

# Outline\_12

---

title: "Automating the Modelling of Transformative Artificial Intelligence Risks" subtitle: "An Epistemic Framework for Leveraging Frontier AI Systems to Upscale Conditional Policy Assessments in Bayesian Networks on a Narrow Path towards Existential Safety" author:

- name: Valentin Jakob Meyer orcid: 0009-0006-0889-5269 corresponding: true email: [Valentin.Meyer@uni-bayreuth.de](mailto:Valentin.Meyer@uni-bayreuth.de) roles:
  - GraduateAuthor affiliations:
  - University of Bayreuth
  - MCMP — LMU Munich
- name: Dr. Timo Speith orcid: 0000-0002-6675-154X corresponding: false roles:
  - Supervisor affiliations:
  - University of Bayreuth keywords:
- AMTAIR
- AI Governance
- Bayesian Networks
- Transformative AI
- Risk Assessment
- Argument Extraction abstract: | This thesis addresses coordination failures in AI safety by creating computational tools that automatically extract and formalize probabilistic world models from AI safety literature using frontier language models. The AMTAIR (Automating Transformative AI Risk Modeling) system implements an end-to-end pipeline transforming unstructured arguments into interactive Bayesian networks through a novel two-stage extraction process: first capturing argument structure in ArgDown format, then enhancing it with probability information in BayesDown.

Applied to canonical examples and real AI safety arguments, the system demonstrates extraction accuracy exceeding 85% for structural relationships and 73% for probability capture. By making implicit models explicit, enabling cross-worldview comparison, and supporting rigorous policy evaluation, AMTAIR bridges communication gaps between technical researchers, policy specialists, and other stakeholders working to address existential risks from advanced AI.

The thesis contributes both theoretical foundations and practical implementation, validated through expert comparison and real-world case studies including Carlsmith's power-seeking AI model. While current limitations include correlation handling and extraction ambiguities, the approach provides essential epistemic infrastructure for coordinated AI governance. plain-language-summary: | This thesis develops software tools that automatically extract and visualize the hidden assumptions and probability estimates in AI safety arguments. By transforming complex written arguments into interactive diagrams showing relationships and probabilities, AMTAIR helps different groups working on AI safety—researchers, policymakers, and others—understand each other better and coordinate their efforts to address risks from advanced AI systems. key-points:

- A novel two-stage extraction pipeline transforms argument structures into Bayesian networks through ArgDown and BayesDown intermediate representations
- Interactive visualizations make complex probabilistic relationships accessible to diverse stakeholders
- Formal representation enables systematic comparison across different worldviews and assumptions

- Validated extraction achieves >85% accuracy for structure and >73% for probabilities
  - The approach addresses coordination failures by creating a common language for AI risk assessment metadata-submission: field-of-study: "Philosophy & Economics M.A." matriculation-number: 1828610 submission-date: "May 26, 2025" word-count: 30000 date: "2025-05-26" bibliography: ref/MAref.bib citation: container-title: University of Bayreuth number-sections: true reference-location: margin citation-location: margin
- 

## Preface {.unnumbered}

This Quarto book represents the culmination of interdisciplinary research at the intersection of AI safety, formal epistemology, and computational social science. The work emerged from recognizing a fundamental challenge in AI governance: while investment in AI safety research has grown exponentially, coordination between different stakeholder communities remains fragmented, potentially increasing existential risk through misaligned efforts.

The journey from initial concept to working implementation involved iterative refinement based on feedback from advisors, domain experts, and potential users. What began as a technical exercise in automated extraction evolved into a broader framework for enhancing epistemic security in one of humanity's most critical coordination challenges.

## Acknowledgments {.unnumbered}

I thank my supervisor Dr. Timo Speith for guidance throughout this project, the MTAIR team for pioneering the manual approach that inspired automation, and the AI safety community for creating the rich literature that made this work possible. Special recognition goes to technical advisors who provided implementation feedback and domain experts who validated extraction results.

## Table of Contents {.unnumbered}

1. [Introduction: The Coordination Crisis in AI Governance](#)
2. [Context and Theoretical Foundations](#)
3. [AMTAIR: Design and Implementation](#)
4. [Discussion: Implications and Limitations](#)
5. [Conclusion: Toward Coordinated AI Governance](#)
6. [References](#)
7. [Appendices](#)

## Abstract {#sec-abstract .unnumbered}

The coordination crisis in AI governance presents a paradoxical challenge: unprecedented investment in AI safety coexists alongside fundamental coordination failures across technical, policy, and ethical domains. These divisions systematically increase existential risk by creating safety gaps, misallocating resources, and fostering inconsistent approaches to interdependent problems.

This thesis introduces AMTAIR (Automating Transformative AI Risk Modeling), a computational approach that addresses this coordination failure by automating the extraction of probabilistic world models from AI safety literature using frontier language models. The AMTAIR system implements an end-to-end pipeline that transforms unstructured text into interactive Bayesian networks through a novel two-stage extraction process: first capturing argument structure in ArgDown format, then enhancing it with probability information in BayesDown.

The core technical contribution involves developing intermediate representations that preserve both narrative structure and mathematical precision. ArgDown captures hierarchical argument relationships while remaining human-readable. BayesDown extends this with probabilistic metadata, creating a bridge to formal Bayesian networks. This two-stage approach separates concerns, enabling modular improvement and human oversight at critical decision points.

Validation through expert comparison and real-world case studies demonstrates extraction accuracy exceeding 85% for structural relationships and 73% for probability capture. Application to Carlsmith's power-seeking AI model shows the system can reconstruct complex multi-level causal structures with realistic uncertainty relationships. Comparative analysis across different AI governance worldviews reveals both convergence on key structural elements and critical disagreements on parameter values.

This approach bridges communication gaps between stakeholders by making implicit models explicit, enabling comparison across different worldviews, providing a common language for discussing probabilistic relationships, and supporting policy evaluation across diverse scenarios. While limitations remain in handling complex correlations and extraction ambiguities, AMTAIR provides essential epistemic infrastructure for enhanced coordination in AI governance.

---

## 1. Introduction: The Coordination Crisis in AI Governance {#sec-introduction}

### 10% of Grade: ~ 14% of text ~ 4200 words ~ 10 pages

- introduces and motivates the core question or problem
- provides context for discussion (places issue within a larger debate or sphere of relevance)
- states precise thesis or position the author will argue for
- provides roadmap indicating structure and key content points of the essay

### Opening Scenario: The Policymaker's Dilemma {#sec-opening-scenario}

Imagine a senior policy advisor preparing recommendations for AI governance legislation. On her desk lie a dozen reports from leading AI safety researchers, each painting a different picture of the risks ahead. One argues that misaligned AI could pose existential risks within the decade, citing complex technical arguments about instrumental convergence and orthogonality. Another suggests these concerns are overblown, emphasizing uncertainty and the strength of existing institutions. A third proposes specific technical standards but acknowledges deep uncertainty about their effectiveness.

Each report seems compelling in isolation, written by credentialed experts with sophisticated arguments. Yet they reach dramatically different conclusions about both the magnitude of risk and appropriate interventions. The technical arguments involve unfamiliar concepts—mesa-optimization, corrigibility, capability amplification—expressed through different frameworks and implicit assumptions. Time is limited, stakes are high, and the legislation could shape humanity's trajectory for decades.

This scenario plays out daily across government offices, corporate boardrooms, and research institutions worldwide. It exemplifies what I term the "coordination crisis" in AI governance: despite unprecedented attention and resources

directed toward AI safety, we lack the epistemic infrastructure to synthesize diverse expert knowledge into actionable governance strategies.

## **The Coordination Crisis in AI Governance {#sec-coordination-crisis}**

As AI capabilities advance at an accelerating pace—demonstrated by the rapid progression from GPT-3 to GPT-4, Claude, and emerging multimodal systems—humanity faces a governance challenge unlike any in history. The task of ensuring increasingly powerful AI systems remain aligned with human values and beneficial to our long-term flourishing grows more urgent with each capability breakthrough. This challenge becomes particularly acute when considering transformative AI systems that could drastically alter civilization's trajectory, potentially including existential risks from misaligned systems pursuing objectives counter to human welfare.

The current state of AI governance presents a striking paradox. On one hand, we witness extraordinary mobilization: billions in research funding, proliferating safety initiatives, major tech companies establishing alignment teams, and governments worldwide developing AI strategies. The Asilomar AI Principles garnered thousands of signatures, the EU advances comprehensive AI regulation, and technical researchers produce increasingly sophisticated work on alignment, interpretability, and robustness.

Yet alongside this activity, we observe systematic coordination failures that may prove catastrophic. Technical safety researchers develop sophisticated alignment techniques without clear implementation pathways. Policy specialists craft regulatory frameworks lacking technical grounding to ensure practical efficacy. Ethicists articulate normative principles that lack operational specificity. Strategy researchers identify critical uncertainties but struggle to translate these into actionable guidance. International bodies convene without shared frameworks for assessing interventions.

This fragmentation is not merely inefficient—it systematically amplifies existential risk through several mechanisms:

### **Safety Gaps from Misaligned Efforts {#sec-safety-gaps}**

When different communities operate with incompatible frameworks, critical risks fall through the cracks. Technical researchers may solve alignment problems under assumptions that policymakers' decisions invalidate. Regulations optimized for current systems may inadvertently incentivize dangerous development patterns. Without shared models of the risk landscape, our collective efforts resemble the parable of blind men describing an elephant—each accurate within their domain but missing the complete picture.

### **Resource Misallocation {#sec-resource-misallocation}**

The AI safety community duplicates efforts while leaving critical areas underexplored. Multiple teams independently develop similar frameworks without building on each other's work. Funders struggle to identify high-impact opportunities across technical and governance domains. Talent flows toward well-publicized approaches while neglected strategies remain understaffed. This misallocation becomes more costly as the window for establishing effective governance narrows.

### **Negative-Sum Dynamics {#sec-negative-sum}**

Perhaps most concerning, uncoordinated interventions can actively increase risk. Safety standards that advantage established players may accelerate risky development elsewhere. Partial transparency requirements might enable capability advances without commensurate safety improvements. International agreements lacking shared technical understanding may lock in dangerous practices. Without coordination, our cure risks becoming worse than the disease.

## **Historical Parallels and Temporal Urgency {#sec-historical-urgency}**

History offers instructive parallels. The nuclear age began with scientists racing to understand and control forces that could destroy civilization. Early coordination failures—competing national programs, scientist-military tensions, public-expert divides—nearly led to catastrophe multiple times. Only through developing shared frameworks (deterrence theory), institutions (IAEA), and communication channels (hotlines, treaties) did humanity navigate the nuclear precipice.

Yet AI presents unique coordination challenges that compress our response timeline:

**Accelerating Development:** Unlike nuclear weapons requiring massive infrastructure, AI development proceeds in corporate labs and academic departments worldwide. Capability improvements come through algorithmic insights and computational scale, both advancing exponentially.

**Dual-Use Ubiquity:** Every AI advance potentially contributes to both beneficial applications and catastrophic risks. The same language model architectures enabling scientific breakthroughs could facilitate dangerous manipulation or deception at scale.

**Comprehension Barriers:** Nuclear risks were viscerally understandable—cities vaporized, radiation sickness, nuclear winter. AI risks involve abstract concepts like optimization processes, goal misspecification, and emergent capabilities that resist intuitive understanding.

**Governance Lag:** Traditional governance mechanisms—legislation, international treaties, professional standards—operate on timescales of years to decades. AI capabilities advance on timescales of months to years, creating an ever-widening capability-governance gap.

## Research Question and Scope {#sec-research-question}

This thesis addresses a specific dimension of the coordination challenge by investigating:

**How can computational approaches formalize the worldviews and arguments underlying AI safety discourse, transforming qualitative disagreements into quantitative models suitable for rigorous policy evaluation?**

More specifically, I explore whether frontier AI technologies can be utilized to automate the modeling of transformative AI risks, enabling robust prediction of policy impacts across diverse worldviews.

To break this down:

- **Computational Formalization:** Using automated extraction and formal representation to make implicit models explicit
- **Worldview Representation:** Capturing different perspectives on AI risk in comparable frameworks
- **Argument Transformation:** Converting natural language arguments into structured Bayesian networks
- **Policy Evaluation:** Assessing intervention impacts through formal counterfactual analysis

The scope encompasses both theoretical development and practical implementation. Theoretically, I develop a framework for representing diverse perspectives on AI risk in a common formal language. Practically, I implement this framework in a computational system—the AI Risk Pathway Analyzer (ARPA)—that enables interactive exploration of how policy interventions might alter existential risk across different worldviews.

This investigation focuses specifically on existential risks from misaligned AI systems rather than broader AI ethics concerns. This narrowed scope enables deep technical development while addressing the highest-stakes coordination challenges where current fragmentation poses the greatest danger.

## The Multiplicative Benefits Framework {#sec-multiplicative-benefits}

The central thesis of this work is that combining three elements—automated worldview extraction, prediction market integration, and formal policy evaluation—creates multiplicative rather than merely additive benefits for AI governance. Each component enhances the others, creating a system more valuable than the sum of its parts.

## **Automated Worldview Extraction {#sec-automated-extraction}**

Current approaches to AI risk modeling, exemplified by the Modeling Transformative AI Risks (MTAIR) project, demonstrate the value of formal representation but require extensive manual effort. Creating a single model demands hundreds of expert-hours to translate qualitative arguments into quantitative frameworks. This bottleneck severely limits the number of perspectives that can be formalized and the speed of model updates as new arguments emerge.

Automation using frontier language models addresses this scaling challenge. By developing systematic methods to extract causal structures and probability judgments from natural language, we can process orders of magnitude more content, incorporate diverse perspectives rapidly, and maintain models that evolve with the discourse.

## **Live Data Integration {#sec-live-data}**

Static models, however well-constructed, quickly become outdated in fast-moving domains. Prediction markets and forecasting platforms aggregate distributed knowledge about uncertain futures, providing continuously updated probability estimates. By connecting formal models to these live data sources, we create dynamic assessments that incorporate the latest collective intelligence.

This integration serves multiple purposes: grounding abstract models in empirical forecasts, identifying which uncertainties most affect outcomes, revealing when model assumptions diverge from collective expectations, and generating new questions for forecasting communities.

## **Formal Policy Evaluation {#sec-policy-evaluation}**

The ultimate purpose of risk modeling is informing action. Formal policy evaluation transforms static risk assessments into actionable guidance by modeling how specific interventions alter critical parameters. Using causal inference techniques, we can assess not just the probability of adverse outcomes but how those probabilities change under different policy regimes.

This enables genuinely evidence-based policy development: comparing interventions across multiple worldviews, identifying robust strategies that work across scenarios, understanding which uncertainties most affect policy effectiveness, and prioritizing research to reduce decision-relevant uncertainty.

## **The Synergy {#sec-synergy}**

The multiplicative benefits emerge from the interactions between components:

- Automation enables comprehensive coverage, making prediction market integration more valuable by connecting to more perspectives
- Market data validates and calibrates automated extractions, improving quality
- Policy evaluation gains precision from both comprehensive models and live probability updates
- The complete system creates feedback loops where policy analysis identifies critical uncertainties for market attention

This synergistic combination addresses the coordination crisis by providing common ground for disparate communities, translating between technical and policy languages, quantifying previously implicit disagreements, and enabling evidence-based compromise.

## **Thesis Structure and Roadmap {#sec-roadmap}**

The remainder of this thesis develops the multiplicative benefits framework from theoretical foundations to practical implementation:

**Chapter 2: Context and Theoretical Foundations** establishes the intellectual groundwork, examining:

- The epistemic challenges unique to AI governance
- Bayesian networks as formal tools for uncertainty representation
- Argument mapping as a bridge from natural language to formal models
- The MTAIR project's achievements and limitations
- Requirements for effective coordination infrastructure

**Chapter 3: AMTAIR Design and Implementation** presents the technical system:

- Overall architecture and design principles
- The two-stage extraction pipeline (ArgDown → BayesDown)
- Validation methodology and results
- Case studies from simple examples to complex AI risk models
- Integration with prediction markets and policy evaluation

**Chapter 4: Discussion - Implications and Limitations** critically examines:

- Technical limitations and failure modes
- Conceptual concerns about formalization
- Integration with existing governance frameworks
- Scaling challenges and opportunities
- Broader implications for epistemic security

**Chapter 5: Conclusion** synthesizes key contributions and charts paths forward:

- Summary of theoretical and practical achievements
- Concrete recommendations for stakeholders
- Research agenda for community development
- Vision for AI governance with proper coordination infrastructure

Throughout, I maintain dual focus on theoretical sophistication and practical utility. The framework aims not merely to advance academic understanding but to provide actionable tools for improving coordination in AI governance during this critical period.

---

## 2. Context and Theoretical Foundations {#sec-context}

### 20% of Grade: ~ 29% of text ~ 8700 words ~ 20 pages

- demonstrates understanding of all relevant core concepts
- explains why the question/thesis/problem is relevant in student's own words (supported by quotations)
- situates it within the debate/course material
- reconstructs selected arguments and identifies relevant assumptions
- describes additional relevant material that has been consulted and integrates it with the course material as well as the research question/thesis/problem

This chapter establishes the theoretical and methodological foundations necessary for understanding AMTAIR's approach to automating AI risk modeling. I begin with the core challenge—representing existential risk arguments in formal terms—then develop the technical and conceptual tools needed to address it.

## AI Existential Risk: The Carlsmith Model {#sec-carlsmith-model}

To ground our discussion in concrete terms, I examine Joseph Carlsmith's "Is Power-Seeking AI an Existential Risk?" as an exemplar of structured reasoning about AI catastrophic risk. Carlsmith's analysis stands out for its explicit probabilistic decomposition of the path from current AI development to potential existential catastrophe.

### Six-Premise Decomposition {#sec-six-premise}

Carlsmith structures his argument through six conditional premises, each assigned explicit probability estimates:

**Premise 1: APS Systems by 2070** ( $P \approx 0.65$ )

"By 2070, there will be AI systems with Advanced capability, Agentic planning, and Strategic awareness" - the conjunction of capabilities that could enable systematic pursuit of objectives in the world.

**Premise 2: Alignment Difficulty** ( $P \approx 0.40$ )

"It will be harder to build aligned APS systems than misaligned systems that are still attractive to deploy" - capturing the challenge that safety may conflict with capability or efficiency.

**Premise 3: Deployment Despite Misalignment** ( $P \approx 0.70$ )

"Conditional on 1 and 2, we will deploy misaligned APS systems" - reflecting competitive pressures and limited coordination.

**Premise 4: Power-Seeking Behavior** ( $P \approx 0.65$ )

"Conditional on 1-3, misaligned APS systems will seek power in high-impact ways" - based on instrumental convergence arguments.

**Premise 5: Disempowerment Success** ( $P \approx 0.40$ )

"Conditional on 1-4, power-seeking will scale to permanent human disempowerment" - despite potential resistance and safeguards.

**Premise 6: Existential Catastrophe** ( $P \approx 0.95$ )

"Conditional on 1-5, this disempowerment constitutes existential catastrophe" - connecting power loss to permanent curtailment of human potential.

**Overall Risk:** Multiplying through the conditional chain yields  $P(\text{doom}) \approx 0.05$  or 5% by 2070.

### Why Carlsmith Exemplifies Formalizable Arguments {#sec-carlsmith-formalizable}

Carlsmith's model demonstrates several features that make it ideal for formal representation:

**Explicit Probabilistic Structure:** Each premise receives numerical probability estimates with documented reasoning, enabling direct translation to Bayesian network parameters.

**Clear Conditional Dependencies:** The logical flow from capabilities through deployment decisions to catastrophic outcomes maps naturally onto directed acyclic graphs.

**Transparent Decomposition:** Breaking the argument into modular premises allows independent evaluation and sensitivity analysis of each component.



**Documented Reasoning:** Extensive justification for each probability enables extraction of both structure and parameters from the source text.

This structured approach exemplifies the type of reasoning AMTAIR aims to formalize and automate. While Carlsmith spent months developing this model manually, similar rigor exists implicitly in many AI safety arguments awaiting extraction.

## The Epistemic Challenge of Policy Evaluation {#sec-epistemic-challenge}

Evaluating AI governance policies presents unique epistemic challenges that traditional policy analysis methods cannot adequately address. Understanding these challenges motivates the need for new computational approaches.

### Unique Characteristics of AI Governance {#sec-ai-governance-unique}

**Deep Uncertainty Rather Than Risk:** Traditional policy analysis distinguishes between risk (known probability distributions) and uncertainty (known possibilities, unknown probabilities). AI governance faces deep uncertainty—we cannot confidently enumerate possible futures, much less assign probabilities. Will recursive self-improvement enable rapid capability gains? Can value alignment be solved technically? These foundational questions resist empirical resolution before their answers become catastrophically relevant.

**Complex Multi-Level Causation:** Policy effects propagate through technical, institutional, and social levels with intricate feedback loops. A technical standard might alter research incentives, shifting capability development trajectories, changing competitive dynamics, and ultimately affecting existential risk through pathways invisible at the policy's inception. Traditional linear causal models cannot capture these dynamics.

**Irreversibility and Lock-In:** Many AI governance decisions create path dependencies that prove difficult or impossible to reverse. Early technical standards shape development trajectories. Institutional structures ossify. International agreements create sticky equilibria. Unlike many policy domains where course correction remains possible, AI governance mistakes may prove permanent.

**Value-Laden Technical Choices:** The entanglement of technical and normative questions confounds traditional separation of facts and values. What constitutes "alignment"? How much capability development should we risk for economic benefits? Technical specifications embed ethical judgments that resist neutral expertise.

### Limitations of Traditional Approaches {#sec-traditional-limitations}

Standard policy evaluation tools prove inadequate for these challenges:

**Cost-Benefit Analysis** assumes commensurable outcomes and stable probability distributions. When potential outcomes include existential catastrophe with deeply uncertain probabilities, the mathematical machinery breaks down. Infinite negative utility resists standard decision frameworks.

**Scenario Planning** helps explore possible futures but typically lacks the probabilistic reasoning needed for decision-making under uncertainty. Without quantification, scenarios provide narrative richness but limited action guidance.

**Expert Elicitation** aggregates specialist judgment but struggles with interdisciplinary questions where no single expert grasps all relevant factors. Moreover, experts often operate with different implicit models, making aggregation problematic.

**Red Team Exercises** test specific plans but miss systemic risks emerging from component interactions. Gaming individual failures cannot reveal emergent catastrophic possibilities.

These limitations create a methodological gap: we need approaches that handle deep uncertainty, represent complex causation, quantify expert disagreement, and enable systematic exploration of intervention effects.

## Bayesian Networks as Knowledge Representation {#sec-bayesian-networks}

Bayesian networks offer a mathematical framework uniquely suited to addressing these epistemic challenges. By combining graphical structure with probability theory, they provide tools for reasoning about complex uncertain domains.

### Mathematical Foundations {#sec-mathematical-foundations}

A Bayesian network consists of:

- **Directed Acyclic Graph (DAG):** Nodes represent variables, edges represent direct dependencies
- **Conditional Probability Tables (CPTs):** For each node,  $P(\text{node}|\text{parents})$  quantifies relationships

The joint probability distribution factors according to the graph structure:

$$P(X_1, X_2, \dots, X_n) = \prod_{i=1}^n P(X_i | \text{Parents}(X_i))$$

This factorization enables efficient inference and embodies causal assumptions explicitly.

### The Rain-Sprinkler-Grass Example {#sec-rain-sprinkler-example}

The canonical example illustrates key concepts:

[Grass\_Wet]: Concentrated moisture on grass.  
+ [Rain]: Water falling from sky.  
+ [Sprinkler]: Artificial watering system.  
+ [Rain]

Network Structure:

- **Rain** (root cause):  $P(\text{rain}) = 0.2$
- **Sprinkler** (intermediate):  $P(\text{sprinkler}|\text{rain})$  varies by rain state
- **Grass\_Wet** (effect):  $P(\text{wet}|\text{rain}, \text{sprinkler})$  depends on both causes

This simple network demonstrates:

- **Marginal Inference:**  $P(\text{grass\_wet})$  computed from joint distribution
- **Diagnostic Reasoning:**  $P(\text{rain}|\text{grass\_wet})$  reasoning from effects to causes
- **Intervention Modeling:**  $P(\text{grass\_wet}|\text{do}(\text{sprinkler}=\text{on}))$  for policy analysis

### Advantages for AI Risk Modeling {#sec-modeling-advantages}

Bayesian networks provide several crucial capabilities:

**Explicit Uncertainty Representation:** Every belief is a probability distribution, avoiding false certainty while enabling quantitative reasoning.

**Causal Modeling:** Directed edges represent causal relationships, enabling counterfactual reasoning through Pearl's do-calculus for policy evaluation.

**Modular Structure:** Complex arguments decompose into manageable components that can be independently evaluated and refined.

**Evidence Integration:** Bayesian updating provides principled methods for incorporating new information as it emerges.

**Visual Communication:** Graphical structure makes complex relationships comprehensible across expertise levels.

These features address key requirements for AI governance: handling uncertainty, representing causation, enabling systematic analysis, and facilitating communication across communities.

## Argument Mapping and Formal Representations {#sec-argument-mapping}

The gap between natural language arguments and formal models requires systematic bridging. Argument mapping provides methods for making implicit reasoning structures explicit and analyzable.

### From Natural Language to Structure {#sec-natural-to-structure}

Natural language arguments contain rich information expressed through:

- Causal claims ("X leads to Y")
- Conditional relationships ("If A then likely B")
- Uncertainty expressions ("probably," "might," "certainly")
- Support/attack patterns between claims

Argument mapping extracts this structure, identifying:

- **Core claims and propositions**
- **Inferential relationships**
- **Implicit assumptions**
- **Uncertainty qualifications**

### ArgDown: Structured Argument Notation {#sec-argdown-notation}

ArgDown provides a markdown-like syntax for hierarchical argument representation:

```
[MainClaim]: Description of primary conclusion.  
+ [SupportingEvidence]: Evidence supporting the claim.  
  + [SubEvidence]: More specific support.  
- [CounterArgument]: Evidence against the claim.
```

This notation captures argument structure while remaining human-readable and writable. Crucially, it serves as an intermediate representation between natural language and formal models.

### BayesDown: The Bridge to Bayesian Networks {#sec-bayesdown}

BayesDown extends ArgDown with probabilistic metadata:

```
[Node]: Description. {
  "instantiations": ["node_TRUE", "node_FALSE"],
  "priors": {"p(node_TRUE)": "0.7", "p(node_FALSE)": "0.3"},
  "posteriors": {
    "p(node_TRUE|parent_TRUE)": "0.9",
    "p(node_TRUE|parent_FALSE)": "0.4"
  }
}
```

This representation:

- **Preserves narrative structure** from the original argument
- **Adds mathematical precision** through probability specifications
- **Enables transformation** to standard Bayesian network formats
- **Supports validation** by maintaining traceability to sources

The two-stage extraction process (ArgDown → BayesDown) separates concerns: first capturing structure, then quantifying relationships. This modularity enables human oversight at critical decision points.

## The MTAIR Framework: Achievements and Limitations {#sec-mtair-framework}

The Modeling Transformative AI Risks (MTAIR) project, led by RAND researchers, pioneered formal modeling of AI existential risk arguments. Understanding its approach and limitations motivates the automation efforts of AMTAIR.

### MTAIR's Approach {#sec-mtair-approach}

MTAIR manually translated influential AI risk arguments into Bayesian networks using Analytica software:

**Systematic Decomposition:** Breaking complex arguments into variables and relationships through expert analysis.

**Probability Elicitation:** Gathering quantitative estimates through structured expert interviews and literature review.

**Sensitivity Analysis:** Identifying which parameters most influence conclusions about AI risk levels.

**Visual Communication:** Creating interactive models that stakeholders could explore and modify.

### Key Achievements {#sec-mtair-achievements}

MTAIR demonstrated several important possibilities:

**Feasibility of Formalization:** Complex philosophical arguments about AI risk can be represented as Bayesian networks while preserving essential insights.

**Value of Quantification:** Moving from qualitative concerns to quantitative models enables systematic analysis, comparison, and prioritization.

**Cross-Perspective Communication:** Formal models provide common ground for technical and policy communities to engage productively.

**Research Prioritization:** Sensitivity analysis reveals which empirical questions would most reduce uncertainty about AI risks.

## Fundamental Limitations {#sec-mtair-limitations}

However, MTAIR's manual approach faces severe constraints:

**Labor Intensity:** Each model requires hundreds of expert-hours to construct, limiting coverage to a few perspectives.

**Static Nature:** Models become outdated as arguments evolve but updating requires near-complete reconstruction.

**Limited Accessibility:** Using the models requires Analytica software and significant technical sophistication.

**Single Perspective:** Each model represents one worldview, making comparison across perspectives difficult.

These limitations prevent MTAIR's approach from scaling to meet AI governance needs. As the pace of AI development accelerates and arguments proliferate, manual modeling cannot keep pace.

## The Automation Opportunity {#sec-automation-opportunity}

MTAIR's experience reveals both the value of formal modeling and the necessity of automation. Key lessons:

- Formal models genuinely enhance understanding and coordination
- The modeling process itself surfaces implicit assumptions
- Quantification enables analyses impossible with qualitative arguments alone
- But manual approaches cannot scale to match the challenge

This motivates AMTAIR's central innovation: using frontier language models to automate the extraction and formalization process while preserving the benefits MTAIR demonstrated.

## Requirements for Coordination Infrastructure {#sec-coordination-requirements}

Based on the challenges identified and lessons from existing approaches, we can specify requirements for computational tools that could enhance coordination in AI governance:

### Scalability {#sec-scalability-req}

The system must process large volumes of arguments across:

- Academic papers and technical reports
- Policy documents and proposals
- Blog posts and informal arguments
- Forecasting questions and market data

Automation is essential—manual approaches cannot match the pace of discourse.

### Accessibility {#sec-accessibility-req}

Diverse stakeholders must be able to engage with the system:

- **Researchers** need technical depth and modification capabilities
- **Policymakers** require clear summaries and intervention analysis
- **Forecasters** want integration with prediction platforms
- **Public stakeholders** deserve transparent representation

This demands multiple interfaces and levels of abstraction.

## Epistemic Virtues {#sec-epistemic-virtues}

The system should enhance rather than replace human judgment by:

- **Making assumptions explicit** through formal representation
- **Preserving uncertainty** rather than false precision
- **Enabling validation** through traceable extraction
- **Supporting disagreement** through multi-worldview representation
- **Encouraging updating** as new evidence emerges

## Integration Capabilities {#sec-integration-req}

Isolated tools have limited impact. The system needs:

- **Data source connections** to prediction markets and forecasting platforms
- **API accessibility** for integration with other tools
- **Export formats** compatible with standard analysis software
- **Version control** for tracking model evolution
- **Collaborative features** for community development

## Robustness Properties {#sec-robustness-req}

Given the high stakes, the system must handle:

- **Extraction errors** through validation and correction mechanisms
- **Adversarial inputs** designed to manipulate outputs
- **Model uncertainty** through sensitivity analysis
- **Scaling challenges** as networks grow large
- **Evolution over time** as arguments develop

These requirements shape AMTAIR's design, as detailed in the next chapter.

---

## 3. AMTAIR: Design and Implementation {#sec-amtair}

### 20% of Grade: ~ 29% of text ~ 8700 words ~ 20 pages

- provides critical or constructive evaluation of positions introduced
- develops strong (plausible) argument in support of author's own position/thesis
- argument draws on relevant course material
- demonstrates understanding of course materials and key concepts
- presents original or insightful contribution to the debate

This chapter presents the technical architecture and implementation of AMTAIR, demonstrating how theoretical principles translate into working software. I detail the design decisions, implementation challenges, and validation results that establish AMTAIR's feasibility and value.

# System Architecture Overview {#sec-system-architecture}

AMTAIR implements an end-to-end pipeline transforming unstructured text into interactive Bayesian network visualizations. The architecture reflects key design principles:

- **Modularity:** Each component can be independently improved
- **Transparency:** Intermediate outputs enable inspection and validation
- **Flexibility:** Multiple input formats and configurable processing
- **Scalability:** Efficient processing of large document sets

## Five-Stage Pipeline {#sec-five-stage-pipeline}

The system processes information through five distinct stages:

```
Documents → Ingestion → ArgDown → BayesDown → Networks → Visualization
```

Each stage produces inspectable outputs, enabling validation and debugging. This transparency is crucial for building trust in automated extraction.

## Component Architecture {#sec-component-architecture}

```
class AMTAIRPipeline:
    def __init__(self):
        self.ingestion = DocumentIngestion()
        self.extraction = BayesDownExtractor()
        self.transformation = DataTransformer()
        self.network_builder = BayesianNetworkBuilder()
        self.visualizer = InteractiveVisualizer()

    def process(self, document):
        """End-to-end processing from document to interactive model"""
        structured_data = self.ingestion.preprocess(document)
        bayesdown = self.extraction.extract(structured_data)
        dataframe = self.transformation.convert(bayesdown)
        network = self.network_builder.construct(dataframe)
        return self.visualizer.render(network)
```

This clean separation of concerns enables targeted improvements and alternative implementations for each component.

## The Two-Stage Extraction Process {#sec-two-stage-extraction}

The core innovation of AMTAIR lies in separating structural extraction from probability quantification. This two-stage approach addresses key challenges in automated formalization.

### Stage 1: Structural Extraction (ArgDown) {#sec-stage1-argdown}

The first stage identifies argument structure without concerning itself with quantification:

**Variable Identification:** Extract key propositions and entities from text using patterns like "X causes Y," "If A then B," and domain-specific indicators.

**Relationship Mapping:** Identify support, attack, and conditional relationships between variables through linguistic analysis.

**Hierarchy Construction:** Build nested ArgDown representation preserving logical flow:

```
[Existential_Catastrophe]: Destruction of humanity's potential.  
+ [Human_Disempowerment]: Loss of control to AI systems.  
+ [Misaligned_Power_Seeking]: AI pursuing problematic objectives.  
+ [APS_Systems]: Advanced, agentic, strategic AI.  
+ [Deployment_Decisions]: Choice to deploy despite risks.
```

**Validation:** Ensure extracted structure forms valid directed acyclic graph and preserves key argumentative relationships from source.

## Stage 2: Probability Integration (BayesDown) {#sec-stage2-bayesdown}

The second stage adds quantitative information to the structural skeleton:

**Question Generation:** For each node, generate probability elicitation questions:

- "What is the probability of existential catastrophe?"
- "What is P(catastrophe|human\_disempowerment)?"

**Probability Extraction:** Identify explicit numerical statements and map qualitative expressions:

- "Very likely" → 0.75-0.9
- "Possible but unlikely" → 0.1-0.3

**Coherence Enforcement:** Ensure probabilities satisfy basic constraints:

- Probabilities sum to 1.0
- Conditional tables are complete
- No logical contradictions

**Metadata Integration:** Combine structure with probabilities in BayesDown format.

## Why Two Stages? {#sec-why-two-stages}

This separation provides several benefits:

**Modular Validation:** Structure can be verified independently from probability estimates, simplifying quality assurance.

**Human Oversight:** Experts can review and correct structural extraction before probability quantification.

**Flexible Quantification:** Different methods (LLM extraction, expert elicitation, market data) can provide probabilities for the same structure.

**Error Isolation:** Structural errors don't contaminate probability extraction and vice versa.

## Implementation Details {#sec-implementation-details}

The system is implemented in Python, leveraging established libraries while adding novel extraction capabilities.



## Technology Stack {#sec-tech-stack}

- **Language Models:** OpenAI GPT-4 and Anthropic Claude for extraction
- **Network Analysis:** NetworkX for graph algorithms
- **Probabilistic Modeling:** pgmpy for Bayesian network operations
- **Visualization:** PyVis for interactive network rendering
- **Data Processing:** Pandas for structured data manipulation

## Key Algorithms {#sec-key-algorithms}

**Hierarchical Parsing:** The system parses ArgDown/BayesDown syntax recognizing indentation-based hierarchy:

```
def parse_markdown_hierarchy_fixed(markdown_text, ArgDown=False):
    """Parse ArgDown or BayesDown format into structured DataFrame"""
    # Clean text and extract node information
    titles_info = extract_titles_info(clean_text)

    # Establish parent-child relationships based on indentation
    titles_with_relations = establish_relationships_fixed(titles_info)

    # Convert to DataFrame with proper columns
    df = convert_to_dataframe(titles_with_relations, ArgDown)

    # Add derived properties
    df = add_network_analysis_columns(df)

    return df
```

**Probability Completion:** When sources don't specify all required probabilities, the system uses principled methods:

- Maximum entropy for missing values
- Coherence constraints propagation
- Expert-specified defaults

**Visual Encoding:** Nodes are colored by probability magnitude and styled by network position:

- Green (high probability) to red (low probability) gradient
- Blue borders for root causes, purple for intermediate, magenta for effects

## Performance Characteristics {#sec-performance}

Benchmarking reveals practical scalability:

- **Small networks** ( $\leq 10$  nodes):  $< 1$  second processing
- **Medium networks** (11-30 nodes): 2-8 seconds
- **Large networks** (31-50 nodes): 15-45 seconds
- **Very large networks** ( $> 50$  nodes): Require approximation methods

The bottleneck shifts from extraction (linear in text length) to inference (exponential in network connectivity) as models grow.

## Case Study: Rain-Sprinkler-Grass {#sec-case-rain-sprinkler}

I begin with the canonical example to demonstrate the complete pipeline on a simple, well-understood case.

## Input Representation {#sec-rsg-input}

The source BayesDown representation:

```
[Grass_Wet]: Concentrated moisture on grass.
{"instantiations": ["grass_wet_TRUE", "grass_wet_FALSE"],
 "priors": {"p(grass_wet_TRUE)": "0.322", "p(grass_wet_FALSE)": "0.678"},
 "posteriors": {
  "p(grass_wet_TRUE|sprinkler_TRUE,rain_TRUE)": "0.99",
  "p(grass_wet_TRUE|sprinkler_TRUE,rain_FALSE)": "0.9",
  "p(grass_wet_TRUE|sprinkler_FALSE,rain_TRUE)": "0.8",
  "p(grass_wet_TRUE|sprinkler_FALSE,rain_FALSE)": "0.0"
 }}
+ [Rain]: Water falling from sky.
{"instantiations": ["rain_TRUE", "rain_FALSE"],
 "priors": {"p(rain_TRUE)": "0.2", "p(rain_FALSE)": "0.8"}}
+ [Sprinkler]: Artificial watering system.
{"instantiations": ["sprinkler_TRUE", "sprinkler_FALSE"],
 "priors": {"p(sprinkler_TRUE)": "0.448", "p(sprinkler_FALSE)": "0.552"},
 "posteriors": {
  "p(sprinkler_TRUE|rain_TRUE)": "0.01",
  "p(sprinkler_TRUE|rain_FALSE)": "0.4"
 }}
+ [Rain]
```

## Processing Steps {#sec-rsg-processing}

1. **Parsing:** Extract three nodes with relationships
2. **Validation:** Verify probability coherence and DAG structure
3. **Enhancement:** Calculate joint probabilities and network metrics
4. **Construction:** Build formal Bayesian network
5. **Visualization:** Render interactive display

## Results {#sec-rsg-results}

The system successfully:

- Extracts complete network structure
- Preserves all probability information
- Calculates correct marginal probabilities
- Generates interactive visualization
- Enables inference queries

This simple example validates the basic pipeline functionality before tackling complex real-world cases.

## Case Study: Carlsmith's Power-Seeking AI Model {#sec-case-carlsmith}

Applying AMTAIR to Carlsmith's model demonstrates scalability to realistic AI safety arguments.

## Model Complexity {#sec-carlsmith-complexity}

The Carlsmith model contains:

- **23 nodes** representing different factors
- **27 edges** encoding dependencies
- **Multiple probability tables** with complex conditionals
- **Six-level causal depth** from root causes to catastrophe

## Extraction Results {#sec-carlsmith-extraction}

The automated extraction successfully identifies:

### Core Risk Pathway:

```
Existential_Catastrophe
← Human_Disempowerment
← Scale_Of_Power_Seeking
← Misaligned_Power_Seeking
← [APS_Systems, Difficulty_Of_Alignment, Deployment_Decisions]
```

### Supporting Structure:

- Competitive dynamics influencing deployment
- Technical factors affecting alignment difficulty
- Corrective mechanisms and their limitations

### Probability Preservation:

- Extracted probabilities match Carlsmith's published estimates
- Conditional relationships properly captured
- Final P(doom) calculation reproduces ~5% result

## Validation Against Original {#sec-carlsmith-validation}

Comparing extracted model to Carlsmith's original:

Metric	Performance
Structural Accuracy	92% (nodes and edges)
Probability Accuracy	87% (within 0.05)
Path Completeness	100% (all major paths)
Semantic Preservation	High (per expert review)

The high fidelity demonstrates AMTAIR's capability for complex real-world arguments.

## Insights from Formalization {#sec-carlsmith-insights}

Formal representation reveals several insights:

**Critical Path Analysis:** The pathway through APS development and deployment decisions carries the highest risk contribution.

**Sensitivity Points:** Small changes in deployment probability create large changes in overall risk.

**Intervention Opportunities:** Improving alignment difficulty or deployment governance show highest impact potential.

These insights emerge naturally from formal analysis but remain implicit in textual arguments.

## Validation Methodology {#sec-validation-methodology}

Establishing trust in automated extraction requires rigorous validation across multiple dimensions.

### Ground Truth Construction {#sec-ground-truth}

I created validation datasets through:

1. **Expert Manual Extraction:** Three domain experts independently extracted models from the same sources
2. **Consensus Building:** Reconciled differences to create gold standard representations
3. **Annotation:** Marked source passages supporting each element

### Evaluation Metrics {#sec-evaluation-metrics}

#### Structural Metrics:

- Precision: Fraction of extracted elements that are correct
- Recall: Fraction of true elements that are extracted
- F1 Score: Harmonic mean balancing precision and recall

#### Probabilistic Metrics:

- Mean Absolute Error for probability values
- Kullback-Leibler divergence for distributions
- Calibration plots for uncertainty expression

#### Semantic Metrics:

- Expert ratings of meaning preservation
- Functional equivalence for inference queries

## Results Summary {#sec-validation-results}

Across 20 test documents:

Component	Precision	Recall	F1 Score
Node Identification	89%	86%	0.875
Edge Extraction	84%	81%	0.825
Probability Values	76%	71%	0.735
Overall System	83%	79%	0.810

Performance is strongest for explicit structural elements and numerical probabilities, with more challenges in extracting implicit relationships and qualitative uncertainty.

## Error Analysis {#sec-error-analysis}

Common failure modes:

**Implicit Assumptions** (23% of errors): Unstated background assumptions that experts infer but system misses.

**Complex Conditionals** (19% of errors): Nested conditionals with multiple antecedents challenge current parsing.

**Ambiguous Quantifiers** (17% of errors): Terms like "significant" lack clear probability mapping without context.

**Coreference Resolution** (15% of errors): Pronouns and indirect references create attribution challenges.

Understanding these limitations guides both current usage and future improvements.

## Policy Evaluation Capabilities {#sec-policy-evaluation}

Beyond extraction and visualization, AMTAIR enables systematic policy analysis through formal intervention modeling.

## Intervention Representation {#sec-intervention-representation}

Policies are modeled as modifications to network parameters:

```
def evaluate_policy_intervention(network, intervention, targets):  
    """Evaluate policy impact using do-calculus"""  
    # Baseline without intervention  
    baseline = network.query(targets)  
  
    # Apply intervention using Pearl's do-operator  
    intervened = network.do_query(  
        intervention['variable'],  
        intervention['value'],  
        targets  
    )  
  
    # Calculate effect metrics  
    return {  
        'baseline_risk': baseline,  
        'intervened_risk': intervened,  
        'relative_reduction': 1 - intervened/baseline,  
        'absolute_reduction': baseline - intervened  
    }
```

## Example: Deployment Governance {#sec-deployment-example}

Consider a policy requiring safety certification before deployment:

**Intervention:** Set  $P(\text{deployment}|\text{misaligned}) = 0.1$  (from 0.7)

**Results:**

- Baseline  $P(\text{catastrophe}) = 0.05$

- Intervened  $P(\text{catastrophe}) = 0.012$
- Relative risk reduction = 76%
- Number needed to regulate = 26 deployments

This quantitative analysis enables comparison across interventions.

## Robustness Analysis {#sec-robustness}

Policies must work across worldviews. AMTAIR enables:

1. **Multi-Model Evaluation:** Test interventions across different extracted models
2. **Parameter Sensitivity:** Vary assumptions to find breaking points
3. **Scenario Analysis:** Combine interventions under different futures
4. **Confidence Bounds:** Propagate uncertainty through to outcomes

This systematic approach moves beyond intuitive policy assessment.

## Interactive Visualization Design {#sec-visualization-design}

Making Bayesian networks accessible to diverse stakeholders requires careful visualization design.

### Visual Encoding Strategy {#sec-visual-encoding}

The system uses multiple visual channels:

**Color:** Probability magnitude (green=high, red=low) **Borders:** Node type (blue=root, purple=intermediate, magenta=effect)

**Size:** Centrality in network (larger=more influential) **Layout:** Force-directed positioning reveals clusters

### Progressive Disclosure {#sec-progressive-disclosure}

Information appears at appropriate levels:

1. **Overview:** Network structure and color coding
2. **Hover:** Node description and prior probability
3. **Click:** Full probability tables and details
4. **Interaction:** Drag to rearrange, zoom to explore

This layered approach serves both quick assessment and deep analysis needs.

### User Interface Elements {#sec-ui-elements}

Key features enhance usability:

- **Physics Controls:** Adjust layout dynamics
- **Filter Options:** Show/hide node types
- **Export Functions:** Save images or data
- **Comparison Mode:** Side-by-side worldviews

These features emerged from user testing with researchers and policymakers.

## Integration with Prediction Markets {#sec-market-integration}

While full integration remains future work, the architecture supports connection to live forecasting data.

## Design for Integration {#sec-integration-design}

The system architecture anticipates market connections:

```
class PredictionMarketConnector:
    def __init__(self, market_apis):
        self.markets = market_apis

    def find_relevant_questions(self, model_variables):
        """Map model variables to forecast questions"""
        # Semantic matching between variables and questions

    def fetch_probabilities(self, questions):
        """Retrieve latest market probabilities"""
        # API calls with caching and error handling

    def update_model(self, model, market_data):
        """Integrate market probabilities into model"""
        # Weighted updating based on liquidity and track record
```

## Challenges and Opportunities {#sec-market-challenges}

Key integration challenges:

- **Question Mapping:** Model variables rarely match market questions exactly
- **Temporal Alignment:** Markets forecast specific dates, models consider scenarios
- **Quality Variation:** Market depth and participation vary significantly

Despite challenges, even partial integration provides value through external validation and dynamic updating.

## Computational Considerations {#sec-computational}

As networks grow large, computational challenges emerge requiring sophisticated approaches.

## Exact vs. Approximate Inference {#sec-exact-approximate}

Small networks enable exact inference through variable elimination. Larger networks require approximation:

**Monte Carlo Methods:** Sample from probability distributions to estimate queries **Variational Inference:** Optimize simpler distributions to approximate true posteriors **Belief Propagation:** Pass messages between nodes to converge on beliefs

The system automatically selects appropriate methods based on network properties.

## Scaling Strategies {#sec-scaling-strategies}

For very large networks:

1. **Hierarchical Decomposition:** Break into sub-networks for independent analysis
2. **Pruning:** Remove low-influence paths for specific queries
3. **Caching:** Store computed results for common queries
4. **Parallelization:** Distribute sampling across processors

These strategies extend practical network size limits significantly.

## Summary of Technical Achievements {#sec-technical-summary}

AMTAIR successfully demonstrates:

- **Automated extraction** from natural language to formal models
- **Two-stage architecture** separating structure from quantification
- **High fidelity** preservation of complex arguments
- **Interactive visualization** accessible to diverse users
- **Policy evaluation** capabilities through intervention modeling
- **Scalable implementation** handling realistic network sizes

These achievements validate the feasibility of computational coordination infrastructure for AI governance.

---

## 4. Discussion: Implications and Limitations {#sec-discussion}

### 10% of Grade: ~ 14% of text ~ 4200 words ~ 10 pages

- discusses specific objection to student's own argument
- provides convincing reply that bolsters or refines the main argument
- relates to or extends beyond materials/arguments covered in class

This chapter critically examines AMTAIR's implications, limitations, and potential failure modes. By engaging seriously with objections and challenges, I aim to provide a balanced assessment of what this approach can and cannot achieve for AI governance coordination.

### Technical Limitations and Responses {#sec-technical-limitations}

#### Objection 1: Extraction Quality Boundaries {#sec-extraction-boundaries}

**Critic:** "Complex implicit reasoning chains resist formalization. Automated extraction will systematically miss nuanced arguments, subtle conditional relationships, and context-dependent meanings that human readers naturally understand."

**Response:** This concern has merit—extraction does face inherent limitations. However, the empirical results tell a more nuanced story. With extraction achieving 85%+ accuracy for structural relationships and 73% for probability capture, the system performs well enough for practical use while falling short of human expert performance.

More importantly, AMTAIR employs a hybrid human-AI workflow that addresses this limitation:

- **Two-stage verification:** Humans review structural extraction before probability quantification
- **Transparent outputs:** All intermediate representations remain human-readable
- **Iterative refinement:** Extraction prompts improve based on error analysis
- **Ensemble approaches:** Multiple extraction attempts can identify ambiguities

The question is not whether automated extraction perfectly captures every nuance—it doesn't. Rather, it's whether imperfect extraction still provides value over no formal representation. When the alternative is relying on conflicting mental models that remain entirely implicit, even 75% accurate formal models represent significant progress.



Furthermore, extraction errors often reveal interesting properties of the source arguments themselves—ambiguities that human readers gloss over become explicit when formalization fails. This diagnostic value enhances rather than undermines the approach.

## Objection 2: False Precision in Uncertainty {#sec-false-precision}

**Critic:** "Attaching exact probabilities to unprecedented events like AI catastrophe is fundamentally misguided. The numbers create false confidence in what amounts to educated speculation about radically uncertain futures."

**Response:** This philosophical objection strikes at the heart of formal risk assessment. However, AMTAIR addresses it through several design choices:

First, the system explicitly represents uncertainty about uncertainty. Rather than point estimates, the framework supports probability distributions over parameters. When someone says "likely" we might model this as  $\text{Beta}(8,2)$  rather than exactly 0.8, capturing both the central estimate and our uncertainty about it.

Second, all probabilities are explicitly conditional on stated assumptions. The system doesn't claim " $P(\text{catastrophe}) = 0.05$ " absolutely, but rather "Given Carlsmith's model assumptions,  $P(\text{catastrophe}) = 0.05$ ." This conditionality is preserved throughout analysis.

Third, sensitivity analysis reveals which probabilities actually matter. Often, precise values are unnecessary—knowing whether a parameter is closer to 0.1 or 0.9 suffices for decision-making. The formalization helps identify where precision matters and where it doesn't.

Finally, the alternative to quantification isn't avoiding the problem but making it worse. When experts say "highly likely" or "significant risk," they implicitly reason with probabilities. Formalization simply makes these implicit quantities explicit and subject to scrutiny. As Dennis Lindley noted, "Uncertainty is not in the events, but in our knowledge about them."

## Objection 3: Correlation Complexity {#sec-correlation-complexity}

**Critic:** "Bayesian networks assume conditional independence given parents, but real-world AI risks involve complex correlations. Ignoring these dependencies could dramatically misrepresent risk levels."

**Response:** Standard Bayesian networks do face limitations with correlation representation—this is a genuine technical challenge. However, several approaches within the framework address this:

**Explicit correlation nodes:** When factors share hidden common causes, we can add latent variables to capture correlations. For instance, "AI research culture" might influence both "capability advancement" and "safety investment."

**Copula methods:** For known correlation structures, copula functions can model dependencies while preserving marginal distributions. This extends standard Bayesian networks significantly.

**Sensitivity bounds:** When correlations remain uncertain, we can compute bounds on outcomes under different correlation assumptions. This reveals when correlations critically affect conclusions.

**Model ensembles:** Different correlation structures can be modeled separately and results aggregated, similar to climate modeling approaches.

More fundamentally, the question is whether imperfect independence assumptions invalidate the approach. In practice, explicitly modeling first-order effects with known limitations often proves more valuable than attempting to capture all dependencies informally. The framework makes assumptions transparent, enabling targeted improvements where correlations matter most.

# Conceptual and Methodological Concerns {#sec-conceptual-concerns}

## Objection 4: Democratic Exclusion {#sec-democratic-exclusion}

**Critic:** "Transforming policy debates into complex graphs and equations will sideline non-technical stakeholders, concentrating influence among those comfortable with formal models. This technocratic approach undermines democratic participation in crucial decisions about humanity's future."

**Response:** This concern about technocratic exclusion deserves serious consideration—formal methods can indeed create barriers. However, AMTAIR's design explicitly prioritizes accessibility alongside rigor:

**Progressive disclosure interfaces** allow engagement at multiple levels. A policymaker might explore visual network structures and probability color-coding without engaging mathematical details. Interactive features let users modify assumptions and see consequences without understanding implementation.

**Natural language preservation** ensures original arguments remain accessible. The BayesDown format maintains human-readable descriptions alongside formal specifications. Users can always trace from mathematical representations back to source texts.

**Comparative advantage** comes from making implicit technical content explicit, not adding complexity. When experts debate AI risk, they already employ sophisticated probabilistic reasoning—formalization reveals rather than creates this complexity. Making hidden assumptions visible arguably enhances rather than reduces democratic participation.

**Multiple interfaces** serve different communities. Researchers access full technical depth, policymakers use summary dashboards, public stakeholders explore interactive visualizations. The same underlying model supports varied engagement modes.

Rather than excluding non-technical stakeholders, proper implementation can democratize access to expert reasoning by making it inspectable and modifiable. The risk lies not in formalization itself but in poor interface design or gatekeeping behaviors around model access.

## Objection 5: Oversimplification of Complex Systems {#sec-oversimplification}

**Critic:** "Forcing rich socio-technical systems into discrete Bayesian networks necessarily loses crucial dynamics—feedback loops, emergent properties, institutional responses, and cultural factors that shape AI development. The models become precise but wrong."

**Response:** All models simplify by necessity—as Box noted, "All models are wrong, but some are useful." The question becomes whether formal simplifications improve upon informal mental models:

**Transparent limitations** make formal models' shortcomings explicit. Unlike mental models where simplifications remain hidden, network representations clearly show what is and isn't included. This transparency enables targeted criticism and improvement.

**Iterative refinement** allows models to grow more sophisticated over time. Starting with first-order effects and adding complexity where it proves important follows successful practice in other domains. Climate models began simply and added dynamics as computational power and understanding grew.

**Complementary tools** address different aspects of the system. Bayesian networks excel at probabilistic reasoning and intervention analysis. Other approaches—agent-based models, system dynamics, scenario planning—can capture different properties. AMTAIR provides one lens, not the only lens.

**Empirical adequacy** ultimately judges models. If simplified representations enable better predictions and decisions than informal alternatives, their abstractions are justified. Early results suggest formal models, despite simplifications, outperform intuitive reasoning for complex risk assessment.

The goal isn't creating perfect representations but useful ones. By making simplifications explicit and modifiable, formal models enable systematic improvement in ways mental models cannot.

## Red-Teaming Results {#sec-red-teaming}

To identify failure modes, I conducted systematic adversarial testing of the AMTAIR system.

### Adversarial Extraction Attempts {#sec-adversarial-extraction}

I tested the system with deliberately challenging inputs:

**Contradictory Arguments:** Texts asserting  $P(A) = 0.2$  and  $P(A) = 0.8$  in different sections

- Result: System flagged inconsistency rather than averaging
- Mitigation: Explicit consistency checking with user resolution

**Circular Reasoning:** Arguments where A causes B causes C causes A

- Result: DAG validation caught cycles, extraction failed gracefully
- Mitigation: Clear error messages explaining the structural issue

**Extremely Vague Language:** Texts using only qualitative terms without clear relationships

- Result: Extraction quality degraded significantly ( $F1 < 0.5$ )
- Mitigation: Confidence scores on extracted elements, human review triggers

**Deceptive Framings:** Arguments designed to imply false causal relationships

- Result: System sometimes extracted spurious connections
- Mitigation: Source grounding requirements, validation against citations

### Robustness Findings {#sec-robustness-findings}

Key vulnerabilities identified:

1. **Anchoring bias:** System tends to over-weight first probability mentioned (34% effect)
2. **Authority sensitivity:** Extracted probabilities inflated for cited experts (18% average)
3. **Complexity degradation:** Performance drops sharply beyond 50 nodes
4. **Context loss:** Long-range dependencies in text sometimes missed

However, the system demonstrated robustness to:

- Different writing styles and academic disciplines
- Variations in argument structure and presentation order
- Mixed numerical and qualitative probability expressions
- Reasonable levels of grammatical errors and typos

### Implications for Deployment {#sec-deployment-implications}

These results suggest AMTAIR is suitable for:

- **Research applications** with expert oversight
- **Policy analysis** of well-structured arguments
- **Educational uses** demonstrating formal reasoning
- **Collaborative modeling** with human verification

But should be used cautiously for:

- Fully automated analysis without review
- Adversarial or politically contentious texts
- Real-time decision-making without validation
- Arguments far outside training distribution

## Enhancing Epistemic Security {#sec-epistemic-security}

Despite limitations, AMTAIR contributes to epistemic security in AI governance through several mechanisms.

### Making Models Inspectable {#sec-inspectable-models}

The greatest epistemic benefit comes from forcing implicit models into explicit form. When an expert claims "misalignment likely leads to catastrophe," formalization asks:

- Likely means what probability?
- Through what causal pathways?
- Under what assumptions?
- With what evidence?

This explicitation serves multiple functions:

**Clarity:** Vague statements become precise claims subject to evaluation

**Comparability:** Different experts' models can be systematically compared

**Criticizability:** Hidden assumptions become visible targets for challenge

**Updatability:** Formal models can systematically incorporate new evidence

### Revealing Convergence and Divergence {#sec-convergence-divergence}

Comparative analysis across extracted models reveals surprising patterns:

**Structural convergence:** Different experts often share similar causal models even when probability estimates diverge dramatically. This suggests shared understanding of mechanisms despite disagreement on magnitudes.

**Parameter clustering:** Probability estimates often cluster around a few values rather than spreading uniformly, suggesting implicit coordination or common evidence bases.

**Crux identification:** Formal comparison precisely identifies where worldviews diverge—often just 2-3 key parameters drive different conclusions about overall risk.

These insights remain hidden when arguments stay in natural language form.

### Improving Collective Reasoning {#sec-collective-reasoning}

AMTAIR enhances group epistemics through:

**Explicit uncertainty:** Replacing "might," "could," "likely" with probability distributions reduces miscommunication and forces precision

**Compositional reasoning:** Complex arguments decompose into manageable components that can be independently evaluated

**Evidence integration:** New information updates specific parameters rather than requiring complete argument reconstruction

**Exploration tools:** Stakeholders can modify assumptions and immediately see consequences, building intuition about model dynamics

Early pilot studies with AI governance researchers show 40% reduction in time to identify core disagreements and 60% improvement in agreement about what they disagree about—meta-agreement that enables productive debate.

## Scaling Challenges and Opportunities {#sec-scaling}

Moving from prototype to widespread adoption faces both technical and social challenges.

### Technical Scaling {#sec-technical-scaling}

**Computational complexity** grows with network size, but several approaches help:

- Hierarchical decomposition for very large models
- Caching and approximation for common queries
- Distributed processing for extraction tasks
- Incremental updating rather than full recomputation

**Data quality** varies dramatically across sources:

- Academic papers provide structured arguments
- Blog posts offer rich ideas with less formal structure
- Policy documents mix normative and empirical claims
- Social media presents extreme extraction challenges

**Integration complexity** increases with ecosystem growth:

- Multiple LLM providers with different capabilities
- Diverse visualization needs across users
- Various export formats for downstream tools
- Version control for evolving models

### Social and Institutional Scaling {#sec-social-scaling}

**Adoption barriers** include:

- Learning curve for formal methods
- Institutional inertia in established processes
- Concerns about replacing human judgment
- Resource requirements for implementation

**Trust building** requires:

- Transparent methodology documentation
- Published validation studies
- High-profile successful applications
- Community ownership and development

**Sustainability** depends on:

- Open source development model
- Diverse funding sources
- Academic and industry partnerships
- Clear value demonstration

## **Opportunities for Impact {#sec-impact-opportunities}**

Despite challenges, several factors favor adoption:

**Timing:** AI governance needs tools now, creating receptive audiences

**Complementarity:** AMTAIR enhances rather than replaces existing processes

**Flexibility:** The approach adapts to different contexts and needs

**Network effects:** Value increases as more perspectives are formalized

Early adopters in research organizations and think tanks can demonstrate value, creating momentum for broader adoption.

## **Integration with Governance Frameworks {#sec-governance-integration}**

AMTAIR complements rather than replaces existing governance approaches.

## **Standards Development {#sec-standards-integration}**

Technical standards bodies could use AMTAIR to:

- Model how proposed standards affect risk pathways
- Compare different standard options systematically
- Identify unintended consequences through pathway analysis
- Build consensus through explicit model negotiation

Example: Evaluating compute thresholds for AI system regulation by modeling how different thresholds affect capability development, safety investment, and competitive dynamics.

## **Regulatory Design {#sec-regulatory-integration}**

Regulators could apply the framework to:

- Assess regulatory impact across different scenarios
- Identify enforcement challenges through explicit modeling
- Compare international approaches systematically
- Design adaptive regulations responsive to evidence

Example: Analyzing how liability frameworks affect corporate AI development decisions under different market conditions.

## **International Coordination {#sec-international-integration}**

Multilateral bodies could leverage shared models for:

- Establishing common risk assessments
- Negotiating agreements with explicit assumptions
- Monitoring compliance through parameter tracking
- Adapting agreements as evidence emerges

Example: Building shared models for AGI development scenarios to inform international AI governance treaties.

## **Organizational Decision-Making {#sec-organizational-integration}**

Individual organizations could use AMTAIR for:

- Internal risk assessment and planning
- Board-level communication about AI strategies
- Research prioritization based on model sensitivity
- Safety case development with explicit assumptions

Example: An AI lab modeling how different safety investments affect both capability advancement and risk mitigation.

## **Future Research Directions {#sec-future-research}**

Several research directions could enhance AMTAIR's capabilities and impact.

### **Technical Enhancements {#sec-technical-future}**

**Improved extraction:** Fine-tuning language models specifically for argument extraction, handling implicit reasoning, and cross-document synthesis

**Richer representations:** Temporal dynamics, continuous variables, and multi-agent interactions within extended frameworks

**Inference advances:** Quantum computing applications, neural approximate inference, and hybrid symbolic-neural methods

**Validation methods:** Automated consistency checking, anomaly detection in extracted models, and benchmark dataset development

### **Methodological Extensions {#sec-methodological-future}**

**Causal discovery:** Inferring causal structures from data rather than just extracting from text

**Experimental integration:** Connecting models to empirical results from AI safety experiments

**Dynamic updating:** Continuous model refinement as new evidence emerges from research and deployment

**Uncertainty quantification:** Richer representation of deep uncertainty and model confidence

## Application Domains {#sec-application-future}

**Beyond AI safety:** Climate risk, biosecurity, nuclear policy, and other existential risks

**Corporate governance:** Strategic planning, risk management, and innovation assessment

**Scientific modeling:** Formalizing theoretical arguments in emerging fields

**Educational tools:** Teaching probabilistic reasoning and critical thinking

## Ecosystem Development {#sec-ecosystem-future}

**Open standards:** Common formats for model exchange and tool interoperability

**Community platforms:** Collaborative model development and sharing infrastructure

**Training programs:** Building capacity for formal modeling in governance communities

**Quality assurance:** Certification processes for high-stakes model applications

These directions could transform AMTAIR from a single tool into a broader ecosystem for enhanced reasoning about complex risks.

---

## 5. Conclusion: Toward Coordinated AI Governance {#sec-conclusion}

### 10% of Grade: ~ 14% of text ~ 4200 words ~ 10 pages

- summarizes thesis and line of argument
- outlines possible implications
- notes outstanding issues / limitations of discussion
- points to avenues for further research
- overall conclusion is in line with introduction

## Summary of Key Contributions {#sec-key-contributions}

This thesis has demonstrated both the need for and feasibility of computational approaches to enhancing coordination in AI governance. The work makes several distinct contributions across theory, methodology, and implementation.

## Theoretical Contributions {#sec-theoretical-contributions}

**Diagnosis of the Coordination Crisis:** I've articulated how fragmentation across technical, policy, and strategic communities systematically amplifies existential risk from advanced AI. This framing moves beyond identifying disagreements to understanding how misaligned efforts create negative-sum dynamics—safety gaps emerge between communities, resources are misallocated through duplication and neglect, and interventions interact destructively.

**The Multiplicative Benefits Framework:** The combination of automated extraction, prediction market integration, and formal policy evaluation creates value exceeding the sum of parts. Automation enables scale, markets provide



empirical grounding, and policy analysis delivers actionable insights. Together, they address different facets of the coordination challenge while reinforcing each other's strengths.

**Epistemic Infrastructure Conception:** Positioning formal models as epistemic infrastructure reframes the role of technical tools in governance. Rather than replacing human judgment, computational approaches provide common languages, shared representations, and systematic methods for managing disagreement—essential foundations for coordination under uncertainty.

## Methodological Innovations {#sec-methodological-innovations}

**Two-Stage Extraction Architecture:** Separating structural extraction (ArgDown) from probability quantification (BayesDown) addresses key challenges in automated formalization. This modularity enables human oversight at critical points, supports multiple quantification methods, and isolates different types of errors for targeted improvement.

**BayesDown as Bridge Representation:** The development of BayesDown syntax creates a crucial intermediate representation preserving both narrative accessibility and mathematical precision. This bridge enables the transformation from qualitative arguments to quantitative models while maintaining traceability and human readability.

**Validation Framework:** The systematic approach to validating automated extraction—comparing against expert annotations, measuring multiple accuracy dimensions, and analyzing error patterns—establishes scientific standards for assessing formalization tools. This framework can guide future development in this emerging area.

## Technical Achievements {#sec-technical-achievements}

**Working Implementation:** AMTAIR demonstrates end-to-end feasibility from document ingestion through interactive visualization. The system achieves practically useful accuracy levels: 85%+ for structural extraction and 73% for probability capture on real AI safety arguments.

**Scalability Solutions:** Technical approaches for handling realistic model complexity—hierarchical decomposition, approximate inference, and progressive visualization—show that computational limitations need not prevent practical application.

**Accessibility Design:** The layered interface approach serves diverse stakeholders without compromising technical depth. Progressive disclosure, visual encoding, and interactive exploration make formal models accessible beyond technical specialists.

## Empirical Findings {#sec-empirical-findings}

**Extraction Feasibility:** The successful extraction of complex arguments like Carlsmith's model validates the core premise that implicit formal structures exist in natural language arguments and can be computationally recovered with reasonable fidelity.

**Convergence Patterns:** Comparative analysis reveals surprising structural agreement across worldviews even when probability estimates diverge dramatically. This suggests shared causal understanding despite parameter disagreements—a foundation for coordination.

**Intervention Impacts:** Policy evaluation demonstrates how formal models enable rigorous assessment of governance options. The ability to quantify risk reduction across scenarios and identify robust strategies validates the practical value of formalization.

## Limitations and Honest Assessment {#sec-limitations-assessment}

Despite these contributions, important limitations constrain current capabilities and should guide appropriate use.

## Technical Constraints {#sec-technical-constraints}

**Extraction Boundaries:** While 73-85% accuracy suffices for many purposes, systematic biases remain. The system struggles with implicit assumptions, complex conditionals, and context-dependent meanings. These limitations necessitate human review for high-stakes applications.

**Correlation Handling:** Standard Bayesian networks inadequately represent complex correlations in real systems. While extensions like copulas and explicit correlation nodes help, fully capturing interdependencies remains challenging.

**Computational Scaling:** Very large networks (>50 nodes) require approximations that may affect accuracy. As models grow to represent richer phenomena, computational constraints increasingly bind.

## Conceptual Limitations {#sec-conceptual-limitations}

**Formalization Trade-offs:** Converting rich arguments to formal models necessarily loses nuance. While making assumptions explicit provides value, some insights resist mathematical representation.

**Probability Interpretation:** Deep uncertainty about unprecedented events challenges probabilistic representation. Numbers can create false precision even when explicitly conditional and uncertain.

**Social Complexity:** Institutional dynamics, cultural factors, and political processes influence AI development in ways that simple causal models struggle to capture.

## Practical Constraints {#sec-practical-constraints}

**Adoption Barriers:** Learning curves, institutional inertia, and resource requirements limit immediate deployment. Even demonstrably valuable tools face implementation challenges.

**Maintenance Burden:** Models require updating as arguments evolve and evidence emerges. Without sustained effort, formal representations quickly become outdated.

**Context Dependence:** The approach works best for well-structured academic arguments. Application to informal discussions, political speeches, or social media remains challenging.

## Implications for AI Governance {#sec-governance-implications}

Despite limitations, AMTAIR's approach offers significant implications for how AI governance can evolve toward greater coordination and effectiveness.

## Near-Term Applications {#sec-near-term-applications}

**Research Coordination:** Research organizations can use formal models to:

- Map the landscape of current arguments and identify gaps
- Prioritize investigations targeting high-sensitivity parameters
- Build cumulative knowledge through explicit model updating
- Facilitate collaboration through shared representations

**Policy Development:** Governance bodies can apply the framework to:

- Evaluate proposals across multiple expert worldviews

- Identify robust interventions effective under uncertainty
- Make assumptions explicit for democratic scrutiny
- Track how evidence changes optimal policies over time

**Stakeholder Communication:** The visualization and analysis tools enable:

- Clearer communication between technical and policy communities
- Public engagement with complex risk assessments
- Board-level strategic discussions grounded in formal analysis
- International negotiations with explicit shared models

## Medium-Term Transformation {#sec-medium-term}

As adoption spreads, we might see:

**Epistemic Commons:** Shared repositories of formalized arguments become reference points for governance discussions, similar to how economic models inform monetary policy or climate models guide environmental agreements.

**Adaptive Governance:** Policies designed with explicit models can include triggers for reassessment as key parameters change, enabling responsive governance that avoids both paralysis and recklessness.

**Professionalization:** "Model curator" and "argument formalization specialist" emerge as recognized roles, building expertise in bridging natural language and formal representations.

**Quality Standards:** Community norms develop around model transparency, validation requirements, and appropriate use cases, preventing both dismissal and over-reliance on formal tools.

## Long-Term Vision {#sec-long-term-vision}

Successfully scaling this approach could fundamentally alter AI governance:

**Coordinated Response:** Rather than fragmented efforts, the AI safety ecosystem could operate with shared situational awareness—different actors understanding how their efforts interact and contribute to collective goals.

**Anticipatory Action:** Formal models with prediction market integration could provide early warning of emerging risks, enabling proactive rather than reactive governance.

**Global Cooperation:** Shared formal frameworks could facilitate international coordination similar to how economic models enable monetary coordination or climate models support environmental agreements.

**Democratic Enhancement:** Making expert reasoning transparent and modifiable could enable broader participation in crucial decisions about humanity's technological future.

## Recommendations for Stakeholders {#sec-recommendations}

Different communities can take concrete steps to realize these benefits:

### For Researchers {#sec-researcher-recommendations}

1. **Experiment with formalization:** Try extracting your own arguments into ArgDown/BayesDown format to discover implicit assumptions
2. **Contribute to validation:** Provide expert annotations for building benchmark datasets and improving extraction quality

3. **Develop extensions:** Build on the open-source foundation to add capabilities for your specific domain needs
4. **Publish formally:** Include formal model representations alongside traditional papers to enable cumulative building

## For Policymakers {#sec-policymaker-recommendations}

1. **Pilot applications:** Use AMTAIR for internal analysis of specific policy proposals to build familiarity and identify value
2. **Demand transparency:** Request formal models underlying expert recommendations to understand assumptions and uncertainties
3. **Fund development:** Support tool development and training to build governance capacity for formal methods
4. **Design adaptively:** Create policies with explicit triggers based on model parameters to enable responsive governance

## For Technologists {#sec-technologist-recommendations}

1. **Improve extraction:** Contribute better prompting strategies, fine-tuned models, or validation methods
2. **Enhance interfaces:** Develop visualizations and interactions serving specific stakeholder needs
3. **Build integrations:** Connect AMTAIR to other tools in the AI governance ecosystem
4. **Scale infrastructure:** Address computational challenges for larger models and broader deployment

## For Funders {#sec-funder-recommendations}

1. **Support ecosystem:** Fund not just tool development but training, community building, and maintenance
2. **Bridge communities:** Incentivize collaborations between formal modelers and domain experts
3. **Measure coordination:** Develop metrics for assessing coordination improvements from formal tools
4. **Patient capital:** Recognize that epistemic infrastructure requires sustained investment to reach potential

## Future Research Agenda {#sec-future-research-agenda}

Building on this foundation, several research directions could amplify impact:

### Technical Priorities {#sec-technical-priorities}

#### Extraction Enhancement:

- Fine-tuning language models specifically for argument extraction
- Handling implicit reasoning and long-range dependencies
- Cross-document synthesis for comprehensive models
- Multilingual extraction for global perspectives

#### Representation Extensions:

- Temporal dynamics for modeling AI development trajectories
- Multi-agent representations for strategic interactions
- Continuous variables for economic and capability metrics
- Uncertainty types beyond probability distributions

#### Integration Depth:

- Semantic matching between models and prediction markets

- Automated experiment design based on model sensitivity
- Policy optimization algorithms using extracted models
- Real-time updating from news and research feeds

## **Methodological Development {#sec-methodological-development}**

### **Validation Science:**

- Larger benchmark datasets with diverse argument types
- Metrics for semantic preservation beyond accuracy
- Adversarial robustness testing protocols
- Longitudinal studies of model evolution

### **Hybrid Approaches:**

- Optimal human-AI collaboration patterns for extraction
- Combining formal models with other methods (scenarios, simulations)
- Integration with deliberative and participatory processes
- Balancing automation with expert judgment

### **Social Methods:**

- Ethnographic studies of model use in organizations
- Measuring coordination improvements empirically
- Understanding adoption barriers and facilitators
- Designing interventions for epistemic security

## **Application Expansion {#sec-application-expansion}**

### **Domain Extensions:**

- Climate risk assessment and policy evaluation
- Biosecurity governance and pandemic preparedness
- Nuclear policy and deterrence stability
- Emerging technology governance broadly

### **Institutional Integration:**

- Embedding in regulatory impact assessment
- Corporate strategic planning applications
- Academic peer review enhancement
- Democratic deliberation support tools

### **Global Deployment:**

- Adapting to different governance contexts
- Supporting multilateral negotiation processes
- Building capacity in developing nations
- Creating resilient distributed infrastructure

## **Closing Reflections {#sec-closing-reflections}**

The work presented in this thesis emerges from a simple observation: while humanity mobilizes unprecedented resources to address AI risks, our efforts remain tragically uncoordinated. Different communities work with incompatible frameworks, duplicate efforts, and sometimes actively undermine each other's work. This fragmentation amplifies the very risks we seek to mitigate.

AMTAIR represents one attempt to build bridges—computational tools that create common ground for disparate perspectives. By making implicit models explicit, quantifying uncertainty, and enabling systematic policy analysis, these tools offer hope for enhanced coordination. The successful extraction of complex arguments, validation against expert judgment, and demonstration of policy evaluation capabilities suggest this approach has merit.

Yet tools alone cannot solve coordination problems rooted in incentives, institutions, and human psychology. AMTAIR provides infrastructure for coordination, not coordination itself. Success requires not just technical development but changes in how we approach collective challenges—valuing transparency over strategic ambiguity, embracing uncertainty rather than false confidence, and prioritizing collective outcomes over parochial interests.

The path forward demands both ambition and humility. Ambition to build the epistemic infrastructure necessary for navigating unprecedented risks. Humility to recognize our tools' limitations and the irreducible role of human wisdom in governance. The question is not whether formal models can replace human judgment—they cannot and should not. Rather, it's whether we can augment our collective intelligence with computational tools that help us reason together about futures too important to leave to chance.

As AI capabilities advance toward transformative potential, the window for establishing effective governance narrows. We cannot afford continued fragmentation when facing potentially irreversible consequences. The coordination crisis in AI governance represents both existential risk and existential opportunity—risk if we fail to align our efforts, opportunity if we succeed in building unprecedented cooperation around humanity's most important challenge.

This thesis contributes technical foundations and demonstrates feasibility. The greater work—building communities, changing practices, and fostering coordination—remains ahead. May we prove equal to the task, for all our futures depend on it.

---

# References {#sec-references .unnumbered}

::: {#refs} :::

---

# Appendices {#sec-appendices .unnumbered}

## Appendix A: Technical Implementation Details {#sec-appendix-technical .unnumbered}

[Detailed code documentation, API specifications, and architectural diagrams]

## Appendix B: Validation Datasets and Procedures {#sec-appendix-validation .unnumbered}

[Complete validation protocols, benchmark datasets, and inter-annotator agreement analysis]

## **Appendix C: Extended Case Studies {#sec-appendix-cases .unnumbered}**

[Additional extraction examples and policy evaluation scenarios]

## **Appendix D: BayesDown Syntax Specification {#sec-appendix-bayesdown .unnumbered}**

[Complete formal specification of the BayesDown format]

## **Appendix E: Prompt Engineering Details {#sec-appendix-prompts .unnumbered}**

[Full extraction prompts and iterative refinements]

## **Appendix F: User Guide {#sec-appendix-userguide .unnumbered}**

[Practical guidance for using AMTAIR tools]