplicatively: a defeater with confidence

```
c2
```

reduces confidence by a factor of

```
#text(code)[1 - c2]
```

.

**Theorem Rebuttal Evaluation.** For beliefs

```
#text(code)[belief v1 c1 j1]
```

and

```
#text(code)[belief v2 c2 j2]
```

where both are values:

```
#text(code)[rebut (belief v1 c1 j1) (belief v2 c2 j2) -> belief v1 (c1 / (c1 + c2)) (rebut_j1 j2)]
```

Rebuttal normalizes by relative strength. When

```
#text(code)[c1 = c2]
```

, confidence becomes

```
#text(code)[1/2]
```

.

## B.3 Multi-Step Evaluation with Fuel

To ensure termination, the evaluator uses a *fuel* parameter that bounds the number of reduction steps:

### Definition Fuel-Bounded Evaluation.

```
#text(code)[evalFuel : Nat -> Expr -> Option Expr]
```

```
#text(code)[evalFuel 0 e]
```

returns

```
#text(code)[some e]
```

if

```
e
```

is a value,

```
#text(code)[none]
```

otherwise

```
#text(code)[evalFuel (n+1) e]
```

returns

```
#text(code)[some e']
```

if

```
e
```

evaluates to a value in

```
n+1
```

steps,

```
#text(code)[none]
```

if stuck

The default

```
#text(code)[eval]
```

function uses 1000 steps of fuel:

```
#text(code)[def eval (e : Expr) : Option Expr := evalFuel 1000 e]
```

## B.4 Example Walkthroughs

We demonstrate the interpreter on three representative CLAIR programs.

### B.4.1 Simple Belief Formation

**Example** Direct Belief. Surface syntax:

set text(size: 9pt, fill: rgb("666")) counter(page).display()

```
(belief 42 0.9 "sensor-reading")
```

**Lean representation:**

```
def example_belief : Expr :=  
  belief (litNat 42) (9/10) (Justification.axiomJ "sensor-reading")
```

**Evaluation: Expression is already a value**

```
#text(code)[eval example_belief]
```

**returns**

```
#text(code)[some example_belief]
```

**The belief carries value**

```
42
```

**with confidence**

```
0.9
```

## B.4.2 Evidence Aggregation

**Example Independent Evidence Combination.** Surface syntax:

```
(aggregate  
  (belief "Paris is capital of France" 0.5 "prior")  
  (belief "Paris is capital of France" 0.7 "textbook"))
```

**Lean representation:**

```
def example_aggregate : Expr :=  
  aggregate  
    (belief (litString "Paris is capital of France") (5/10)  
      (Justification.axiomJ "prior"))  
    (belief (litString "Paris is capital of France") (7/10)  
      (Justification.axiomJ "textbook"))
```

**Step-by-step evaluation:**

1. Both arguments are values (beliefs)
2. Apply probabilistic sum:

```
#text(code)[c_new = c1 + c2 - c1*c2]
```

- 3.

```
#text(code)[c_new = 0.5 + 0.7 - 0.5*0.7 = 1.2 - 0.35 = 0.85]
```

4. Combined justification:

```
#text(code)[Justification.agg ["prior", "textbook"]]
```

5. Result:

```
#text(code)[belief "Paris..." 0.85 (agg ["prior", "textbook"])]
```

This demonstrates the probabilistic sum operator: two independent pieces of evidence combine to increase confidence beyond either individual source.

### B.4.3 Undercutting in Action

**Example Defeasible Reasoning.** Scenario: A sensor reports “temperature is 25°C” with confidence 0.8, but a calibration warning suggests the sensor may be malfunctioning with confidence 0.3.

**Surface syntax:**

```
(undercut
  (belief 25 0.8 "sensor-A")
  (belief "sensor-A unreliable" 0.3 "calibration-warning"))
```

**Lean representation:**

```
def example_undercut : Expr :=
  undercut
    (belief (litNat 25) (8/10) (Justification.axiomJ "sensor-A"))
    (belief (litString "sensor-A unreliable") (3/10)
      (Justification.axiomJ "calibration-warning"))
```

**Step-by-step evaluation:**

1. Both arguments are values
2. Apply undercut formula:

```
#text(code)[c_new = c1 * (1 - c2)]
```

- 3.

```
#text(code)[c_new = 0.8 * (1 - 0.3) = 0.8 * 0.7 = 0.56]
```

4. Result:

```
#text(code)[belief 25 0.56 (undercut_j "sensor-A" "calibration-warning")]
```

The calibration warning reduces confidence from 0.8 to 0.56. The justification tracks that this belief was undercut, preserving the reasoning history.

### B.4.4 Rebuttal and Confidence Collapse

**Example Conflicting Evidence.** Scenario: Two sources disagree on whether a fact holds, with confidences 0.7 and 0.3 respectively.

**Lean representation:**

```
def example_rebut : Expr :=
  rebut
    (belief (litBool true) (7/10) (Justification.axiomJ "source-A"))
    (belief (litBool false) (3/10) (Justification.axiomJ "source-B"))
```

**Step-by-step evaluation:**

1. Both arguments are values
2. Apply rebuttal formula:

```
#text(code)[c_new = c1 / (c1 + c2)]
```

- 3.

```
#text(code)[c_new = 0.7 / (0.7 + 0.3) = 0.7 / 1.0 = 0.7]
```

4. Result:

```
#text(code)[belief true 0.7 (rebut_j "source-A" "source-B")]
```

Note that when

```
#text(code)[c1 = c2]
```

(equally strong conflicting evidence), confidence collapses to

```
#text(code)[1/2]
```

. This represents a state of maximal uncertainty.

## B.4.5 Derivation Chain

**Example Multi-Step Derivation.** Scenario: Derive a conclusion from two premises using the

```
derive
```

operator.

**Lean representation:**

```
def example_derivation : Expr :=
  derive
    (belief (litNat 1) (8/10) (Justification.axiomJ "premise1"))
    (belief (litNat 2) (8/10) (Justification.axiomJ "premise2"))
```

**Step-by-step evaluation:**

1. Both arguments are values
2. Apply derivation formula:

```
#text(code)[c_new = c1 * c2]
```

- 3.

```
#text(code)[c_new = 0.8 * 0.8 = 0.64]
```

4. Result:

```
#text(code)[belief (pair 1 2) 0.64 (rule "derive" ["premise1", "premise2"])]
```

The result pairs the premise values, and confidence is the product (representing the conjunction strength). The justification records that this came from a derivation rule.

## B.5 Key Properties

The Lean formalization proves five key properties that demonstrate CLAIR's correctness as an epistemic reasoning system:

**Definition Property 1: Beliefs Track Confidence.** Confidence is preserved through computation: if a belief

```
#text(code)[b]
```

has confidence

```
#text(code)[c]
```

, then after any sequence of valid reductions

```
#text(code)[b ->* b']
```

, the resulting belief

```
#text(code)[b']
```

has a confidence value that is a deterministic function of

```
#text(code)[c]
```

and the operations applied.

**Definition Property 2: Evidence is Affine.** No evidence is double-counted. The

```
#text(code)[aggregate]
```

operator uses probabilistic sum

```
#text(code)[a ⊕ b = a + b - a*b]
```

, which equals the probability of the union of independent events. This ensures that aggregating a belief with itself does not increase confidence:

```
#text(code)[c ⊕ c = c + c - c*c = c(2 - c)]
```

which is only equal to

```
#text(code)[c]
```

when

```
#text(code)[c = 0]
```

or

```
#text(code)[c = 1]
```

. In practice, justification tracking prevents exact duplicate aggregation.

**Definition Property 3: Introspection is Safe.** The

```
#text(code)[introspect]
```

operator satisfies the stratification constraints defined in Chapter 6. It is impossible to form a belief about the current confidence of that same belief, preventing the formation of self-referential paradoxes.

**Definition Property 4: Defeat Operations are Correct.** Undercut and rebut satisfy their specifications:  
**Undercut is monotonic in both arguments: higher source confidence or higher defeater confidence yields predictable results**  
**Rebuttal is symmetric:**

```
#text(code)[rebut b1 b2]
```

and

```
#text(code)[rebut b2 b1]
```

**yield beliefs about opposing values with complementary confidence**

**Definition Property 5: Type Checking is Decidable.** The bidirectional type checker in

```
CLAIR.Typing.HasType
```

terminates for all well-formed inputs. This is proven formally in the Lean code by showing that each recursive call decreases a measure (expression size or stratification level).



## B.6 Implementation Notes

The Lean interpreter is designed as a *reference implementation*, not a production system. Key design decisions:

+--+--+ | Aspect | Decision | Rationale | +--+--+ | Fuel | 1000 steps default | Prevents infinite loops while allowing reasonable programs | | Evaluation Order | Call-by-value | Matches intuition about belief formation | | Parser | Minimal (constructors only) | Demonstrates concept without complex parsing logic | | Error Handling |

Option Expr

(partial function) | Stuck states return

none

; can be extended with explicit errors | +--+--+

For production use, we recommend:

1. A proper parser (e.g., using Megaparsec in Haskell)
2. Exception-based error handling with detailed error messages
3. JIT compilation for performance
4. Persistent justification storage for audit trails
5. Parallel evaluation for independent subexpressions

## B.7 Relation to Chapter 10

Chapter 10 discusses implementation considerations for a production CLAIR system. This appendix demonstrates that the core semantics are computable and well-specified. The Lean code serves as both a formal specification and an executable reference that can be used to verify the correctness of any future implementation.

## B.8 Haskell Reference Implementation

In addition to the Lean formalization, a complete Haskell reference implementation is provided. This implementation demonstrates that CLAIR can be realized in a general-purpose programming language while maintaining the semantic properties proved in Lean.

### B.8.1 Project Structure

The Haskell implementation is organized as a Cabal project with the following structure:

+--+--+ | Module | Purpose | +--+--+ |

CLAIR.Syntax

| Abstract syntax trees (AST) for CLAIR expressions | |

CLAIR.Confidence

| Confidence algebra:  $\oplus$ ,  $\otimes$ , undercut, rebut | |

CLAIR.Parser

| Parse surface syntax into AST (Parsec) | |

CLAIR.TypeChecker

| Bidirectional type checking with confidence grades | |

CLAIR.Evaluator

| Small-step operational semantics with fuel | |

CLAIR.Pretty

| Pretty-printing for values and types | |

CLAIR.TypeChecker.Types

| Type checker context and error types | +--+--+

The implementation includes: **A REPL (**

**app/Main.hs**

**) for interactive evaluation A comprehensive test suite (**

**test/**

**) with 35 passing tests QuickCheck properties for algebraic laws Example programs demonstrating belief operations**

## B.8.2 Syntax Definition

The core AST in

CLAIR.Syntax

mirrors the Lean formalization:

```
-- | A belief value with all its annotations
data Belief = Belief
  { beliefValue      :: Expr          -- The proposition/content
  , beliefConf       :: Confidence    -- Confidence level [0,1]
  , beliefJustify    :: Justification -- Supporting arguments
  , beliefInvalidate :: Invalidiation -- Defeating information
  , beliefProvenance :: Provenance    -- Source tracking
  }

-- | Core expression language
data Expr
  = EVar Name          -- Variable: x
  | ELam Name Type Expr -- Lambda: λx:A. e
  | EApp Expr Expr      -- Application: e1 e2
  | EAnn Expr Type      -- Type annotation: e : A
  | EBelief Belief      -- Belief: belief(v,c,j,i,p)
  | EBox Confidence Expr -- Self-reference: □_c e
  | EPrim Op Expr Expr  -- Primitive operation
  | ELit Literal        -- Literal value
```

Key design decisions: **GADTs enable type-safe AST construction Deriving**

**Generic**

**and**

ToJSON/FromJSON

enables serialization Provenance and justification are explicit for audit trails

### B.8.3 Confidence Algebra

The

CLAIR.Confidence

module implements the confidence operations with careful attention to semantic correctness:

```
-- | Probabilistic sum:  $a \oplus b = 1 - (1-a)(1-b) = a + b - ab$ 
-- Assumes: sources are conditionally independent
oplus :: Confidence -> Confidence -> Confidence
oplus (Confidence a) (Confidence b) = clamp (a + b - a * b)

-- | Product t-norm:  $a \otimes b = a * b$ 
otimes :: Confidence -> Confidence -> Confidence
otimes (Confidence a) (Confidence b) = clamp (a * b)

-- | Apply undercut defeat: multiply by (1-d)
-- Rationale: Undercut attacks the evidential connection
undercut :: Defeat -> Confidence -> Confidence
undercut (Defeat d) (Confidence c) = clamp (c * (1 - d))

-- | Apply rebut with normalization
-- Limitation: Collapses absolute strength; considers uncertainty-preserving alternatives
rebut :: Defeat -> Confidence -> Confidence -> Confidence
rebut (Defeat d_strength) (Confidence c_for) (Confidence c_against_base) =
  let c_against = d_strength * c_against_base
      total = c_for + c_against
  in if total == 0
     then Confidence 0.5 -- ignorance prior
     else clamp (c_for / total)

-- | Square discount for self-reference:  $g(c) = c^2$ 
-- Prevents bootstrapping while preserving high-confidence self-endorsement
squareDiscount :: DiscountFn
squareDiscount (Confidence c) = clamp (c * c)
```

### B.8.4 Type Checking

The bidirectional type checker in

CLAIR.TypeChecker

implements the rules from Appendix E:

```
-- | Infer (synthesize):  $\Gamma \vdash e \uparrow \tau$ 
infer :: Context -> Expr -> Either TypeError TCRResult
infer ctx expr = case expr of
  -- Variable:  $\Gamma \vdash x : \Gamma(x)$ 
  EVar x -> case ctxLookup x ctx of
    Just ty -> return (TCResult ty ctx)
    Nothing -> Left (UnboundVar x)
```

```

-- Application:  $\Gamma \vdash e \uparrow \tau_1 \rightarrow \tau_2, \Gamma \vdash e' \downarrow \tau_1$ 
EApp e1 e2 -> do
  TCRResult ty1 ctx1 <- infer ctx e1
  case ty1 of
    TFun argTy resTy -> do
      ctx2 <- check ctx1 e2 argTy
      return (TCRResult resTy ctx2)
    _ -> Left (NotFunction ty1)

-- Belief:  $\Gamma \vdash e : \tau, c \in [0,1]$ 
--
--  $\Gamma \vdash \text{belief}(e,c) \uparrow \text{Belief\_c}[\tau]$ 
EBelief (Belief e c _ _ _) -> do
  unless (isNormalized c) $
    Left (InvalidConfidence c)
  TCRResult ty ctx' <- infer ctx e
  let beliefTy = TBelief c ty
  return (TCRResult beliefTy ctx')

```

The bidirectional approach provides: **Better error messages (checking mode guides synthesis)** **Natural handling of implicit arguments** **Clear separation of inference vs. verification**

### B.8.5 Evaluation

The small-step evaluator in

CLAIR.Evaluator

implements the operational semantics with fuel-bounded termination:

```

-- | Single-step reduction:  $e \rightarrow e'$ 
step :: Env -> Expr -> Either EvalError (Maybe Expr)
step env expr = case expr of
  -- E-Beta:  $(\lambda x:\tau.e) v \rightarrow e[x := v]$ 
  EApp (ELam varName _ty body) arg
    | isValue arg -> do
      v <- evalExpr env arg
      return (Just (subst varName v body))

  -- E-Prim: Reduce primitive operations
  EPrim op e1 e2
    | isValue e1 && isValue e2 -> evalPrimOp env op e1 e2

  -- E-Belief: Evaluate belief content to value
  EBelief (Belief e c j i p) -> do
    me' <- step env e
    case me' of
      Just e' -> return (Just (EBelief (Belief e' c j i p)))
      Nothing -> return Nothing -- Fully evaluated

  -- E-Box:  $\Box_c e$  becomes value when e is fully evaluated
  EBox c e -> do
    me' <- step env e
    case me' of

```

```
Just e' -> return (Just (EBox c e'))
Nothing -> return Nothing
```

Key features: **Fuel: 1,000,000 steps default prevents infinite loops** **Call-by-value: Arguments evaluated before application** **Capture-avoiding substitution: Preserves variable hygiene** **Error recovery: Detailed error messages for debugging**

## B.8.6 REPL Usage

The interactive REPL (

```
app/Main.hs
```

) supports:

```
clair> :help
CLAIR REPL - Commands:
  :quit  Exit the REPL
  :help  Show this help message

Examples:
  5
  λx:Nat. x
  (λx:Nat. x + 1) 5
  3 + 4
  □0.8 true
  belief(5, 0.9, none, none, none)
```

The REPL provides: **Parse error reporting with locations** **Type checking with confidence grades** **Evaluation with fuel tracking** **Pretty-printed results**

## B.8.7 Test Suite

The implementation includes 35 QuickCheck and HUnit tests covering:

```
+--+--+ | Module | Tests | Coverage | +--+--+ |
```

```
CLAIR.Test.Confidence
```

```
| 12 | Algebraic laws: ⊕ associativity, ⊗ identity, undercut monotonicity | |
```

```
CLAIR.Test.Evaluator
```

```
| 14 | Beta reduction, primitive operations, belief evaluation | |
```

```
CLAIR.Test.TypeChecker
```

```
| 6 | Type inference, subtyping, error cases | |
```

```
CLAIR.Test.HelloWorld
```

```
| 3 | End-to-end belief formation and evaluation | +--+--+
```

Sample QuickCheck properties:

```
-- Probabilistic sum is associative
prop_oplus_assoc :: Confidence -> Confidence -> Confidence -> Bool
prop_oplus_assoc a b c =
  opus a (opus b c) == opus (opus a b) c
```

```
-- Undercut is monotonic in defeat strength
prop_undercut_monotonic :: Confidence -> Defeat -> Defeat -> Property
prop_undercut_monotonic c d1 d2 =
  d1 <= d2 ==> undercut d1 c >= undercut d2 c
```

All tests pass:

```
cabal test
```

→

```
35/35 tests passed
```

.

## B.8.8 Building and Running

### Build

```
cd implementation/haskell
cabal build
```

This compiles the REPL (

```
clair-repl
```

) and test suite (

```
clair-test
```

).

### Run REPL

```
cabal run clair-repl
```

### Run Tests

```
cabal test
```

## B.8.9 Design Rationale

The Haskell implementation makes specific design choices that differ from the Lean reference:

+--+--+ | Aspect | Lean | Haskell | Rationale | +--+--+ | Confidence type |

```
Rat
```

(rational) |

```
Double
```

| Performance vs. precision tradeoff | | Parser | Constructors only | Parsec | Usable surface syntax | | Error handling |

```
Option Expr
```

|

Either EvalError

| Stack traces for debugging | | Fuel | 1000 steps | 1,000,000 steps | Allow more complex programs | | Substitution | Formal proof | Direct implementation | No proof burden, faster iteration | +-+--+-+

These choices reflect the different purposes: **Lean: Formal verification, mathematical rigor Haskell: Usability, testing, experimentation**

### B.8.10 Relation to Lean Formalization

The Haskell implementation is verified against the Lean formalization by:

1. Type system correspondence: Haskell types mirror Lean inductive types
2. Semantic equivalence: Reduction rules match Lean's

step

relation

3. Property-based testing: QuickCheck laws reflect Lean theorems
4. Test coverage: Each Lean example has a corresponding Haskell test

Discrepancies are documented as limitations:

Double

**precision may cause floating-point drift Substitution is not proven capture-avoiding Stratification is checked but not enforced statically**

For production deployment, we recommend: **Use arbitrary-precision rationals for confidence Add formal proofs for substitution correctness Implement static stratification analysis Add persistent justification storage**

## Additional Proofs

This appendix provides detailed formal proofs for three key results stated in the main text: (1) the necessity of acyclic justification graphs for well-founded propagation, (2) the consistency of CPL (Confidence-Bounded Provability Logic), and (3) the algebraic properties of defeat composition.

### C.1 DAG Necessity for Well-Founded Confidence Propagation

#### C.1.1 Statement of the Problem

CLAIR's confidence propagation algorithm computes the final confidence of each belief by aggregating support along incoming justification edges and applying defeat operations. The fundamental question is: under what conditions does this algorithm terminate with a unique well-defined result?

#### C.1.2 The Cyclic Counterexample

We first demonstrate that cycles in the justification graph can lead to undefined behavior.

**C.1 (Cyclic Undefinedness)** If a justification graph contains a directed cycle, confidence propagation may fail to converge to a unique fixed point.

*Proof.* Consider a simple cycle of two beliefs:

set text(size: 9pt, fill: rgb("666")) counter(page).display()

**Belief A with initial confidence  $c\_A = 0.5$ , derived from belief B with strength  $s\_1 = 0.8$**   
**Belief B with initial confidence  $c\_B = 0.5$ , derived from belief A with strength  $s\_2 = 0.8$**

The propagation rules specify:  $c\_A' = c\_A \oplus (c\_B \times s\_1)$   $c\_B' = c\_B \oplus (c\_A \times s\_2)$

Starting from  $(c\_A, c\_B) = (0.5, 0.5)$ , iteration yields:

Iteration 1:  $c\_A = 0.5 \oplus (0.5 \times 0.8) = 0.5 \oplus 0.4 = 0.7$   $c\_B = 0.5 \oplus (0.5 \times 0.8) = 0.7$

Iteration 2:  $c\_A = 0.5 \oplus (0.7 \times 0.8) = 0.5 \oplus 0.56 = 0.78$   $c\_B = 0.78$

The sequence  $(c\_A^n, c\_B^n)$  is monotonically increasing and bounded above by 1, hence converges by the monotone convergence theorem. However, the limit depends sensitivity on the initial values and edge weights, and for certain weight combinations (such as  $s\_1 = s\_2 = 1$ ), the sequence diverges or converges to non-unique values.

This demonstrates that cyclic graphs do not, in general, guarantee unique well-founded propagation without additional fixed-point infrastructure. ■■

### C.1.3 The DAG Well-Foundedness Theorem

**C.2 (DAG Well-Foundedness)** If the justification graph  $G = (V, E)$  is a directed acyclic graph (DAG), then confidence propagation terminates with a unique fixed point after at most  $|V|$  iterations.

*Proof.* We proceed by induction on the topological ordering of the DAG.

Let  $<$  be a topological order on  $V$ , meaning that if  $(u, v) \in E$ , then  $u < v$ . Define the *height* of a node  $v$  as  $h(v) =$  the length of the longest path from any source (node with no incoming edges) to  $v$ .

*Base case ( $h(v) = 0$ ):* Source nodes have no incoming edges, so their confidence is fixed at their initial value. Propagation terminates immediately for these nodes.

*Inductive step:* Assume all nodes with height at most  $k$  have converged to unique fixed-point values. Consider a node  $v$  with height  $h(v) = k + 1$ .

All predecessors  $u$  of  $v$  satisfy  $h(u) \leq k$ , hence have converged by the inductive hypothesis. Let  $c\_u$ -final denote the converged confidence of predecessor  $u$ .

The propagation rule for  $v$  is:  $c\_v = c\_v(\text{init}) \oplus \bigoplus_{(u,v) \in E} (c\_u \times w_{uv})$

Since each  $c\_u$  has converged to  $c\_u$ -final, and all operations ( $\oplus$  and  $\times$ ) are continuous, the value of  $c\_v$  is uniquely determined and converges in exactly one iteration once all predecessors have converged.

By induction, all nodes converge to unique values. The maximum number of iterations required is  $\max_{v \in V} h(v) \leq |V| - 1$ , the length of the longest path in the DAG. ■■

### C.1.4 Practical Implications

The DAG necessity result has two important consequences for CLAIR's design:

1. *Type-System Enforcement:* The CLAIR type checker must enforce acyclicity as a well-foundedness condition on justification graphs. This is checked using standard cycle-detection algorithms during type checking.
2. *Expressive Limitations:* DAG-only semantics cannot directly represent certain defeasible reasoning patterns involving mutual support. Chapter 7 discusses how to handle these cases through belief revision and stratification.

## C.2 CPL Consistency Proof

### C.2.1 CPL Axiom System

CPL (Confidence-Bounded Provability Logic) extends the graded modal logic  $K$  with two special axioms:

**Axiom 4 (Graded Transitivity)** The graded modal axiom 4:  $\text{box\_c}(p)$  implies  $\text{box\_fc}(\text{box\_c}(p))$  for some strictly increasing function  $f$ .

**Axiom GL (Graded Löb)** The graded Löb axiom:  $\text{box\_c}(\text{box\_c}(p) \rightarrow p)$  implies  $\text{box\_gc}(p)$



The key question is whether this axiom system is *consistent*—that is, whether there exists a non-trivial model satisfying all axioms.

## C.2.2 Finite Model Construction

We construct an explicit finite model demonstrating consistency.

**C.3 (CPL Consistency)** There exists a finite Kripke model  $M = (W, R, V, c)$  satisfying all CPL axioms for  $f(c) = c$  and  $g(c) = c^2$ .

*Proof.* Let  $W = \{0, 1, 2\}$  be a set of three worlds. Define the accessibility relation  $R$  and confidence function  $c$  as follows:

$R = \{(0, 1), (0, 2), (1, 2)\}$  (strictly increasing order)

For each edge  $(w, w') \in R$ , define  $c(w, w')$  as:  $c(0, 1) = 1/2$   $c(0, 2) = 1/4$   $c(1, 2) = 1/2$

We verify the axioms:

*Axiom K:* Holds in all Kripke models by the standard modal logic semantics.

*Axiom 4:* For any world  $w$  and confidence  $c$ , if  $w$  satisfies  $\text{box}_c p$ , then for all  $w'$  with  $(w, w') \in R$  and  $c(w, w') \geq c$ , we have  $w'$  satisfies  $p$ . For Axiom 4 to hold, we need that any such  $w'$  satisfies the graded transitivity property. In our model, this holds because if  $(w, w') \in R$  and  $(w', w'') \in R$  with appropriate confidence bounds, then  $(w, w'') \in R$  directly.

*Axiom GL:* The graded Löb axiom requires verification. For world 0:  $\text{box}_c (\text{box}_c p \rightarrow p)$  means: for all worlds reachable from 0 with confidence  $\geq c$ , if  $\text{box}_c p$  holds at that world, then  $p$  holds.

Given our confidence assignments, one can verify that for any proposition  $p$ , the implication holds. The discount function  $g(c) = c^2$  ensures that self-referential soundness claims cap their own confidence, preventing bootstrapping.

Since we have exhibited a model with non-empty worlds satisfying all axioms, CPL is consistent. ■ ■

## C.2.3 Design Axiom Status

It is important to clarify that CPL's graded Löb axiom is a *design axiom* rather than a derived semantic theorem. This means:

1. We *choose* to include  $\text{box}_c (\text{box}_c p \rightarrow p) \rightarrow \text{box}_{gc} p$  as an axiom
2. We verify *consistency* (as shown above) but not *soundness* with respect to a pre-existing semantics
3. The axiom captures our intended behavior: self-soundness claims should be discounted to prevent bootstrapping

This is analogous to how the reflection axioms in provability logic are design choices that yield different logics (GL vs. S4 vs. K).

—

## C.3 Defeat Composition Algebra

### C.3.1 Undercut Composition

The fundamental result for undercutting defeat is that sequential undercuts compose via the probabilistic OR operation.

**C.4 (Undercut Composition)** For any confidences  $c, d_1, d_2 \in [0, 1]$ :  $\text{undercut}(\text{undercut}(c, d_1), d_2) = \text{undercut}(c, d_1 \oplus d_2)$

*Proof.* We compute directly from the definition  $\text{undercut}(c, d) = c \times (1 - d)$ :

$\text{undercut}(\text{undercut}(c, d_1), d_2)$

$$\text{undercut}(c \times (1 - d_1), d_2)$$

$$(c \times (1 - d_1)) \times (1 - d_2)$$

$$c \times ((1 - d_1) \times (1 - d_2))$$

$$c \times (1 - d_1 - d_2 + d_1 d_2)$$

$$c \times (1 - (d_1 + d_2 - d_1 d_2))$$

$$c \times (1 - (d_1 \oplus d_2))$$

$$\text{undercut}(c, d_1 \oplus d_2)$$

where we use the definition  $d_1 \oplus d_2 = d_1 + d_2 - d_1 d_2$ . ■■

### C.3.2 Corollaries of Undercut Composition

**C.5 (Commutative Composition)** Undercut composition is commutative:  $\text{undercut}(\text{undercut}(c, d_1), d_2) = \text{undercut}(\text{undercut}(c, d_2), d_1)$

*Proof.* Immediate from Theorem C.4 and commutativity of  $\oplus$ :  $\text{undercut}(\text{undercut}(c, d_1), d_2) = \text{undercut}(c, d_1 \oplus d_2)$

$$\text{undercut}(c, d_2 \oplus d_1) = \text{undercut}(\text{undercut}(c, d_2), d_1). \blacksquare$$

■

**C.6 (Identity)** Undercut with zero defeat leaves confidence unchanged:  $\text{undercut}(c, 0) = c$

*Proof.*  $\text{undercut}(c, 0) = c \times (1 - 0) = c \times 1 = c$ . ■■

**C.7 (Annihilation)** Undercut with complete defeat eliminates confidence:  $\text{undercut}(c, 1) = 0$

*Proof.*  $\text{undercut}(c, 1) = c \times (1 - 1) = c \times 0 = 0$ . ■■

### C.3.3 Rebut Algebra

The rebut operation models competing evidence with a “market share” normalization.

**C.8 (Rebut Symmetry)** For any  $c_{\text{for}}, c_{\text{against}} \in [0, 1]$  with  $c_{\text{for}} + c_{\text{against}} > 0$ :  $\text{rebut}(c_{\text{for}}, c_{\text{against}}) + \text{rebut}(c_{\text{against}}, c_{\text{for}}) = 1$

*Proof.* From the definition  $\text{rebut}(c_{\text{for}}, c_{\text{against}}) = c_{\text{for}} / (c_{\text{for}} + c_{\text{against}})$ :  
 $\text{rebut}(c_{\text{for}}, c_{\text{against}}) + \text{rebut}(c_{\text{against}}, c_{\text{for}})$

$$\frac{c\_for}{c\_for + c\_against} + \frac{c\_against}{c\_against + c\_for}$$

$$\frac{(c\_for + c\_against)}{(c\_for + c\_against)}$$

1 ■

■

**C.9 (Rebut Monotonicity)** Rebut is monotone in the first argument and antitone in the second: **If  $c\_1 \leq c\_2$ , then  $\text{rebut}(c\_1, c) \leq \text{rebut}(c\_2, c)$**  **If  $d\_1 \leq d\_2$ , then  $\text{rebut}(c, d\_2) \leq \text{rebut}(c, d\_1)$**

*Proof.* For the first claim:  $\text{rebut}(c\_1, c) = c\_1 / (c\_1 + c) \leq c\_2 / (c\_2 + c) = \text{rebut}(c\_2, c)$  since the function  $f(x) = x / (x + c)$  is increasing in  $x$  for  $c > 0$ .

The second claim follows similarly by noting  $\text{rebut}(c, d) = 1 - \text{rebut}(d, c)$  from Theorem C.8. ■ ■

### C.3.4 Interaction Between Undercut and Rebut

Undercut and rebut represent two fundamentally different types of defeat: **Undercut attacks the inferential link between premises and conclusion** **Rebut attacks the conclusion directly with counter-evidence**

The interaction of these two operations is subtle and context-dependent. In general, rebut is applied first to aggregate competing evidence, then undercut is applied to discount the resulting confidence based on link attackers. This ordering is reflected in CLAIR’s evaluation semantics (Chapter 4).

### C.3.5 Limitation: Rebut Normalization

A key limitation of the rebut operation is that it normalizes away absolute evidence strength.

**C.10 (Rebut Collapse)** For any scaling factor  $\lambda > 0$ :  $\text{rebut}(\lambda c\_for, \lambda c\_against) = \text{rebut}(c\_for, c\_against)$

*Proof.*  $\text{rebut}(\lambda c\_for, \lambda c\_against) = \lambda c\_for / (\lambda c\_for + \lambda c\_against)$

$$\frac{\lambda c\_for}{\lambda (c\_for + c\_against)}$$

$$\frac{c\_for}{c\_for + c\_against}$$

$$\text{rebut}(c\_for, c\_against). \blacksquare$$

■

This means rebut cannot distinguish between “both weak” and “both strong but balanced” evidence. Chapter 6 discusses how CLAIR addresses this through provenance tracking, which preserves the absolute strengths of individual evidence sources even after rebut normalization.

## Glossary

This appendix provides concise definitions of key terms, notation, and acronyms used throughout this dissertation. Terms are marked with their primary chapter of introduction.

## D.1 Term Definitions

### Epistemic Terms

+— +Term | Definition | Chapter +— +Belief | A first-class value in CLAIR consisting of a proposition, confidence, and justification. Beliefs can be constructed, combined, defeated, and inspected through language operations. | 1 +Confidence | A value in the unit interval  $[0,1]$  representing degree of epistemic commitment. Confidence is *not* probability: rather than quantifying frequency or subjective chance, confidence tracks how justified a belief is given available evidence. | 3 +Commitment (Epistemic) | The stance of endorsing a proposition as true to a specified degree. High confidence (close to 1) indicates strong commitment; low confidence (close to 0) indicates weak or no commitment. | 3 +Justification | A data structure tracking the derivation history of a belief. Justifications form a directed acyclic graph (DAG) where nodes represent beliefs and edges represent inference rules or evidence sources. | 4 +Invalidation | A condition associated with a belief specifying circumstances under which the belief should be reconsidered. Invalidation enables defeasible reasoning—beliefs can be defeated without logical contradiction. | 4 +Provenance | The origin or source of evidence for a belief. CLAIR tracks provenance through justification graphs, enabling audit trails for how conclusions were reached. | 4 +Reliability | In CLAIR’s semantics, the tendency of a source to produce true beliefs. Reliability is the *semantic interpretation* of confidence: a belief with confidence  $c$  is interpreted as coming from a source with reliability  $c$ . | 3 +—

### Operations and Relations

+— +Term | Definition | Chapter +— +Probabilistic OR ( $\oplus$ ) | Aggregation operator for independent evidence:  $\mathbf{a} \oplus \mathbf{b} = \mathbf{a} + \mathbf{b} - \mathbf{ab}$ . This equals the probability that at least one of two independent events occurs. Satisfies commutativity, associativity, and has identity 0. | 3 +Undercut ( $\neg\rightarrow$ ) | Defeat operation that reduces confidence multiplicatively:  $\mathbf{undercut}(\mathbf{c}, \mathbf{d}) = \mathbf{c} \times (\mathbf{1}-\mathbf{d})$ . A defeater with confidence  $\mathbf{d}$  reduces target confidence by factor  $(\mathbf{1}-\mathbf{d})$ . | 4 +Rebuttal | Defeat operation for conflicting evidence:  $\mathbf{rebut}(\mathbf{c}_1, \mathbf{c}_2) = \mathbf{c}_1 / (\mathbf{c}_1 + \mathbf{c}_2)$ . Normalizes competing confidences to  $[0,1]$ ; equal confidences yield  $1/2$ . | 4 +Derivation | Combining beliefs via rule application:  $\mathbf{derive}(\mathbf{b}_1, \mathbf{b}_2)$  pairs the values and multiplies confidences  $\mathbf{c}_1 \times \mathbf{c}_2$ . Tracks justification through rule application. | 4 +Aggregation | Combining independent evidence using  $\oplus$ :  $\mathbf{aggregate}(\mathbf{b}_1, \mathbf{b}_2)$  produces belief with confidence  $\mathbf{c}_1 \oplus \mathbf{c}_2$ . Tracks justification through aggregation node. | 4 +Subtyping | Confidence ordering: belief at confidence  $\mathbf{c}_1$  can be used where belief at  $\mathbf{c}_2$  is required iff  $\mathbf{c}_1 \geq \mathbf{c}_2$ . Enables confidence weakening. | 10 +—

### Structural Properties

+— +Term | Definition | Chapter +— +Directed Acyclic Graph (DAG) | A graph with directed edges and no cycles. CLAIR requires justification graphs to be DAGs to ensure well-foundedness and prevent circular reasoning. | 4 +Stratification | Layering beliefs into levels  $\mathbf{0}, \mathbf{1}, \mathbf{2}, \dots$  where level  $\mathbf{n}$  can only refer to levels  $<\mathbf{n}$ . Enforces Tarski’s hierarchy to prevent self-referential paradoxes. | 6 +Well-formedness | Constraints on CLAIR programs: (1) justification graphs must be acyclic, (2) stratification constraints on introspection, (3) confidences in  $[0,1]$ . | 10 +—

### Logical and Modal Terms

+— +Term | Definition | Chapter +— +CPL (Confidence-Bounded Provability Logic) | Graded extension of Gödel-Löb provability logic. Adds confidence grades to provability operator  $\Box$ . Axiomatizes self-referential reasoning with confidence discounts. | 5 +Löb’s Theorem | Modal logic theorem:  $\Box(\Box\phi \rightarrow \phi) \rightarrow \Box\phi$ . In provability logic, enables self-reference; in CLAIR, motivates anti-bootstrapping constraint. | 5 +Graded Modality | Modal operators with quantitative gradess. CLAIR’s  $\Box_c$  means “provable with confidence at least  $c$ .” | 5 +Anti-bootstrapping |

Principle that self-validating claims cannot increase confidence. Formally:  $c \leq g(c)$  for some discount function  $g$ . CLAIR uses  $g(c) = c^2$ . | 5 +Kripke Semantics | Possible worlds framework for modal logic. CPL uses graded Kripke models where accessibility relations track confidence thresholds. | 5 +—

## Computational Terms

+— +Term | Definition | Chapter +— +Fuel | Bound on computation steps in the Lean evaluator: **evalFuel n e** evaluates for at most  $n$  steps. Prevents infinite loops while ensuring termination. | B +Call-by-value | Evaluation strategy: function arguments are evaluated before application. CLAIR uses call-by-value to match intuition about belief formation. | B +Small-step semantics | Reduction relation  $e \rightarrow e'$  defining one step of computation. Composed to form multi-step evaluation  $e \rightarrow e'$ . | 10 +Bidirectional Type Checking | Type checking algorithm with synthesis (infer type from expression) and checking (verify expression matches type). Enables practical implementation. | 10 +De Bruijn Indices | Variable representation using natural numbers: var 0 = most recent binder, var 1 = next recent, etc. Enables formal proofs about substitution. | A +—

## Argumentation and Belief Revision

+— +Term | Definition | Chapter +— +Defeasible Reasoning | Reasoning where conclusions can be defeated by new information without logical contradiction. CLAIR supports defeasibility through undercut and rebut operations. | 4 +Underminer | An argument that attacks the connection between evidence and conclusion (e.g., “the sensor is miscalibrated”). Reduces confidence via **undercut(c,d) =  $c \times (1-d)$** . | 4 +Rebuttal | An argument providing conflicting evidence for the opposite conclusion. Normalizes via **rebut(c1,c2) =  $c1/(c1+c2)$** . | 4 +AGM Theory | Classic belief revision framework (Alchourrón, Gärdenfors, Makinson). CLAIR extends AGM to graded, DAG-structured beliefs. | 7 +Contraction | Belief revision operation: remove a belief from a belief set. In CLAIR, achieved by setting confidence to 0. | 7 +Revision | Belief revision operation: add a belief while maintaining consistency. In CLAIR, achieved by aggregating with existing beliefs. | 7 +—

## Impossibility Results

+— +Term | Definition | Chapter +— +Gödel’s Incompleteness | Any consistent formal system capable of arithmetic contains true but unprovable statements. Motivates CPL’s design restrictions. | 12 +Church’s Undecidability | First-order logic validity is undecidable. CLAIR restricts to decidable fragments (CPL-finite, CPL-0). | 12 +Tarski’s Undefinability | Truth cannot be defined within the same language. Motivates stratification: level  $n$  cannot quantify over level  $n$ . | 12 +Löb’s Paradox | Curious proposition using Löb’s theorem that leads to contradiction if not carefully restricted. Resolved via stratification. | 6 +—

## D.2 Notation Table

+— +Symbol | Meaning | Type | Chapter +— +**c** | Confidence value in  $[0,1]$  | Confidence | 3 +⊕ | Probabilistic OR:  $a \oplus b = a + b - ab$  | Binary operation | 3 +— | Undercut:  $c \multimap d = c \times (1-d)$  | Binary operation | 4 +× | Multiplication (conjunctive combination) | Binary operation | 3 +**g(c)** | Löb discount function; CLAIR uses  $c^2$  | Unary function | 5 +□<sub>c</sub> φ | Necessarily φ with confidence at least c | Modal operator | 5 +◇<sub>c</sub> φ | Possibly φ with confidence at least c | Modal operator | 5 +⊢ | Typing judgment:  $\Gamma \vdash e : A @ c$  | Relation | 10 +→ | Small-step reduction:  $e \rightarrow e'$  | Relation | 10 +→ | **Multi-step reduction (reflexive transitive closure) | Relation | 10 +Γ | Typing context (list of variable: type pairs) | Context | 10 +A ⇒ B | Function type from A to B | Type | 10 +Belief<A> | Belief type holding value of type A | Type constructor | 10 +b.value | Extract value from belief b | Projection | 4**

+b.confidence | Extract confidence from belief b | Projection | 4 +b.justification | Extract justification from belief b | Projection | 4 +m <n | Stratification constraint: level m below n | Ordering | 6 +∀ | Universal quantifier | Quantifier | Appendix A +∃ | Existential quantifier | Quantifier | Appendix A +∈ | Set membership | Relation | Appendix A +⊆ | Subset relation | Relation | Appendix A +∪ | Set union | Binary operation | Appendix A +∩ | Set intersection | Binary operation | Appendix A +λx:A. e | Lambda abstraction (anonymous function) | Expression | 10 +e\*1 e\*2 | Function application | Expression | 10 +—

## D.3 Acronyms

+— +Acronym | Full Name | Definition +— +CLAIR | Comprehensible LLM AI Intermediate Representation | The formal system presented in this dissertation +CPL | Confidence-Bounded Provability Logic | Graded modal logic extending Gödel-Löb +CPL-finite | CPL with finite confidence set {0, 1/n, 2/n, ..., 1} | Decidable fragment suitable for implementation +CPL-o | CPL with only confidence o (ungraded provability) | Collapses to standard provability logic GL +AGM | Alchourrón-Gärdenfors-Makinson | Classic belief revision theory +DAG | Directed Acyclic Graph | Graph structure for justifications +LLM | Large Language Model | AI systems CLAIR is designed to augment +IR | Intermediate Representation | Compiler representation between source and machine code +AI | Artificial Intelligence | Field of study +—

## D.4 Type System Summary

### Base Types

+— +Type | Description | Example Values +— +Nat | Natural numbers (non-negative integers) | 0, 1, 2, 42, 128 +Bool | Boolean values | true, false +String | Text strings | “hello”, “sensor-reading” +Unit | Unit type (single value) | unit +Pair(A, B) | Ordered pair | (1, true), (“x”, 42) +Sum(A, B) | Tagged union (either A or B) | inl(5), inr(true) +A ⇒ B | Function type from A to B | λx:Nat. x + 1 +Belief<A> | Belief holding value of type A | belief(42, 0.9, j) +—

### Confidence Operations

+— +Operation | Notation | Formula | Identity +— +Probabilistic OR |  $\mathbf{a} \oplus \mathbf{b}$  |  $\mathbf{a} + \mathbf{b} - \mathbf{a} \times \mathbf{b}$  | 0 +Multiplication |  $\mathbf{a} \times \mathbf{b}$  |  $\mathbf{a} \times \mathbf{b}$  | 1 +Undercut | **undercut**(c, d) |  $\mathbf{c} \times (\mathbf{1} - \mathbf{d})$  | N/A (d=0 gives c) +Rebuttal | **rebut**(c1, c2) |  $\mathbf{c}_1 / (\mathbf{c}_1 + \mathbf{c}_2)$  | N/A +Minimum | **min**(a, b) | **min**(a, b) | 0 (if bounded below) +Maximum | **max**(a, b) | **max**(a, b) | 1 (if bounded above) +—

# Complete CLAIR Language Specification

This appendix provides the complete formal specification of the CLAIR language, including syntax (concrete and abstract), static semantics (type system), and dynamic semantics (operational rules). The specification is self-contained and sufficient for implementing a conforming CLAIR interpreter.

## E.1 Syntax

### E.1.1 Type Grammar

The type grammar defines the set of valid types in CLAIR.

+— +Category | Production | Description Base Types | B ::= “Nat” | Natural numbers | | B ::= “Bool” | Boolean values | | B ::= “String” | Text strings | | B ::= “Unit” | Unit type (single value) Confidence | c ∈ Q | Rational in [0,1] Types | A ::= B | Base type | | A ::= A → B | Function type

| |  $A ::= A \times B$  | Product type | |  $A ::= A + B$  | Sum type | |  $A ::= \text{“Belief}<”A[“c”]}$  | Belief type with confidence bound | |  $A ::= \text{“Meta}<”A[“n”][“c”]}$  | Stratified meta-belief type +—

### E.1.2 Expression Grammar

The expression grammar defines the syntactic forms of CLAIR programs.

+— +Category | Production | Description Variables |  $e ::= x$  | Variable reference Lambdas |  $e ::= \lambda x. A$  | Anonymous function Application |  $e ::= e e$  | Function application Pairs |  $e ::= (e, e)$  | Ordered pair Projections |  $e ::= e.1$  | First projection | |  $e ::= e.2$  | Second projection Injections |  $e ::= \text{inl}_B(e)$  | Left injection | |  $e ::= \text{inr}_A(e)$  | Right injection Case |  $e ::= \text{case } e \text{ of } x \Rightarrow e | \text{inr } y \Rightarrow e$  | Sum elimination Literals |  $e ::= n$  | Natural number literal | |  $e ::= \text{“true”} | \text{“false”}$  | Boolean literals | |  $e ::= \text{“s”}$  | String literal | |  $e ::= ()$  | Unit literal Beliefs |  $e ::= \text{belief}(e, c, j)$  | Belief constructor | |  $e ::= \text{val}(e)$  | Extract belief value | |  $e ::= \text{conf}(e)$  | Extract belief confidence | |  $e ::= \text{just}(e)$  | Extract belief justification Derivation |  $e ::= \text{derive}(e, e)$  | Combine beliefs (conjunctive) Aggregation |  $e ::= \text{aggregate}(e, e)$  | Aggregate beliefs (independent) Defeat |  $e ::= \text{undercut}(e, e)$  | Apply undercut | |  $e ::= \text{rebut}(e, e)$  | Apply rebuttal Introspection |  $e ::= \text{introspect}(e)$  | Safe self-reference Let Binding |  $e ::= \text{let } x = e \text{ in } e$  | Local binding +—

### E.1.3 Abstract Syntax

The abstract syntax uses de Bruijn indices for variable representation.

**Definition** The abstract syntax “Expr” is defined inductively with the following constructors:

“var( $n$ )” — Variable at de Bruijn index  $n$  “lam( $A, e$ )” — Lambda abstraction “ $\lambda. A$ .”  $e$  “app( $e_1, e_2$ )” — Function application  $e_1 e_2$  “pair( $e_1, e_2$ )” — Ordered pair “fst( $e$ )” — First projection “snd( $e$ )” — Second projection “inl( $B, e$ )” — Left injection “inr( $A, e$ )” — Right injection “case( $e, e_1, e_2$ )” — Case analysis “litNat( $n$ )” — Natural literal “litBool( $b$ )” — Boolean literal “litString( $s$ )” — String literal “litUnit” — Unit literal “belief( $v, c, j$ )” — Belief constructor “val( $e$ )” — Value extraction “conf( $e$ )” — Confidence extraction “just( $e$ )” — Justification extraction “derive( $e_1, e_2$ )” — Derivation “aggregate( $e_1, e_2$ )” — Aggregation “undercut( $e, d$ )” — Undercut “rebut( $e_1, e_2$ )” — Rebuttal “introspect( $e$ )” — Introspection “letIn( $e_1, e_2$ )” — Let binding

The “Justification” type tracks derivation structure: “axiomJ( $\text{name}$ )” — Named axiom “rule( $\text{name}, \text{js}$ )” — Named rule with premises “agg( $\text{js}$ )” — Aggregation “undercut\_j( $j, d$ )” — Undercut application “rebut\_j( $j_1, j_2$ )” — Rebuttal application

### E.1.4 Well-Formedness

A type is well-formed if all confidence bounds are in  $[0,1]$ .

**Definition** The judgment “wf( $A$ )” holds when:

“wf( $B$ )” for all base types  $B$  “wf( $A$ )”  $\wedge$  “wf( $B$ )”  $\Rightarrow$  “wf( $A \rightarrow B$ )” “wf( $A$ )”  $\wedge$  “wf( $B$ )”  $\Rightarrow$  “wf( $A \times B$ )” “wf( $A$ )”  $\wedge$  “wf( $B$ )”  $\Rightarrow$  “wf( $A + B$ )” “wf( $A$ )”  $\wedge c \in [0,1] \Rightarrow$  “wf(Belief<” $A[“c”]$ )” “wf( $A$ )”  $\wedge c \in [0,1] \Rightarrow$  “wf(Meta<” $A[“n”][“c”]$ )”

## E.2 Static Semantics (Type System)

### E.2.1 Typing Contexts

A typing context  $\Gamma$  maps variable indices to type-confidence pairs.

$\Gamma ::= \emptyset \mid \Gamma, \langle A, c \rangle$

**Definition** The lookup operation  $\Gamma.\text{lookup}(n)$  returns the type-confidence pair at index  $n$  (0-indexed from most recent binding).

### E.2.2 Typing Judgment Form

The primary typing judgment has the form:

$\Gamma \vdash e : A @ c$

“In context  $\Gamma$ , expression  $e$  has type  $A$  with confidence bound  $c$ .”

## E.2.3 Typing Rules

The following rules define well-typed CLAIR expressions. Confidence propagation is explicit in each rule.

**T-Variable:**  $\Gamma \vdash \text{“var}(n)\text{”} \text{“:”} A \text{“@”} c$  if  $\Gamma.\text{lookup}(n) = \langle A, c \rangle$

**T-Abstraction:** If  $\Gamma \text{“,”} \langle A, c_A \rangle \vdash e \text{“:”} B \text{“@”} c_B$ , then  $\Gamma \vdash \text{“lam}(A, e)\text{”} \text{“:”} (A \text{“} \rightarrow \text{”} B) \text{“@”} c_B$

**T-Application:** If  $\Gamma \vdash e_1 \text{“:”} (A \text{“} \rightarrow \text{”} B) \text{“@”} c_1$  and  $\Gamma \vdash e_2 \text{“:”} A \text{“@”} c_2$ , then  $\Gamma \vdash \text{“app}(e_1, e_2)\text{”} \text{“:”} B \text{“@”} (c_1 \times c_2)$

**Confidence multiplies for conjunctive derivation.**

**T-Pair:** If  $\Gamma \vdash e_1 \text{“:”} A \text{“@”} c_1$  and  $\Gamma \vdash e_2 \text{“:”} B \text{“@”} c_2$ , then  $\Gamma \vdash \text{“pair}(e_1, e_2)\text{”} \text{“:”} (A \text{“} \times \text{”} B) \text{“@”} (c_1 \times c_2)$

**T-Fst/T-Snd** preserve confidence.

**T-Inl/T-Inr** preserve confidence.

**T-Case:** If  $\Gamma \vdash e \text{“:”} (A \text{“} + \text{”} B) \text{“@”} c$  and  $\Gamma \text{“,”} \langle A, c \rangle \vdash e_1 \text{“:”} C \text{“@”} c_1$  and  $\Gamma \text{“,”} \langle B, c \rangle \vdash e_2 \text{“:”} C \text{“@”} c_2$ , then  $\Gamma \vdash \text{“case}(e, e_1, e_2)\text{”} \text{“:”} C \text{“@”} (c \oplus c_1 \oplus c_2)$

**Uses probabilistic OR for confidence aggregation.**

**T-Belief:** If  $\Gamma \vdash v \text{“:”} A \text{“@”} 1$  and  $c_{\text{bound}} \leq c_{\text{actual}}$ , then  $\Gamma \vdash \text{“belief}(v, c_{\text{actual}}, j)\text{”} \text{“:”} \text{“Belief}<\text{”} A \text{“@”} c_{\text{bound}}$

The belief’s actual confidence  $c_{\text{actual}}$  must meet or exceed the declared bound  $c_{\text{bound}}$ .

**T-Val/T-Conf/T-Just** preserve the enclosing confidence.

**T-Derive:** If  $\Gamma \vdash e_1 \text{“:”} \text{“Belief}<\text{”} A \text{“@”} c_{e_1}$  and  $\Gamma \vdash e_2 \text{“:”} \text{“Belief}<\text{”} B \text{“@”} c_{e_2}$ , then  $\Gamma \vdash \text{“derive}(e_1, e_2)\text{”} \text{“:”} \text{“Belief}<\text{”} A \text{“} \times \text{”} B \text{“@”} (c_{e_1} \times c_{e_2})$

**Confidence multiplies: both premises must be true.**

**T-Aggregate:** If  $\Gamma \vdash e_1 \text{“:”} \text{“Belief}<\text{”} A \text{“@”} c_{e_1}$  and  $\Gamma \vdash e_2 \text{“:”} \text{“Belief}<\text{”} A \text{“@”} c_{e_2}$ , then  $\Gamma \vdash \text{“aggregate}(e_1, e_2)\text{”} \text{“:”} \text{“Belief}<\text{”} A \text{“@”} (c_{e_1} \oplus c_{e_2})$

**Uses probabilistic OR: independent evidence.**

**T-Undercut:** If  $\Gamma \vdash e \text{“:”} \text{“Belief}<\text{”} A \text{“@”} c_e$  and  $\Gamma \vdash d \text{“:”} \text{“Belief}<\text{”} d_c \text{“@”} c_d$ , then  $\Gamma \vdash \text{“undercut}(e, d)\text{”} \text{“:”} \text{“Belief}<\text{”} A \text{“@”} (c_e \times c_d)$  where  $\text{“undercut}(c, d)\text{”} = c \times (1 - d)$ .

**T-Rebut:** If  $\Gamma \vdash e_{\text{for}} \text{“:”} \text{“Belief}<\text{”} A \text{“@”} c_{e_1}$  and  $\Gamma \vdash e_{\text{against}} \text{“:”} \text{“Belief}<\text{”} A \text{“@”} c_{e_2}$ , then  $\Gamma \vdash \text{“rebut}(e_{\text{for}}, e_{\text{against}})\text{”} \text{“:”} \text{“Belief}<\text{”} A \text{“@”} (c_{e_1} \times c_{e_2})$  where  $\text{“rebut}(c_{\text{for}}, c_{\text{against}})\text{”} = c_{\text{for}}$

**T-Introspect:** If  $\Gamma \vdash e \text{“:”} \text{“Meta}<\text{”} A \text{“@”} c_e$  and  $m < n$ , then:

The resulting type is:  $\text{“Meta}<\text{”} \text{“Meta}<\text{”} A \text{“@”} [m][c] > [n][g(c)]$  (a meta-belief about a meta-belief).

Thus:  $\Gamma \vdash \text{“introspect}(e)\text{”} \text{“:”} \text{“Meta}<\text{”} \text{“Meta}<\text{”} A \text{“@”} [m][c] > [n][g(c)]$

where  $g(c) = c^2$  is the Löb discount function (to prevent bootstrapping).

**Requires level constraint:** only higher levels can introspect lower levels.

## E.2.4 Subtyping

CLAIR supports subtyping based on confidence bounds.

**S-Belief:**  $\text{“Belief}<\text{”} A \text{“@”} c < \text{“Belief}<\text{”} A \text{“@”} c'$  if  $c \geq c'$

Higher confidence is a subtype of lower confidence.

**S-Meta** follows the same rule.

**T-Sub:** If  $\Gamma \vdash e \text{“:”} A \text{“@”} c$  and  $A < A'$  and  $c \geq c'$ , then  $\Gamma \vdash e \text{“:”} A' \text{“@”} c'$

Allows weakening both type and confidence.



## E.3 Dynamic Semantics

### E.3.1 Values

A value is a fully evaluated expression.

**Definition** The predicate “ $\text{IsValue}(e)$ ” holds for:

All lambda abstractions “ $\text{lam}(A, e)$ ” All pairs “ $\text{pair}(v_1, v_2)$ ” where  $v_1, v_2$  are values All injections “ $\text{inl}(B, v)$ ” and “ $\text{inr}(A, v)$ ” where  $v$  is a value All literals (“ $\text{litNat}$ ”, “ $\text{litBool}$ ”, “ $\text{litString}$ ”, “ $\text{litUnit}$ ”) All belief constructors “ $\text{belief}(v, c, j)$ ” where  $v$  is a value

### E.3.2 Small-Step Operational Semantics

The small-step relation  $e \rightarrow e'$  defines single-step reduction using call-by-value evaluation order.

**Beta Reduction:** “ $(\lambda : A. e) v \rightarrow e[x := v]$ ” if “ $\text{IsValue}(v)$ ”

**Context Rules:** “ $\text{app}(e_1, e_2) \rightarrow \text{app}(e_1', e_2)$ ” if  $e_1 \rightarrow e_1'$  “ $\text{app}(v_1, e_2) \rightarrow \text{app}(v_1, e_2')$ ” if  $e_2 \rightarrow e_2'$  and “ $\text{IsValue}(v_1)$ ”

**Projections:** “ $\text{fst}(\text{pair}(v_1, v_2)) \rightarrow v_1$ ” if “ $\text{IsValue}(v_1)$ ”, “ $\text{IsValue}(v_2)$ ” “ $\text{snd}(\text{pair}(v_1, v_2)) \rightarrow v_2$ ” if “ $\text{IsValue}(v_1)$ ”, “ $\text{IsValue}(v_2)$ ”

Context rules for “ $\text{fst}$ ” and “ $\text{snd}$ ” evaluate subexpressions first.

**Case Analysis:** When “ $\text{case}(\text{inl}(v), e_1, e_2)$ ” evaluates and  $v$  is a value, it reduces to “ $e_1[x := v]$ ”. When “ $\text{case}(\text{inr}(v), e_1, e_2)$ ” evaluates and  $v$  is a value, it reduces to “ $e_2[y := v]$ ”.

Context rule for “ $\text{case}$ ” evaluates the scrutinee first.

Context rule for “ $\text{case}$ ” evaluates the scrutinee first.

**Let Binding:** “ $\text{letIn}(v, e) \rightarrow e[x := v]$ ” if “ $\text{IsValue}(v)$ ”

Context rule for “ $\text{letIn}$ ” evaluates the binding first.

**Belief Projections:** “ $\text{val}(\text{belief}(v, c, j)) \rightarrow v$ ” if “ $\text{IsValue}(v)$ ” “ $\text{conf}(\text{belief}(v, c, j)) \rightarrow \text{belief}(v, c, j)$ ” if “ $\text{IsValue}(v)$ ” “ $\text{just}(\text{belief}(v, c, j)) \rightarrow \text{litString}(\text{toString}(j))$ ” if “ $\text{IsValue}(v)$ ”

Context rules evaluate subexpressions to values first.

**Derivation, Aggregation, Defeat:**

These operations evaluate to values using the confidence operations defined in the typing rules. The small-step semantics provides context rules that evaluate both subexpressions to values before computing the result.

For example, when “ $\text{derive}(v_1, v_2)$ ” has both subexpressions as values, the evaluator computes a new belief with: **Value:** “ $\text{pair}(v_1.\text{value}, v_2.\text{value})$ ” **Confidence:** “ $v_1.\text{confidence} \times v_2.\text{confidence}$ ” **Justification:** “ $\text{rule}(\text{"derive"}, [v_1.\text{just}, v_2.\text{just}])$ ”

**Introspection:** “ $\text{introspect}(v) \rightarrow v$ ” if “ $\text{IsValue}(v)$ ”

Context rule evaluates subexpression first.

### E.3.3 Multi-Step Reduction

The multi-step reduction relation  $e \rightarrow^* e'$  is the reflexive-transitive closure of “ $\rightarrow$ ”.

$e \rightarrow^* e'$  if  $e = e'$  or  $e \rightarrow e''$  and  $e'' \rightarrow^* e'$

### E.3.4 Evaluation Function

The evaluation function “ $\text{eval}(e)$ ” returns the result of reducing  $e$  to a value, or fails if:

1. The expression gets stuck (no applicable rule)
2. The expression exceeds the fuel bound (default: 1000 steps)

**Definition** “ $\text{eval}(e) = \text{some}(v)$ ” if  $e \rightarrow^* v$  and “ $\text{IsValue}(v)$ ” “ $\text{eval}(e) = \text{none}$ ” otherwise (stuck or out of fuel)

## E.4 Well-Formedness Constraints

### E.4.1 Acyclicity of Justification Graphs

For deterministic evaluation in the DAG semantics, justification graphs must be acyclic.

**Definition** A justification graph  $G = (V, E)$  is **well-formed** iff:

1. For all nodes  $v \in V$ , the transitive closure of  $E$  starting from  $v$  contains no cycles.
2. Equivalently: there is no path  $v \rightarrow^* v$  for any  $v \in V$ .

**Enforcement:** In the reference interpreter, this is checked during type checking by tracking provenance and ensuring no belief can transitively depend on itself.

### E.4.2 Stratification Constraints

For safe introspection, meta-belief levels must form a strict hierarchy.

**Definition** The **stratification constraint** requires:

1. Every “introspect(e)” operation has proof  $m \prec n$  where:  **$m$  is the level of the source belief  $n$  is the level of the resulting meta-belief**
2. This is enforced at compile time: the type checker requires a proof term of type  $m \prec n$ .

**Consequence:** No belief can introspect itself or any belief that transitively introspects it.

### E.4.3 Confidence Bounds

All confidence values must satisfy  $0 \leq c \leq 1$ .

**Enforcement:** The type system checks this at all belief construction points.

## E.5 Example Programs

**Basic Belief** Basic Belief:

```
belief(42, 0.9, axiomJ("sensor-reading"))
```

A belief in the value 42 with confidence 0.9, traced to a sensor reading axiom.

**Conjunctive Derivation** Derivation:

```
let p = belief("it is raining", 0.8, axiomJ("weather-report")) in
let q = belief("I have an umbrella", 0.9, axiomJ("visual-check")) in
derive(p, q)
```

Combines two beliefs by multiplication, yielding confidence 0.72.

**Aggregation** Independent Evidence Aggregation:

```
let e1 = belief("stock will rise", 0.6, axiomJ("analyst-a")) in
let e2 = belief("stock will rise", 0.7, axiomJ("analyst-b")) in
aggregate(e1, e2)
```

Uses probabilistic OR:  $0.6 \oplus 0.7 = 0.6 + 0.7 - 0.6 \times 0.7 = 0.88$ .

**Undercut** Defeat:

```
let claim = belief("system is secure", 0.95, axiomJ("vendor-claim")) in
let vuln = belief("critical CVE found", 0.8, axiomJ("security-audit")) in
undercut(claim, vuln)
```

Applies undercut:  $0.95 \times (1 - 0.8) = 0.19$ .

## E.6 Summary

This specification provides:

1. **Complete type grammar:** Base types, functions, products, sums, beliefs, and meta-beliefs
2. **Complete expression grammar:** All CLAIR syntactic forms
3. **Static semantics:** 20 typing rules covering all constructs
4. **Dynamic semantics:** Small-step operational semantics with call-by-value evaluation
5. **Well-formedness constraints:** Acyclicity, stratification, and confidence bounds

The specification is sufficient for implementing a conforming CLAIR interpreter and type checker. The Lean 4 formalization in Appendix A provides a machine-checked version of this specification.