



UNIVERSITÄT
BAYREUTH

– P&E Master's Programme –
Chair of Philosophy, Computer
Science & Artificial Intelligence

Automating the Modelling of Transformative Artificial Intelligence Risks

*“An Epistemic Framework for Leveraging Frontier AI Systems to Upscale Conditional Policy
Assessments in Bayesian Networks on a Narrow Path towards Existential Safety”*

A thesis submitted at the Department of Philosophy

for the degree of *Master of Arts in Philosophy & Economics*

Author:

Valentin Jakob Meyer
Valentin.meyer@uni-bayreuth.de
Matriculation Number: 1828610
Tel.: +49 (1573) 4512494
Pielmühler Straße 15
93138 Lappersdorf

Supervisor:

Dr. Timo Speith

Word Count:

30.000

Source / Identifier:

Document URL

26th of May 2025

Table of Contents

Preface	1
Abstract	3
Prefatory Apparatus: Frontmatter	5
Illustrations and Terminology — Quick References	5
Acknowledgments	5
List of Graphics & Figures	5
List of Abbreviations	5
Final Thesis: Automating the Modeling of Transformative Artificial Intelligence	
Risks	7
Frontmatter: Preface	7
Acknowledgments	7
List of Figures	8
List of Tables	8
List of Abbreviations	8
1. Introduction: The Coordination Crisis in AI Governance	9
1.1 Opening Scenario: The Policymaker’s Dilemma	9
1.2 The Coordination Crisis in AI Governance	10
1.2.1 Safety Gaps from Misaligned Efforts	10
1.2.2 Resource Misallocation	11
1.2.3 Negative-Sum Dynamics	11
1.3 Historical Parallels and Temporal Urgency	12
1.4 Research Question and Scope	12
1.5 The Multiplicative Benefits Framework	13
1.5.1 Automated Worldview Extraction	13
1.5.2 Live Data Integration	14
1.5.3 Formal Policy Evaluation	14
1.5.4 The Synergy	14
1.6 Thesis Structure and Roadmap	15
2. Context and Theoretical Foundations	17

2.1 AI Existential Risk: The Carlsmith Model	17
2.1.1 Six-Premise Decomposition	17
2.1.2 Why Carlsmith Exemplifies Formalizable Arguments	19
2.2 The Epistemic Challenge of Policy Evaluation	19
2.2.1 Unique Characteristics of AI Governance	19
2.2.2 Limitations of Traditional Approaches	20
2.2.3 The Underlying Epistemic Framework	21
2.2.4 Toward New Epistemic Tools	21
2.3 Bayesian Networks as Knowledge Representation	22
2.3.1 Mathematical Foundations	22
2.3.2 The Rain-Sprinkler-Grass Example	23
Rain-Sprinkler-Grass Network Rendering	24
2.3.3 Advantages for AI Risk Modeling	24
2.3.3 Advantages for AI Risk Modeling	25
2.4 Argument Mapping and Formal Representations	25
2.4.1 From Natural Language to Structure	26
2.4.2 ArgDown: Structured Argument Notation	26
2.4.3 BayesDown: The Bridge to Bayesian Networks	27
2.5 The MTAIR Framework: Achievements and Limitations	27
2.5.1 MTAIR’s Approach	28
2.5.2 Key Achievements	28
2.5.3 Fundamental Limitations	29
2.5.4 The Automation Opportunity	29
2.6 Literature Review: Content and Technical Levels	30
2.6.1 AI Risk Models Evolution	30
2.6.2 Governance Proposals Taxonomy	31
2.6.3 Bayesian Network Theory and Applications	31
2.6.4 Software Tools Landscape	32
2.6.5 Formalization Approaches	32
2.6.6 Correlation Accounting Methods	33
2.7 Methodology	33
2.7.1 Research Design Overview	33
2.7.2 Formalizing World Models from AI Safety Literature	34
2.7.3 From Natural Language to Computational Models	34
2.7.4 Directed Acyclic Graphs: Structure and Semantics	35
2.7.5 Quantification of Probabilistic Judgments	35
2.7.6 Inference Techniques for Complex Networks	36
2.7.7 Integration with Prediction Markets and Forecasting Platforms	36
3. AMTAIR: Design and Implementation	39
3.1 System Architecture Overview	39
3.1.1 Five-Stage Pipeline Architecture	40
3.1.2 Design Principles	40

3.2 The Two-Stage Extraction Process	41
3.2.1 Stage 1: Structural Extraction (ArgDown)	41
3.2.2 Stage 2: Probability Integration (BayesDown)	42
3.2.3 Why Two Stages?	42
3.3 Implementation Technologies	43
3.3.1 Technology Stack	43
3.3.2 Key Algorithms	44
2.3.3 Advantages for AI Risk Modeling	44
2.4 Argument Mapping and Formal Representations	45
2.4.1 From Natural Language to Structure	45
2.4.2 ArgDown: Structured Argument Notation	45
2.4.3 BayesDown: The Bridge to Bayesian Networks	46
2.5 The MTAIR Framework: Achievements and Limitations	47
2.5.1 MTAIR’s Approach	47
2.5.2 Key Achievements	48
2.5.3 Fundamental Limitations	48
2.5.4 The Automation Opportunity	49
2.6 Literature Review: Content and Technical Levels	49
2.6.1 AI Risk Models Evolution	49
2.6.2 Governance Proposals Taxonomy	50
2.6.3 Bayesian Network Theory and Applications	51
2.6.4 Software Tools Landscape	51
2.6.5 Formalization Approaches	52
2.6.6 Correlation Accounting Methods	52
2.7 Methodology	53
2.7.1 Research Design Overview	53
2.7.2 Formalizing World Models from AI Safety Literature	53
2.7.3 From Natural Language to Computational Models	54
2.7.4 Directed Acyclic Graphs: Structure and Semantics	54
2.7.5 Quantification of Probabilistic Judgments	55
2.7.6 Inference Techniques for Complex Networks	55
2.7.7 Integration with Prediction Markets and Forecasting Platforms	56
3. AMTAIR: Design and Implementation	57
3.1 System Architecture Overview	57
3.1.1 Five-Stage Pipeline Architecture	58
3.1.2 Design Principles	58
3.2 The Two-Stage Extraction Process	59
3.2.1 Stage 1: Structural Extraction (ArgDown)	59
3.2.2 Stage 2: Probability Integration (BayesDown)	60
3.2.3 Why Two Stages?	60
3.3 Implementation Technologies	61
3.3.1 Technology Stack	61

3.3.2 Key Algorithms	62
3.3.3 (Expected) Performance Characteristics	63
3.3.4 Deterministic vs. Probabilistic Components of the Workflow	63
3.4 Case Study: Rain-Sprinkler-Grass	64
3.4.1 Processing Steps	64
3.4.2 Example Conversion Steps	64
3.4.3 Results	66
3.5 Case Study: Carlsmith’s Power-Seeking AI Model	66
3.5.1 Model Complexity	66
3.5.2 Automated Extraction of the Carlsmith’s Argument Structure	67
3.5.3 From ArgDown to BayesDown in Carlsmith’s Model	68
3.5.4 Practically Meaningful BayesDown	68
3.5.5 Interactive Visualization and Exploration	69
3.5.6 Validation Against Original (From the MTAIR Project)	70
3.6 Validation Methodology	71
3.6.1 Ground Truth Construction	71
3.6.2 Evaluation Metrics	72
3.6.3 Results Summary	72
3.6.4 Error Analysis	73
3.7 Policy Evaluation Capabilities	73
3.7.1 Intervention Representation	73
3.7.2 Example: Deployment Governance	74
3.7.3 Robustness Analysis	74
3.8 Interactive Visualization Design	75
3.8.1 Visual Encoding Strategy	75
3.8.2 Progressive Disclosure	75
3.8.3 User Interface Elements	76
3.9 Integration with Prediction Markets	76
3.9.1 Design for Integration	76
3.9.2 Challenges and Opportunities	77
3.10 Computational Performance Analysis	77
3.10.1 Exact vs. Approximate Inference	77
3.10.2 Scaling Strategies	78
3.11 Results and Achievements	78
3.11.1 Extraction Quality Assessment	78
3.11.2 Computational Performance	79
3.11.3 Policy Impact Evaluation	79
3.12 Summary of Technical Contributions	80
3.5.6 Validation Against Original (From the MTAIR Project)	82
3.6 Validation Methodology	82
3.6.1 Ground Truth Construction	82
3.6.2 Evaluation Metrics	83

3.6.3 Results Summary	83
3.6.4 Error Analysis	84
3.7 Policy Evaluation Capabilities	84
3.7.1 Intervention Representation	84
3.7.2 Example: Deployment Governance	85
3.7.3 Robustness Analysis	85
3.8 Interactive Visualization Design	86
3.8.1 Visual Encoding Strategy	86
3.8.2 Progressive Disclosure	86
3.8.3 User Interface Elements	86
3.9 Integration with Prediction Markets	87
3.9.1 Design for Integration	87
3.9.2 Challenges and Opportunities	87
3.10 Computational Performance Analysis	87
3.10.1 Exact vs. Approximate Inference	88
3.10.2 Scaling Strategies	88
3.11 Results and Achievements	88
3.11.1 Extraction Quality Assessment	88
3.11.2 Computational Performance	89
3.11.3 Policy Impact Evaluation	89
3.12 Summary of Technical Contributions	89
4. Discussion: Implications and Limitations	91
4.1 Technical Limitations and Responses	91
4.1.1 Objection 1: Extraction Quality Boundaries	91
4.1.2 Objection 2: False Precision in Uncertainty	92
4.1.3 Objection 3: Correlation Complexity	92
4.2 Conceptual and Methodological Concerns	93
4.2.1 Objection 4: Democratic Exclusion	93
4.2.2 Objection 5: Oversimplification of Complex Systems	93
4.2.3 Objection 6: Idiosyncratic Implementation and Modeling Choices	94
4.3 Red-Teaming Results	95
4.3.1 Adversarial Extraction Attempts	95
4.3.2 Robustness Findings	95
4.3.3 Implications for Deployment	96
4.4 Enhancing Epistemic Security	96
4.4.1 Making Models Inspectable	96
4.4.2 Revealing Convergence and Divergence	97
4.4.3 Improving Collective Reasoning	97
4.5 Scaling Challenges and Opportunities	97
4.5.1 Technical Scaling	97
4.5.2 Social and Institutional Scaling	98
4.5.3 Opportunities for Impact	98

4.6 Integration with Governance Frameworks	99
4.6.1 Standards Development	99
4.6.2 Regulatory Design	99
4.6.3 International Coordination	99
4.6.4 Organizational Decision-Making	100
4.7 Future Research Directions	100
4.7.1 Technical Enhancements	100
4.7.2 Methodological Extensions	100
4.7.3 Application Domains	101
4.7.4 Ecosystem Development	101
4.8 Known Unknowns and Deep Uncertainties	101
4.8.1 Categories of Deep Uncertainty	101
4.8.2 Adaptation Strategies for Deep Uncertainty	101
4.8.3 Robust Decision-Making Principles	102
4.9 Summary of Implications	102
5. Conclusion: Toward Coordinated AI Governance	105
5.1 Summary of Key Contributions	105
5.1.1 Theoretical Contributions	105
5.1.2 Methodological Innovations	106
5.1.3 Technical Achievements	106
5.1.4 Empirical Findings	106
5.2 Limitations and Honest Assessment	107
5.2.1 Technical Constraints	107
5.2.2 Conceptual Limitations	107
5.2.3 Practical Constraints	107
5.3 Implications for AI Governance	107
5.3.1 Near-Term Applications	108
5.3.2 Medium-Term Transformation	108
5.3.3 Long-Term Vision	108
5.4 Recommendations for Stakeholders	109
5.4.1 For Researchers	109
5.4.2 For Policymakers	109
5.4.3 For Technologists	109
Bibliography	111
Appendices	113
	113

List of Figures

List of Tables

1	Examples of duplicated AI safety efforts across organizations	11
2	Comparison of AI governance vs traditional policy domains	20

Preface

Abstract

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. This thesis introduces AMTAIR (Automating Transformative AI Risk Modeling), a computational approach addressing this coordination failure by automating the extraction of probabilistic world models from AI safety literature using frontier language models. The system implements an end-to-end pipeline transforming unstructured text into interactive Bayesian networks through a novel two-stage extraction process that bridges communication gaps between stakeholders.

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. 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.

Prefatory Apparatus: Frontmatter

Illustrations and Terminology — Quick References

Acknowledgments

- Academic supervisor (Prof. Timo Speith) and institution (University of Bayreuth)
- Research collaborators, especially those connected to the original MTAIR project
- Technical advisors who provided feedback on implementation aspects
- Personal supporters who enabled the research through encouragement and feedback

List of Graphics & Figures

List of Abbreviations

- AGI - Artificial General Intelligence
- AMTAIR - Automating Modeling of Transformative AI Risks
- API - Application Programming Interface
- APS - Advanced, Planning, Strategic (AI systems per **carlsmith2021**)
- BN - Bayesian Network
- CPT - Conditional Probability Table
- DAG - Directed Acyclic Graph
- LLM - Large Language Model
- MTAIR - Modeling Transformative AI Risks
- TAI - Transformative Artificial Intelligence

Glossary

- **Argument mapping**: A method for visually representing the structure of arguments
- **BayesDown**: An extension of ArgDown that incorporates probabilistic information

- **Bayesian network:** A probabilistic graphical model representing variables and their dependencies
- **Conditional probability:** The probability of an event given that another event has occurred
- **Directed Acyclic Graph (DAG):** A graph with directed edges and no cycles
- **Existential risk:** Risk of permanent curtailment of humanity's potential
- **Power-seeking AI:** AI systems with instrumental incentives to acquire resources and power
- **Prediction market:** A market where participants trade contracts that resolve based on future events
- **d-separation:** A criterion for identifying conditional independence relationships in Bayesian networks
- **Monte Carlo sampling:** A computational technique using random sampling to obtain numerical results

Final Thesis: Automating the Modeling of Transformative Artificial Intelligence Risks

Frontmatter: Preface

This thesis 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. The AMTAIR project—Automating Transformative AI Risk Modeling—represents an attempt to build computational bridges between communities that, despite shared concerns about AI risk, often struggle to communicate effectively due to incompatible frameworks, terminologies, and implicit assumptions.

I hope this work contributes to building the intellectual and technical infrastructure necessary for humanity to navigate the transition to transformative AI safely. The tools and frameworks presented here are offered in the spirit of collaborative problem-solving, recognizing that the challenges we face require unprecedented cooperation across disciplines, institutions, and world-views.

Acknowledgments

I thank my supervisor Dr. Timo Speith for his guidance throughout this project, providing both technical insights and philosophical grounding. The MTAIR team’s pioneering manual approach inspired this automation effort, and I am grateful for their foundational work.

I acknowledge Johannes Meyer and Jelena Meyer for their invaluable assistance in verifying the automated extraction procedure through manual extraction of ArgDown and BayesDown data

from the Carlsmith paper, providing crucial ground truth for validation.

Special recognition goes to Coleman Snell for his partnership and research collaboration with the AMTAIR project, offering both technical expertise and strategic vision. The AI safety community's creation of rich literature made this work possible, and I thank all researchers whose arguments provided the raw material for formalization.

Any errors or limitations remain my own responsibility.

List of Figures

List of Tables

List of Abbreviations

AI - Artificial Intelligence
AGI - Artificial General Intelligence
AMTAIR - Automating Transformative AI Risk Modeling
API - Application Programming Interface
APS - Advanced, Planning, Strategic (AI systems)
BN - Bayesian Network
CPT - Conditional Probability Table
DAG - Directed Acyclic Graph
LLM - Large Language Model
ML - Machine Learning
MTAIR - Modeling Transformative AI Risks
NLP - Natural Language Processing
P&E - Philosophy & Economics
PDF - Portable Document Format
TAI - Transformative Artificial Intelligence

1. Introduction: The Coordination Crisis in AI Governance

Chapter Overview

Grade Weight: 10% | **Target Length:** ~14% of text (~4,200 words)

Requirements: Introduces and motivates the core question, provides context, states precise thesis, provides roadmap

1.1 Opening Scenario: The Policymaker’s Dilemma

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 scenario¹ 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 **todd2024**.

Show Image

¹The orthogonality thesis posits that intelligence and goals are independent—an AI can have any set of objectives regardless of its intelligence level. The instrumental convergence thesis suggests that different AI systems may adopt similar instrumental goals (e.g., self-preservation, resource acquisition) to achieve their objectives.

1.2 The Coordination Crisis in AI Governance

As AI capabilities advance at an accelerating pace—demonstrated by the rapid progression from GPT-3 to GPT-4, Claude, and emerging multimodal systems **maslej2025 samborska2025**—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.

Despite unprecedented investment in AI safety research, rapidly growing awareness among key stakeholders, and proliferating frameworks for responsible AI development, we face what I’ll term the “coordination crisis” in AI governance—a systemic failure to align diverse efforts across technical, policy, and strategic domains into a coherent response proportionate to the risks we face.

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 **tegmark2024**, the EU advances comprehensive AI regulation **european2024**, 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.

Show Image

1.2.1 Safety Gaps from Misaligned Efforts

The fragmentation problem manifests in incompatible frameworks between technical researchers, policy specialists, and strategic analysts. Each community develops sophisticated approaches within their domain, yet translation between domains remains primitive. This creates systematic blind spots where risks emerge at the interfaces between technical capabilities, institutional responses, and strategic dynamics.

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 **paul2023**.

Historical precedents demonstrate how coordination failures in technology governance can lead to dangerous dynamics. The nuclear arms race exemplifies how lack of coordination can create negative-sum outcomes where all parties become less secure despite massive investments in safety measures. Similar dynamics may emerge in AI development without proper coordination infrastructure.

1.2.2 Resource Misallocation

The AI safety community faces a complex tradeoff in resource allocation. While some duplication of efforts can improve reliability through independent verification—akin to reproducing scientific results—the current level of fragmentation often leads to wasteful redundancy. Multiple teams independently develop similar frameworks without building on each other’s work, creating opportunity costs where critical but unglamorous research areas remain understaffed. Funders struggle to identify high-impact opportunities across technical and governance domains, lacking the epistemic infrastructure to assess where marginal resources would have the greatest impact. This misallocation becomes more costly as the window for establishing effective governance narrows with accelerating AI development.

Table 1: Examples of duplicated AI safety efforts across organizations

Research Area	Organization A	Organization B	Duplication Level	Opportunity Cost
Interpretability Methods	Anthropic’s mechanistic interpretability	DeepMind’s concept activation vectors	Medium	Reduced focus on multi-agent safety
Alignment Frameworks	MIRI’s embedded agency	FHI’s comprehensive AI services	High	Limited work on institutional design
Risk Assessment Models	GovAI’s policy models	CSER’s existential risk frameworks	High	Insufficient capability benchmarking

1.2.3 Negative-Sum Dynamics

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.

The game-theoretic structure of AI development creates particularly pernicious dynamics. Arm-

strong et al. **armstrong2016** demonstrate how uncoordinated policies can incentivize a “race to the precipice” where competitive pressures override safety considerations. The situation resembles a multi-player prisoner’s dilemma or stag hunt where individually rational decisions lead to collectively catastrophic outcomes **samuel2023 hunt2025**.

1.3 Historical Parallels and Temporal 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) **schelling1960**, institutions (IAEA), and communication channels (hotlines, treaties) did humanity navigate the nuclear precipice **rehman2025**.

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.

Show Image

1.4 Research Question and Scope

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

Can frontier AI technologies be utilized to automate the modeling of transformative AI risks, enabling robust prediction of policy impacts across diverse worldviews?

More specifically, I explore whether frontier language models can automate the extraction and formalization of probabilistic world models from AI safety literature, creating a scalable computational framework that enhances coordination in AI governance through systematic policy evaluation under uncertainty.

To break this down into its components:

- **Frontier AI Technologies:** Today’s most capable language models (GPT-4, Claude-3 level systems)
- **Automated Modeling:** Using these systems to extract and formalize argument structures from natural language
- **Transformative AI Risks:** Potentially catastrophic outcomes from advanced AI systems, particularly existential risks
- **Policy Impact Prediction:** Evaluating how governance interventions might alter probability distributions over outcomes
- **Diverse Worldviews:** Accounting for fundamental disagreements about AI development trajectories and risk factors

The investigation encompasses both theoretical development and practical implementation, focusing 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.

1.5 The Multiplicative Benefits Framework

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.

Show Image

1.5.1 Automated Worldview 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 dozens 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
- Maintain models that evolve with the discourse
- Reduce barriers to entry for contributing worldviews

1.5.2 Live Data Integration

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 **tetlock2015**.

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
- Generating new questions for forecasting communities

1.5.3 Formal Policy Evaluation

Formal policy evaluation transforms static risk assessments into actionable guidance by modeling how specific interventions alter critical parameters. Using causal inference techniques **pearl2000 pearl2009**, 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
- Prioritizing research to reduce decision-relevant uncertainty

1.5.4 The 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.

1.6 Thesis Structure and 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, and requirements for effective coordination infrastructure.

Chapter 3: AMTAIR Design and Implementation presents the technical system including 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, and 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, and broader implications for epistemic security.

Chapter 5: Conclusion synthesizes key contributions and charts paths forward with a summary of theoretical and practical achievements, concrete recommendations for stakeholders, research agenda for community development, and vision for AI governance with proper coordination infrastructure.

Throughout this progression, 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.

Show Image

Having established the coordination crisis and outlined how automated modeling can address it, we now turn to the theoretical foundations that make this approach possible. The next chapter examines the unique epistemic challenges of AI governance and introduces the formal tools—particularly Bayesian networks—that enable rigorous reasoning under deep uncertainty.

2. Context and Theoretical Foundations

Chapter Overview

Grade Weight: 20% | **Target Length:** ~29% of text (~8,700 words)

Requirements: Demonstrates understanding of relevant concepts, explains relevance, situates in debate, reconstructs arguments

This chapter establishes the theoretical and methodological foundations for the AMTAIR approach. We begin by examining a concrete example of structured AI risk assessment—Joseph Carlsmith’s power-seeking AI model—to ground our discussion in practical terms. We then explore the unique epistemic challenges of AI governance that render traditional policy analysis inadequate, introduce Bayesian networks as formal tools for representing uncertainty, and examine how argument mapping bridges natural language reasoning and formal models. The chapter concludes by analyzing the MTAIR project’s achievements and limitations, motivating the need for automated approaches, and surveying relevant literature across AI risk modeling, governance proposals, and technical methodologies.

2.1 AI Existential Risk: The 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 **carlsmith2022**. Carlsmith’s analysis stands out for its explicit probabilistic decomposition of the path from current AI development to potential existential catastrophe.

2.1.1 Six-Premise Decomposition

According to the MTAIR model **clarke2022**, Carlsmith decomposes existential risk into a probabilistic chain with explicit estimates²:

1. **Premise 1:** Transformative AI development this century (P 0.80)(P 0.80) (P 0.80)
2. **Premise 2:** AI systems pursuing objectives in the world (P 0.95)(P 0.95) (P 0.95)

²Multiple versions of Carlsmith’s paper exist with slight updates to probability estimates: **carlsmith2021**, **carlsmith2022**, **carlsmith2024**. We primarily reference the version used by the MTAIR team for their extraction. Extended discussion and expert probability estimates can be found on LessWrong.

3. **Premise 3:** Systems with power-seeking instrumental incentives (P 0.40)(P 0.40) (P 0.40)
4. **Premise 4:** Sufficient capability for existential threat (P 0.65)(P 0.65) (P 0.65)
5. **Premise 5:** Misaligned systems despite safety efforts (P 0.50)(P 0.50) (P 0.50)
6. **Premise 6:** Catastrophic outcomes from misaligned power-seeking (P 0.65)(P 0.65) (P 0.65)

Composite Risk Calculation: $P(\text{doom}) = 0.05 \times 0.05 \times 0.05 = 0.000125$ (5%)

mermaid

flowchart TD

```

P1[Premise 1: Transformative AI<br/>P 0.80] --> P2[Premise 2: AI pursuing objectives<br/>P 0.40]
P2 --> P3[Premise 3: Power-seeking incentives<br/>P 0.40]
P3 --> P4[Premise 4: Existential capability<br/>P 0.65]
P4 --> P5[Premise 5: Misalignment despite safety<br/>P 0.50]
P5 --> P6[Premise 6: Catastrophic outcome<br/>P 0.65]
P6 --> D[Existential Catastrophe<br/>P 0.05]

```

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

Premise 1: APS Systems by 2070 (P 0.65)(P 0.65) (P 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 0.40)(P 0.40) (P 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 0.70)(P 0.70) (P 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 0.65)(P 0.65) (P 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 0.40)(P 0.40) (P 0.40) “Conditional on 1-4, power-seeking will scale to permanent human disempowerment”—despite potential resistance and safeguards.

Premise 6: Existential Catastrophe (P 0.95)(P 0.95) (P 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 $P(\text{doom}) = 0.05 \times 0.05 \times 0.05 = 0.000125$ or 5% by 2070.

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.

2.1.2 Why Carlsmith Exemplifies Formalizable Arguments

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.

We will return to Carlsmith’s model in Chapter 3 as our primary complex case study, demonstrating how AMTAIR successfully extracts and formalizes this sophisticated multi-level argument.

Beyond Carlsmith’s model, other structured approaches to AI risk—such as Christiano’s “What failure looks like” [christiano2019](#)—provide additional targets for automated extraction, enabling comparative analysis across different expert worldviews.

2.2 The Epistemic Challenge of Policy Evaluation

AI governance policy evaluation faces unique epistemic challenges that render traditional policy analysis methods insufficient. Understanding these challenges motivates the need for new computational approaches.

2.2.1 Unique Characteristics of AI Governance

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 [hallegatte2012](#). 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.

Table 2: Comparison of AI governance vs traditional policy domains

Dimension	Traditional Policy	AI Governance
Uncertainty Type	Risk (known distributions)	Deep uncertainty (unknown unknowns)
Causal Structure	Linear, traceable	Multi-level, feedback loops
Reversibility	Course correction possible	Path dependencies, lock-in
Fact-Value Separation	Clear boundaries	Entangled technical-normative
Empirical Grounding	Historical precedents	Unprecedented phenomena
Time Horizons	Years to decades	Months to centuries

2.2.2 Limitations of Traditional Approaches

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.

2.2.3 The Underlying Epistemic Framework

The AMTAIR approach rests on a specific epistemic framework that combines probabilistic reasoning, conditional logic, and possible worlds semantics. This framework provides the philosophical foundation for representing deep uncertainty about AI futures.

Probabilistic Epistemology: Following the Bayesian tradition, we treat probability as a measure of rational credence rather than objective frequency. This subjective interpretation allows meaningful probability assignments even for unique, unprecedented events like AI catastrophe. As E.T. Jaynes demonstrated, probability theory extends deductive logic to handle uncertainty, providing a calculus for rational belief [jaynes2003](#).

Conditional Structure: The framework emphasizes conditional rather than absolute probabilities. Instead of asking “What is $P(\text{catastrophe})$?” we ask “What is $P(\text{catastrophe} \mid \text{specific assumptions})$?” This conditionalization makes explicit the dependency of conclusions on world-view assumptions, enabling productive disagreement about premises rather than conclusions.

Possible Worlds Semantics: We conceptualize uncertainty as distributions over possible worlds—complete descriptions of how reality might unfold. Each world represents a coherent scenario with specific values for all relevant variables. Probability distributions over these worlds capture both what we know and what we don’t know about the future.

This framework enables several key capabilities:

1. **Representing ignorance:** We can express uncertainty about uncertainty itself through hierarchical probability models
2. **Combining evidence:** Bayesian updating provides principled methods for integrating new information
3. **Comparing worldviews:** Different probability distributions over the same space of possibilities enable systematic comparison
4. **Evaluating interventions:** Counterfactual reasoning about how actions change probability distributions

2.2.4 Toward New Epistemic Tools

The inadequacy of traditional methods for AI governance creates an urgent need for new epistemic tools. These tools must:

- **Handle Deep Uncertainty:** Move beyond point estimates to represent ranges of possibilities
- **Capture Complex Causation:** Model multi-level interactions and feedback loops
- **Quantify Disagreement:** Make explicit where experts diverge and why
- **Enable Systematic Analysis:** Support rigorous comparison of policy options

Key Insight: The computational approaches developed in this thesis—particularly Bayesian networks enhanced with automated extraction—directly address each of these requirements by providing formal frameworks for reasoning under uncertainty.

Show Image

Show Image

Show Image

Show Image

Recent work on conditional trees demonstrates the value of structured approaches to uncertainty. McCaslin et al. **mccaslin2024** show how hierarchical conditional forecasting can identify high-value questions for reducing uncertainty about complex topics like AI risk. Their methodology, which asks experts to produce simplified Bayesian networks of informative forecasting questions, achieved nine times higher information value than standard forecasting platform questions.

Tetlock’s work with the Forecasting Research Institute **tetlock2022** exemplifies how prediction markets can provide empirical grounding for formal models. By structuring questions as conditional trees, they enable forecasters to express complex dependencies between events, providing exactly the type of data needed for Bayesian network parameterization.

Gruetzmacher **gruetzmacher2022** evaluates the tradeoffs between full Bayesian networks and conditional trees for forecasting tournaments. While conditional trees offer simplicity, Bayesian networks provide richer representation of dependencies—motivating AMTAIR’s approach of using full networks while leveraging conditional tree insights for question generation.

2.3 Bayesian Networks as Knowledge Representation

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.

2.3.1 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.

Pearl’s foundational work **pearl2014** established Bayesian networks as a principled approach to automated reasoning under uncertainty, providing both theoretical foundations and practical algorithms.

2.3.2 The Rain-Sprinkler-Grass Example

The canonical example illustrates key concepts³:

```
[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

mermaid

```
flowchart TD
    R[Rain<br/>P(rain) = 0.2] --> S[Sprinkler]
    R --> G[Grass_Wet]
    S --> G

    subgraph CPT1 [Sprinkler CPT]
        S1[P(sprinkler|rain) = 0.01]
        S2[P(sprinkler|¬rain) = 0.4]
    end

    subgraph CPT2 [Grass_Wet CPT]
        G1[P(wet|rain,sprinkler) = 0.99]
        G2[P(wet|rain,¬sprinkler) = 0.8]
        G3[P(wet|¬rain,sprinkler) = 0.9]
        G4[P(wet|¬rain,¬sprinkler) = 0.01]
    end
```

python

```
# Basic network representation
nodes = ['Rain', 'Sprinkler', 'Grass_Wet']
edges = [('Rain', 'Sprinkler'), ('Rain', 'Grass_Wet'), ('Sprinkler', 'Grass_Wet')]

# Conditional probability specification
P_wet_given_causes = {
    (True, True): 0.99,    # Rain=T, Sprinkler=T
    (True, False): 0.80,   # Rain=T, Sprinkler=F
    (False, True): 0.90,   # Rain=F, Sprinkler=T
```

³This example, while simple, demonstrates all essential features of Bayesian networks and serves as the foundation for understanding more complex applications

```
(False, False): 0.01    # Rain=F, Sprinkler=F
}
```

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

Show Image

Rain-Sprinkler-Grass Network Rendering

```
#| label: rain_sprinkler_grass_example_network_rendering
#| echo: true
#| eval: true
#| fig-cap: "Dynamic Html Rendering of the Rain-Sprinkler-Grass DAG with Conditional Probabi
#| fig-link: "https://singularitysmith.github.io/AMTAIR_Prototype/bayesian_network.html"
#| fig-alt: "Dynamic Html Rendering of the Rain-Sprinkler-Grass DAG"
```

```
from IPython.display import IFrame
```

```
IFrame(src="https://singularitysmith.github.io/AMTAIR_Prototype/bayesian_network.html", width=
```

2.3.3 Advantages for AI Risk Modeling

These features address key requirements for AI governance:

- **Handling Uncertainty:** Every parameter is a distribution, not a point estimate
- **Representing Causation:** Directed edges embody causal relationships
- **Enabling Analysis:** Formal inference algorithms support systematic evaluation
- **Facilitating Communication:** Visual structure aids cross-domain understanding

2.3.3 Advantages for AI Risk Modeling

Bayesian networks offer several compelling advantages for the peculiar challenge of modeling AI risks—a domain where we’re essentially trying to reason about systems that don’t yet exist, wielding capabilities we can barely imagine, potentially causing outcomes we desperately hope to avoid.

Explicit Uncertainty Representation: Unlike traditional risk assessment tools that often hide uncertainty behind point estimates, Bayesian networks wear their uncertainty on their sleeve. Every node, every edge, every probability is a distribution rather than a false certainty. This matters enormously when discussing AI catastrophe—we’re not pretending to know the unknowable, but rather mapping the landscape of our ignorance with mathematical precision.

Native Causal Reasoning: The directed edges in Bayesian networks aren’t just arrows on a diagram; they encode causal beliefs about how the world works. This enables both forward reasoning (“If we develop AGI, what happens?”) and diagnostic reasoning (“Given that we observe concerning AI behaviors, what does this tell us about underlying alignment?”). Pearl’s do-calculus [pearl2009](#) transforms these networks into laboratories for counterfactual exploration.

Evidence Integration: As new research emerges, as capabilities advance, as governance experiments succeed or fail, Bayesian networks provide a principled framework for updating our beliefs. Unlike static position papers that age poorly, these models can evolve with our understanding—a living document for a rapidly changing field.

Modular Construction: Complex arguments about AI risk involve multiple interacting factors across technical, social, and political domains. Bayesian networks allow us to build these arguments piece by piece, validating each component before assembling the whole. This modularity also enables different experts to contribute their specialized knowledge without needing to understand every aspect of the system.

Visual Communication: Perhaps most importantly for the coordination challenge, Bayesian networks provide a visual language that transcends disciplinary boundaries. A policymaker might not understand the mathematics of instrumental convergence, but they can see how the “power-seeking” node connects to “human disempowerment” in the network diagram. This shared visual vocabulary creates common ground for productive disagreement.

2.4 Argument Mapping and Formal Representations

The journey from a researcher’s intuition about AI risk to a formal probabilistic model resembles translating poetry into mathematics—something essential is always at risk of being lost, yet something equally essential might be gained. Argument mapping provides the crucial middle ground, a structured approach to preserving the logic of natural language arguments while preparing them for mathematical formalization.

2.4.1 From Natural Language to Structure

Natural language arguments about AI risk are rich tapestries woven from causal claims, conditional relationships, uncertainty expressions, and support patterns. When Bostrom writes about the “treacherous turn” **bostrom2014**, he’s not just coining a memorable phrase—he’s encoding a complex causal story about how a seemingly aligned AI system might conceal its true objectives until it gains sufficient power to pursue them without constraint.

The challenge lies in extracting this structure without losing the nuance. Traditional logical analysis might reduce Bostrom’s argument to syllogisms, but this would miss the probabilistic texture, the implicit conditionality, the causal directionality that makes the argument compelling. Argument mapping takes a different approach, seeking to identify:

- **Core claims and propositions:** What exactly is being asserted?
- **Inferential relationships:** How do claims support or challenge each other?
- **Implicit assumptions:** What unstated premises make the argument work?
- **Uncertainty qualifications:** Where does the author express doubt or confidence?

Recent advances in computational argument mining **anderson2007 benn2011 khartabil2021** have shown promise in automating parts of this process. Tools like Microsoft’s Claimify **metropolitansky2025** demonstrate how large language models can extract verifiable claims from complex texts, though the challenge of preserving argumentative structure remains formidable.

2.4.2 ArgDown: Structured Argument Notation

Enter ArgDown **voigt2025**, a markdown-inspired syntax that captures hierarchical argument structure while remaining human-readable. Think of it as the middle child between the wild expressiveness of natural language and the rigid formality of logic—inheriting the best traits of both parents while developing its own personality.

```
[AI_Poses_Risk]: Advanced AI systems may pose existential risk to humanity.
+ [Capability_Growth]: AI capabilities are growing exponentially.
+ [Compute_Scaling]: Available compute doubles every few months.
+ [Algorithmic_Progress]: New architectures show surprising emergent abilities.
+ [Alignment_Difficulty]: Aligning AI with human values is unsolved.
- [Current_Progress]: Some progress on interpretability and oversight.
- [Institutional_Response]: Institutions are mobilizing to address risks.
```

This notation does several clever things simultaneously. The hierarchical structure mirrors how we naturally think about arguments—main claims supported by evidence, which in turn rest on more fundamental observations. The + and – symbols indicate support and opposition relationships, creating a visual flow of argumentative force. Most importantly, it preserves the semantic content of each claim while imposing just enough structure to enable computational processing.

For AMTAIR, we adapt ArgDown specifically for causal arguments, where the hierarchy repre-

sents causal influence rather than logical support. This seemingly small change has profound implications—we’re not just mapping what follows from what, but what causes what.

2.4.3 BayesDown: The Bridge to Bayesian Networks

If ArgDown is the middle child, then BayesDown—developed specifically for this thesis—is the ambitious younger sibling who insists on quantifying everything. By extending ArgDown syntax with probabilistic metadata in JSON format, BayesDown creates a complete specification for Bayesian networks while maintaining human readability.

json

```
[Existential_Catastrophe]: Permanent curtailment of humanity's potential. {
  "instantiations": ["catastrophe_TRUE", "catastrophe_FALSE"],
  "priors": {"p(catastrophe_TRUE)": "0.05", "p(catastrophe_FALSE)": "0.95"},
  "posteriors": {
    "p(catastrophe_TRUE|disempowerment_TRUE)": "0.95",
    "p(catastrophe_TRUE|disempowerment_FALSE)": "0.001"
  }
}

+ [Human_Disempowerment]: Loss of human control over future trajectory. {
  "instantiations": ["disempowerment_TRUE", "disempowerment_FALSE"],
  "priors": {"p(disempowerment_TRUE)": "0.20", "p(disempowerment_FALSE)": "0.80"}
}
```

This representation performs a delicate balancing act. The natural language descriptions preserve the semantic meaning that makes arguments comprehensible. The hierarchical structure maintains the causal relationships that give arguments their logical force. The JSON metadata adds the mathematical precision needed for formal analysis. Together, they create what I call a “hybrid representation”—neither fully natural nor fully formal, but something more useful than either alone.

The two-stage extraction process ($\text{ArgDown} \rightarrow \text{BayesDown}$) mirrors how experts actually think about complex arguments. First, we identify what matters and how things relate causally (structure). Then, we consider how likely different scenarios are based on those relationships (quantification). This separation isn’t just convenient for implementation—it’s psychologically valid.

2.5 The MTAIR Framework: Achievements and Limitations

Understanding AMTAIR requires understanding its intellectual ancestor: the Modeling Transformative AI Risks (MTAIR) project. Like many good ideas in science, MTAIR began with a simple observation and a ambitious goal.

2.5.1 MTAIR’s Approach

The MTAIR project, spearheaded by David Manheim and colleagues [clarke2022](#), emerged from a frustration familiar to anyone who’s attended a conference on AI safety: brilliant people talking past each other, using the same words to mean different things, reaching incompatible conclusions from seemingly shared premises. The diagnosis was elegant—perhaps these disagreements stemmed not from fundamental philosophical differences but from implicit models that had never been made explicit.

Their prescription was equally elegant: manually translate influential AI risk arguments into formal Bayesian networks, making assumptions visible and disagreements quantifiable. Using Analytica software, the team embarked on what can only be described as an intellectual archaeology expedition, carefully excavating the implicit causal models buried in papers, blog posts, and treatises about AI risk.

The process was painstaking:

1. **Systematic Decomposition:** Breaking complex arguments into component claims, identifying variables and relationships through close reading and expert consultation.
2. **Probability Elicitation:** Gathering quantitative estimates through structured expert interviews, literature review, and careful interpretation of qualitative claims.
3. **Sensitivity Analysis:** Testing which parameters most influenced conclusions, revealing where disagreements actually mattered versus where they were merely academic.
4. **Visual Communication:** Creating interactive models that stakeholders could explore, modify, and understand without deep technical training.

The ambition was breathtaking—to create a formal lingua franca for AI risk discussions, enabling productive disagreement and cumulative progress.

2.5.2 Key Achievements

Credit where credit is due: MTAIR demonstrated something many thought impossible. Complex philosophical arguments about AI risk—the kind that sprawl across hundred-page papers mixing technical detail with speculative scenarios—could indeed be formalized without losing their essential insights.

Feasibility of Formalization: The project’s greatest achievement was simply showing it could be done. Arguments from Bostrom, Christiano, and others translated surprisingly well into network form, suggesting that beneath the surface complexity lay coherent causal models waiting to be extracted.

Value of Quantification: Moving from “likely” and “probably” to actual numbers forced precision in a domain often clouded by vague pronouncements. Disagreements that seemed fundamental sometimes evaporated when forced to specify exactly what probability ranges were under dispute.

Cross-Perspective Communication: The formal models created neutral ground where technical AI researchers and policy wonks could meet. Instead of talking past each other in incom-

patible languages, they could point to specific nodes and edges, making disagreements concrete and tractable.

Research Prioritization: Perhaps most practically, sensitivity analysis revealed which empirical questions actually mattered. If changing your belief about technical parameter X from 0.3 to 0.7 doesn’t meaningfully affect the conclusion about AI risk, maybe we should focus our research elsewhere.

2.5.3 Fundamental Limitations

But here’s where the story takes a sobering turn. Despite these achievements, MTAIR faced limitations that prevented it from achieving its full vision—limitations that ultimately motivated the development of AMTAIR.

Labor Intensity: Creating a single model required what can charitably be called a heroic effort. Based on team reports and model complexity, estimates ranged from 200 to 400 expert-hours per formalization⁴. In a field where new influential arguments appear monthly, this pace couldn’t keep up with the discourse.

Static Nature: Once built, these beautiful models began aging immediately. New research emerged, capability assessments shifted, governance proposals evolved—but updating the models required near-complete reconstruction. They were snapshots of arguments at particular moments, not living representations that could evolve.

Limited Accessibility: Using the models required Analytica software and non-trivial technical sophistication. The very experts whose arguments were being formalized often couldn’t directly engage with their formalized representations without intermediation.

Single Perspective: Each model represented one worldview at a time. Comparing different perspectives required building entirely separate models, making systematic comparison across viewpoints labor-intensive and error-prone.

These weren’t failures of execution but fundamental constraints of the manual approach. Like medieval scribes copying manuscripts, the MTAIR team had shown the value of preservation and dissemination, but the printing press had yet to be invented.

2.5.4 The Automation Opportunity

The MTAIR experience revealed a tantalizing possibility: if the bottleneck was human labor rather than conceptual feasibility, perhaps automation could crack open the problem. The rise of large language models capable of sophisticated reasoning about text created a technological moment ripe for exploitation.

Key lessons from MTAIR informed the automation approach:

- Formal models genuinely enhance understanding and coordination—the juice is worth the squeeze

⁴These estimates include time for initial extraction, expert consultation, probability elicitation, validation, and refinement

- The modeling process itself surfaces implicit assumptions—extraction is as valuable as the final product
- Quantification enables analyses impossible with qualitative arguments alone—numbers matter even when uncertain
- But manual approaches cannot scale to match the challenge—we need computational leverage

This set the stage for AMTAIR’s central innovation: using frontier language models to automate the extraction and formalization process while preserving the benefits MTAIR had demonstrated. Not to replace human judgment, but to amplify it—turning what took weeks into what takes hours, enabling comprehensive coverage rather than selective sampling.

2.6 Literature Review: Content and Technical Levels

The intellectual landscape surrounding AI risk resembles a rapidly expanding metropolis—new neighborhoods of thought spring up monthly, connected by bridges of varying stability to the established districts. A comprehensive review would fill volumes, so let me provide a guided tour of the territories most relevant to AMTAIR’s mission.

2.6.1 AI Risk Models Evolution

The evolution of AI risk models traces a path from philosophical speculation to increasingly rigorous formalization—a journey from “what if?” to “how likely?”

Early Phase (2000-2010): The conversation began with broad conceptual arguments. Good’s ultraintelligent machine **good1966** and Vinge’s technological singularity set the stage, but these were more thought experiments than models. Yudkowsky’s early writings **yudkowsky2008** introduced key concepts like recursive self-improvement and orthogonality but remained largely qualitative.

Formalization Phase (2010-2018): Bostrom’s *Superintelligence* **bostrom2014** marked a watershed, providing systematic analysis of pathways, capabilities, and risks. The book’s genius lay not in mathematical formalism but in conceptual clarity—decomposing the nebulous fear of “robot overlords” into specific mechanisms like instrumental convergence and infrastructure profusion.

Quantification Phase (2018-present): Recent years have seen explicit probability estimates entering mainstream discourse. Carlsmith’s power-seeking model **carlsmith2022**, Cotra’s biological anchors, and various compute-based timelines represent attempts to put numbers on previously qualitative claims. The field increasingly recognizes that governance decisions require more than philosophical arguments—they need probability distributions.

This progression reflects a maturing field, though it also creates new challenges. As models become more quantitative, they risk false precision. As they become more complex, they risk inscrutability. AMTAIR attempts to navigate these tensions by preserving the narrative clarity of earlier work while enabling the mathematical rigor of recent approaches.

2.6.2 Governance Proposals Taxonomy

If risk models are the diagnosis, governance proposals are the treatment plans—and like medicine, they range from gentle interventions to radical surgery.

Technical Standards: The “first, do no harm” approach focuses on concrete safety requirements—interpretability benchmarks, robustness testing, capability thresholds. These proposals, exemplified by standard-setting bodies and technical safety organizations, offer specificity at the cost of narrowness.

Regulatory Frameworks: Moving up the intervention ladder, we find comprehensive regulatory proposals like the EU AI Act **europaen2024**. These create institutional structures, liability regimes, and oversight mechanisms, trading broad coverage for implementation complexity.

International Coordination: At the ambitious end, proposals for international AI governance treaties, soft law arrangements, and technical cooperation agreements aim to prevent races to the bottom. Think nuclear non-proliferation but for minds instead of missiles.

Research Priorities: Cutting across these categories, work by Dafoe **dafoe2018** and others maps the research landscape itself—what questions need answering before we can govern wisely? This meta-level analysis shapes funding flows and talent allocation.

A particularly compelling example of conditional governance thinking comes from “A Narrow Path” **miotti2024**, which proposes a phased approach: immediate safety measures to prevent uncontrolled development, international institutions to ensure stability, and long-term scientific foundations for beneficial transformative AI. This temporal sequencing—safety, stability, then flourishing—reflects growing sophistication in governance thinking.

2.6.3 Bayesian Network Theory and Applications

The mathematical machinery underlying AMTAIR rests on decades of theoretical development in probabilistic graphical models. Understanding this foundation helps appreciate both the power and limitations of the approach.

The key insight, crystallized in the work of Pearl **pearl2014** and elaborated by Koller & Friedman **koller2009**, is that independence relationships in complex systems can be read from graph structure. D-separation, the Markov condition, and the relationship between graphs and probability distributions provide the mathematical spine that makes Bayesian networks more than pretty pictures.

Critical concepts for AI risk modeling:

- **Conditional Independence:** Variable A is independent of C given B—encoded through graph separation
- **Markov Condition:** Each variable is independent of its non-descendants given its parents
- **Inference Algorithms:** From exact variable elimination to approximate Monte Carlo methods

- **Causal Interpretation:** When edges represent causal influence, the network supports counterfactual reasoning

These aren’t just mathematical niceties. When we claim that “deployment decisions” mediates the relationship between “capability advancement” and “catastrophic risk,” we’re making a precise statement about conditional independence that has testable implications.

2.6.4 Software Tools Landscape

The gap between Bayesian network theory and practical implementation is bridged by an ecosystem of software tools, each with its own strengths and opinions about how probabilistic reasoning should work.

pgmpy: This Python library provides the computational backbone for AMTAIR, offering both learning algorithms and inference engines. Its object-oriented design maps naturally onto our extraction pipeline.

NetworkX: For graph manipulation and analysis, NetworkX has become the de facto standard in Python, providing algorithms for everything from centrality measurement to community detection.

PyVis: Interactive visualization transforms static networks into explorable landscapes. PyVis’s integration with web technologies enables the rich interactive features that make formal models accessible.

Pandas/NumPy: The workhorses of scientific Python handle data manipulation and numerical computation, providing the infrastructure on which everything else builds.

The integration challenge—making these tools play nicely together while maintaining performance and correctness—shaped many architectural decisions in AMTAIR. Each tool excels in its domain, but the seams between them required careful engineering.

2.6.5 Formalization Approaches

The challenge of formalizing natural language arguments extends far beyond AI risk, touching on fundamental questions in logic, linguistics, and artificial intelligence.

Pollock’s work on cognitive carpentry **pollock1995** provides philosophical grounding, arguing that human reasoning itself involves implicit formal structures that can be computationally modeled. This view—that formalization reveals rather than imposes structure—underlies AMTAIR’s approach.

Key theoretical challenges:

- **Semantic Preservation:** How do we maintain meaning while adding precision?
- **Structural Extraction:** What implicit relationships lurk in natural language?
- **Uncertainty Quantification:** How do we map “likely” to numbers?

Recent work on causal structure learning from text **babakov2025** **ban2023** **bethard2007** offers hope that these challenges can be addressed computationally. The convergence of large lan-

guage models with formal methods creates new possibilities for bridging the semantic-symbolic gap.

2.6.6 Correlation Accounting Methods

One of the most persistent criticisms of Bayesian networks concerns their assumption of conditional independence given parents. In the real world, and especially in complex socio-technical systems like AI development, correlations abound.

Methods for handling these correlations have evolved considerably:

Copula Methods: By separating marginal distributions from dependence structure, copulas [nelson2006](#) allow modeling of complex correlations while preserving the Bayesian network framework. Think of it as adding a correlation layer on top of the basic network.

Hierarchical Models: Introducing latent variables that influence multiple observed variables captures correlations naturally. If “AI research culture” influences both “capability progress” and “safety investment,” their correlation is explained.

Explicit Correlation Nodes: Sometimes the most straightforward approach is best—directly model correlation mechanisms as additional nodes in the network.

Sensitivity Bounds: When correlations remain uncertain, compute best and worst case scenarios. This reveals when independence assumptions critically affect conclusions versus when they’re harmless simplifications.

For AMTAIR, the pragmatic approach dominates: start with independence assumptions, identify where they matter through sensitivity analysis, then selectively add correlation modeling where it most affects conclusions.

2.7 Methodology

The methodology of this research resembles less a linear march from hypothesis to conclusion and more an iterative dance between theory and implementation, vision and reality. Let me walk you through the choreography.

2.7.1 Research Design Overview

This research follows what methodologists might call a “design science” approach—we’re not just studying existing phenomena but creating new artifacts (the AMTAIR system) and evaluating their utility for solving practical problems (the coordination crisis in AI governance).

The overall flow:

1. **Theoretical Development:** Establishing why automated extraction could address the coordination crisis, grounded in epistemic theory and mechanism design
2. **Technical Implementation:** Building working software that demonstrates feasibility, not as a proof-of-concept toy but as a system capable of handling real arguments

3. **Empirical Validation:** Testing extraction quality against expert judgment, measuring not just accuracy but usefulness for downstream tasks
4. **Application Studies:** Applying the system to real AI governance questions, evaluating whether formal models actually enhance decision-making

This isn't waterfall development where each phase completes before the next begins. Rather, insights from implementation fed back into theory, validation results shaped technical improvements, and application attempts revealed new requirements. The methodology itself embodied the iterative refinement it sought to enable.

2.7.2 Formalizing World Models from AI Safety Literature

The core methodological challenge—transforming natural language arguments into formal probabilistic models—requires careful consideration of what we're actually trying to capture.

A “world model” in this context isn't just any formal representation but specifically a causal model embodying beliefs about how different factors influence AI risk. The extraction approach must therefore:

- **Identify key variables:** Not just any entities mentioned, but causally relevant factors
- **Extract causal relationships:** Not mere correlation or co-occurrence, but directed influence
- **Capture uncertainty:** Both structural uncertainty (does A cause B?) and parametric uncertainty (how strongly?)
- **Preserve context:** Maintaining enough semantic information to interpret the formal model

Large language models enable this through sophisticated pattern recognition and reasoning capabilities, but they're tools, not magic wands. The methodology must account for their strengths (recognizing implicit structure) and weaknesses (potential hallucination, inconsistency).

2.7.3 From Natural Language to Computational Models

The journey from text to computation follows a carefully designed pipeline that mirrors human cognitive processes. Just as you wouldn't ask someone to simultaneously parse grammar and solve equations, we separate structural understanding from quantitative reasoning.

The Two-Stage Process:

Stage 1 focuses on structure—what causes what? The LLM reads an argument much as a human would, identifying key claims and their relationships. The prompt design here is crucial, providing enough guidance to ensure consistent extraction while allowing flexibility for different argument styles.

Stage 2 adds quantities—how likely is each outcome? With structure established, the system generates targeted questions about probabilities. This separation enables different approaches to quantification: extracting explicit estimates from text, inferring from qualitative language, or even connecting to external prediction markets.

The magic happens in the interplay. Structure constrains what probabilities are needed. Probability requirements might reveal missing structural elements. The process is a dialogue between qualitative and quantitative understanding.

2.7.4 Directed Acyclic Graphs: Structure and Semantics

At the mathematical heart of Bayesian networks lie Directed Acyclic Graphs (DAGs)—structures that are simultaneously simple enough to analyze and rich enough to capture complex phenomena.

The “directed” part encodes causality or influence—edges have direction, flowing from cause to effect. The “acyclic” part ensures logical coherence—you can’t have A causing B causing C causing A, no matter how much certain political arguments might suggest otherwise.

Key properties for AI risk modeling:

Acyclicity: More than a mathematical convenience, this enforces coherent temporal or causal ordering. In AI risk arguments, this prevents circular reasoning where consequences justify premises that predict those same consequences.

D-separation: This graphical criterion determines conditional independence. If knowing about AI capabilities tells you nothing additional about risk given that you know deployment decisions, then capabilities and risk are d-separated given deployment.

Markov Condition: Each variable depends only on its parents, not on its entire ancestry. This locality assumption makes inference tractable and forces modelers to make intervention points explicit.

Path Analysis: Following paths through the graph reveals how influence propagates. Multiple paths between variables indicate redundancy—important for understanding intervention robustness.

The causal interpretation, following Pearl’s framework, transforms these mathematical objects into tools for counterfactual reasoning. When we ask “what if we prevented deployment of misaligned systems?” we’re performing surgery on the DAG, setting variables and propagating consequences.

2.7.5 Quantification of Probabilistic Judgments

Here we encounter one of the most philosophically fraught aspects of the methodology: turning words into numbers. When an expert writes “highly likely,” what probability should we assign? When they say “significant risk,” what distribution captures their belief?

The methodology embraces rather than elides this challenge:

Calibration Studies: Research on human probability expression shows systematic patterns. “Highly likely” typically maps to 0.8-0.9, “probable” to 0.6-0.8, though individual and cultural variation is substantial.

Extraction Strategies: The system uses multiple approaches:

- Direct extraction: “We estimate 65% probability”
- Linguistic mapping: “Very likely” \rightarrow 0.85 (with uncertainty)
- Comparative extraction: “More likely than X” where $P(X)$ is known
- Bounded extraction: “At least 30%” \rightarrow [0.30, 1.0]

Uncertainty Representation: Rather than false precision, we maintain uncertainty about probabilities themselves. This might seem like uncertainty piled on uncertainty, but it’s honest—and mathematically tractable through hierarchical models.

The goal isn’t perfect extraction but useful extraction. If we can narrow “significant risk” from $[0, 1]$ to $[0.15, 0.45]$, we’ve added information even if we haven’t achieved precision.

2.7.6 Inference Techniques for Complex Networks

Once we’ve built these formal models, we need to reason with them—and here computational complexity rears its exponential head. The number of probability calculations required for exact inference grows exponentially with network connectivity, quickly overwhelming even modern computers.

The methodology employs a portfolio of approaches:

Exact Methods: For smaller networks (<30 nodes), variable elimination and junction tree algorithms provide exact answers. These form the gold standard against which we validate approximate methods.

Sampling Approaches: Monte Carlo methods trade exactness for scalability. By simulating many possible worlds consistent with our probability model, we approximate the true distributions. The law of large numbers is our friend here.

Variational Methods: These turn inference into optimization—find the simplest distribution that approximates our true beliefs. Like finding the best polynomial approximation to a complex curve.

Hybrid Strategies: Different parts of the network might use different methods. Exact inference for critical subgraphs, approximation for peripheral components.

The choice of method affects not just computation time but the types of questions we can meaningfully ask. This creates a methodological feedback loop where feasible inference shapes model design.

2.7.7 Integration with Prediction Markets and Forecasting Platforms

While full integration remains future work, the methodology anticipates connection to live forecasting data as a critical enhancement. The vision is compelling: formal models grounded in collective intelligence, updating as new information emerges.

The planned approach would involve:

Semantic Matching: Model variables rarely align perfectly with forecast questions. “AI causes human extinction” might map to multiple specific forecasts about capabilities, deployment, and

impacts. Developing robust matching algorithms is essential.

Temporal Alignment: Markets predict specific dates (“AGI by 2030”) while models consider scenarios (“given AGI development”). Bridging these requires careful probability conditioning.

Quality Weighting: Not all forecasts are created equal. Platform reputation, forecaster track records, and market depth all affect reliability. The methodology must account for this heterogeneity.

Update Scheduling: Real-time updates would overwhelm users and computation. The system needs intelligent policies about when model updates provide value.

Platforms like Metaculus **tetlock2022** already demonstrate sophisticated conditional forecasting on AI topics. The challenge lies not in data availability but in meaningful integration that enhances rather than complicates decision-making.

With these theoretical foundations and methodological commitments established, we can now turn to the concrete implementation of AMTAIR. The next chapter demonstrates how these abstract principles translate into working software that addresses real governance challenges. The journey from theory to practice always involves surprises—some pleasant, others less so—but that’s what makes it interesting.

3. AMTAIR: Design and Implementation

Chapter Overview

Grade Weight: 20% | **Target Length:** ~29% of text (~8,700 words)

Requirements: Critical evaluation, strong argument for position, original contribution

The moment of truth in any research project comes when elegant theories meet stubborn reality. For AMTAIR, this meant transforming the vision of automated argument extraction into working code that could handle the beautiful messiness of real AI safety arguments. Let me take you through this journey from blueprint to implementation, complete with victories, defeats, and the occasional moment of “well, that’s unexpected.”

3.1 System Architecture Overview

Picture, if you will, a factory for transforming arguments into models. Raw materials (PDFs, blog posts, research papers) enter at one end. Finished products (interactive Bayesian networks) emerge at the other. In between lies a carefully orchestrated pipeline where each stage performs its specialized transformation, passing refined materials to the next.

The AMTAIR architecture embodies a philosophy: complex tasks become manageable when decomposed into focused components. Rather than building a monolithic “argument-to-model” black box, we created a series of specialized modules, each excellent at one thing.

The pipeline consists of five main stages:

1. **Text Ingestion and Preprocessing:** Like a careful librarian, this stage catalogues incoming documents, normalizes their format, extracts metadata, and identifies the argumentative content worth processing.
2. **Argument Extraction:** The intellectual heart of the system, where large language models perform their magic, transforming prose into structured representations.
3. **Data Transformation:** The workshop where extracted arguments are refined, validated, and prepared for mathematical representation.
4. **Network Construction:** The assembly line where formal Bayesian networks are instantiated, complete with conditional probability tables.

5. **Interactive Visualization:** The showroom where complex models become accessible through thoughtful design and interactivity.

3.1.1 Five-Stage Pipeline Architecture

Let’s examine each stage more closely, understanding not just what they do but why they exist as separate components.

Text Ingestion and Preprocessing handles the unglamorous but essential work of standardization. Academic PDFs, with their two-column layouts and embedded figures, differ vastly from blog posts with inline code and hyperlinks. This stage creates a uniform representation while preserving essential structure and metadata. Format normalization strips away presentation while preserving content. Metadata extraction captures authorship, publication date, and citations. Relevance filtering identifies sections containing arguments rather than literature reviews or acknowledgments. Character encoding standardization prevents those maddening replacement characters that plague text processing.

Argument Extraction represents AMTAIR’s core innovation. Using a two-stage process that mirrors human reasoning, it first identifies structural relationships (what influences what) then quantifies those relationships (how likely, how strong). This separation enables targeted prompts optimized for each task, human verification between stages, and modular improvements as LLM capabilities evolve.

Data Transformation bridges the gap between textual representations and mathematical models. It parses the BayesDown syntax into structured data, validates that the resulting network forms a proper DAG, checks probability consistency, and handles missing data intelligently.

Network Construction instantiates the formal mathematical model. This involves creating nodes and edges according to extracted structure, populating conditional probability tables, initializing inference engines, and validating the complete model.

Interactive Visualization makes the complex accessible. Through thoughtful visual encoding of probabilities and relationships, progressive disclosure of detail, interactive exploration capabilities, and multiple export formats, it serves diverse stakeholder needs.

3.1.2 Design Principles

Core Design Philosophy: The architecture embodies several principles that guided countless implementation decisions:

Modularity: Each component has clear inputs, outputs, and responsibilities. This isn’t just good software engineering—it enables independent improvement of components and graceful degradation when parts fail.

Validation Checkpoints: Between each stage, we validate outputs before proceeding. Bad extractions don’t propagate into visualization. Malformed networks trigger re-extraction rather than cryptic errors.

Human-in-the-Loop: While pursuing automation, we recognize that human judgment remains invaluable. The architecture provides natural intervention points where experts can verify and correct.

Extensibility: New document formats, improved extraction prompts, alternative visualization libraries—the architecture accommodates growth without restructuring.

The system emphasizes transparency over black-box efficiency. Users can inspect intermediate representations, understand extraction decisions, and verify transformations. This builds trust—essential for a system handling high-stakes arguments about existential risk.

3.2 The Two-Stage Extraction Process

The heart of AMTAIR beats with a two-stage rhythm: structure, then probability. This separation, which initially seemed like an implementation detail, revealed itself as fundamental to the extraction challenge.

3.2.1 Stage 1: Structural Extraction (ArgDown)

Imagine reading a complex argument about AI risk. Your first pass likely isn’t calculating exact probabilities—you’re mapping the landscape. What are the key claims? How do they relate? What supports what? Stage 1 mirrors this cognitive process.

The extraction begins with pattern recognition. Natural language contains linguistic markers of causal relationships: “leads to,” “results in,” “depends on,” “influences.” The LLM, trained on vast corpora of argumentative text, recognizes these patterns and their variations.

Consider extracting from a passage like: “The development of artificial general intelligence will likely lead to rapid capability gains through recursive self-improvement. This intelligence explosion could result in systems pursuing convergent instrumental goals, potentially including resource acquisition and self-preservation. Without solved alignment, such power-seeking behavior poses existential risks to humanity.”

The system identifies three key variables connected by causal relationships:

- AGI Development → Intelligence Explosion
- Intelligence Explosion → Power-Seeking Behavior
- Power-Seeking Behavior → Existential Risk

But extraction goes beyond simple pattern matching. The system must handle complex linguistic phenomena like coreference (“this,” “such systems”), implicit relationships, conditional statements, and negative statements. The magic lies in prompt engineering that guides the LLM to consistent extraction while remaining flexible enough for diverse argument styles.

The output, formatted in ArgDown syntax, preserves both structure and semantics:

```
[Existential_Risk]: Threat to humanity's continued existence and flourishing.
```

```
+ [Power_Seeking_Behavior]: AI systems pursuing instrumental goals like resource acquisition
```

- + [Intelligence_Explosion]: Rapid recursive self-improvement leading to superintelligence
- + [AGI_Development]: Creation of artificial general intelligence systems.

3.2.2 Stage 2: Probability Integration (BayesDown)

With structure established, Stage 2 adds the quantitative flesh to the qualitative bones. This stage faces a different challenge: extracting numerical beliefs from text that often expresses uncertainty in frustratingly vague terms.

The process begins by generating targeted questions based on the extracted structure. For each node, we need prior probabilities. For each child-parent relationship, we need conditional probabilities. The combinatorics can be daunting—a node with three binary parents requires 8 conditional probability values.

The system employs multiple strategies for probability extraction:

Explicit Extraction: When authors provide numerical estimates (“we assign 70% probability”), extraction is straightforward, though we must handle various formats and contexts.

Linguistic Mapping: Qualitative expressions map to probability ranges based on calibration studies. “Highly likely” becomes approximately 0.85, though we maintain uncertainty about this mapping.

Comparative Reasoning: Statements like “more probable than not” or “at least as likely as X” provide bounds even without exact values.

Coherence Enforcement: Probabilities must sum correctly. If $P(A|B) = 0.7$, then $P(\text{not } A|B)$ must equal 0.3. The system detects and resolves inconsistencies.

The result is a complete BayesDown specification:

json

```
[Existential_Risk]: Threat to humanity's continued existence. {
  "instantiations": ["true", "false"],
  "priors": {"p(true)": "0.10", "p(false)": "0.90"},
  "posteriors": {
    "p(true|power_seeking_true)": "0.65",
    "p(true|power_seeking_false)": "0.001"
  }
}
```

3.2.3 Why Two Stages?

The separation of structure from probability isn’t merely convenient—it’s cognitively valid and practically essential. Let me count the ways this design decision pays dividends:

Cognitive Alignment: Humans naturally separate “what relates to what” from “how likely is it.” The two-stage process mirrors this, making the system’s operation intuitive and interpretable.

Error Isolation: Structural errors (missing a key variable) differ fundamentally from probability errors (estimating 0.7 instead of 0.8). Separating stages allows targeted debugging and improvement.

Modular Validation: Experts can verify structure without needing to evaluate every probability. This enables efficient human oversight at natural checkpoints.

Flexible Quantification: Different probability sources (text extraction, expert elicitation, market data) can feed into the same structure. The architecture accommodates multiple approaches to the probability challenge.

Transparency: Users can inspect ArgDown to understand what was extracted before probabilities were added. This builds trust and enables meaningful correction.

The two-stage approach also revealed an unexpected benefit: ArgDown itself became a valuable output. Researchers began using these structural extractions for qualitative analysis, even without probability quantification. Sometimes, just making argument structure explicit provides sufficient value.

3.3 Implementation Technologies

Choosing technologies for AMTAIR resembled assembling a band—each instrument needed to excel individually while harmonizing with the ensemble. The selection criteria balanced capability, maturity, interoperability, and community support.

3.3.1 Technology Stack

The final ensemble performs beautifully:

Component	Technology	Purpose	Why This Choice
Language Models	GPT-4, Claude	Argument extraction	State-of-the-art reasoning capabilities
Network Analysis	NetworkX	Graph algorithms	Mature, comprehensive, well-documented
Probabilistic Modeling	pgmpy	Bayesian operations	Native Python, active development
Visualization	PyVis	Interactive rendering	Web-based, customizable, responsive
Data Processing	Pandas	Structured manipulation	Industry standard, powerful operations

Language Models form the cognitive core. GPT-4 and Claude demonstrate remarkable ability to understand complex arguments, recognize implicit structure, and maintain coherence across

long extractions. The choice to support multiple models provides robustness and allows leveraging their complementary strengths.

NetworkX handles all graph-theoretic heavy lifting. From basic operations like cycle detection to advanced algorithms like centrality measurement, it provides a comprehensive toolkit that would take years to replicate.

pgmpy bridges the gap between graph structure and probabilistic reasoning. Its clean API design maps naturally onto our extracted representations, while its inference algorithms handle the computational complexity of Bayesian reasoning.

PyVis transforms static networks into living documents. Built on vis.js, it provides smooth physics simulations, rich interactivity, and extensive customization options—all accessible through Python.

Pandas might seem mundane compared to its companions, but it’s the reliable rhythm section that keeps everything together. Its ability to reshape, merge, and transform structured data makes the complex data transformations tractable.

3.3.2 Key Algorithms

2.3.3 Advantages for AI Risk Modeling

Bayesian networks offer several compelling advantages for the peculiar challenge of modeling AI risks—a domain where we’re essentially trying to reason about systems that don’t yet exist, wielding capabilities we can barely imagine, potentially causing outcomes we desperately hope to avoid.

Explicit Uncertainty Representation: Unlike traditional risk assessment tools that often hide uncertainty behind point estimates, Bayesian networks wear their uncertainty on their sleeve. Every node, every edge, every probability is a distribution rather than a false certainty. This matters enormously when discussing AI catastrophe—we’re not pretending to know the unknowable, but rather mapping the landscape of our ignorance with mathematical precision.

Native Causal Reasoning: The directed edges in Bayesian networks aren’t just arrows on a diagram; they encode causal beliefs about how the world works. This enables both forward reasoning (“If we develop AGI, what happens?”) and diagnostic reasoning (“Given that we observe concerning AI behaviors, what does this tell us about underlying alignment?”). Pearl’s do-calculus **pearl2009** transforms these networks into laboratories for counterfactual exploration.

Evidence Integration: As new research emerges, as capabilities advance, as governance experiments succeed or fail, Bayesian networks provide a principled framework for updating our beliefs. Unlike static position papers that age poorly, these models can evolve with our understanding—a living document for a rapidly changing field.

Modular Construction: Complex arguments about AI risk involve multiple interacting factors across technical, social, and political domains. Bayesian networks allow us to build these arguments piece by piece, validating each component before assembling the whole. This mod-

ularity also enables different experts to contribute their specialized knowledge without needing to understand every aspect of the system.

Visual Communication: Perhaps most importantly for the coordination challenge, Bayesian networks provide a visual language that transcends disciplinary boundaries. A policymaker might not understand the mathematics of instrumental convergence, but they can see how the “power-seeking” node connects to “human disempowerment” in the network diagram. This shared visual vocabulary creates common ground for productive disagreement.

2.4 Argument Mapping and Formal Representations

The journey from a researcher’s intuition about AI risk to a formal probabilistic model resembles translating poetry into mathematics—something essential is always at risk of being lost, yet something equally essential might be gained. Argument mapping provides the crucial middle ground, a structured approach to preserving the logic of natural language arguments while preparing them for mathematical formalization.

2.4.1 From Natural Language to Structure

Natural language arguments about AI risk are rich tapestries woven from causal claims, conditional relationships, uncertainty expressions, and support patterns. When Bostrom writes about the “treacherous turn” **bostrom2014**, he’s not just coining a memorable phrase—he’s encoding a complex causal story about how a seemingly aligned AI system might conceal its true objectives until it gains sufficient power to pursue them without constraint.

The challenge lies in extracting this structure without losing the nuance. Traditional logical analysis might reduce Bostrom’s argument to syllogisms, but this would miss the probabilistic texture, the implicit conditionality, the causal directionality that makes the argument compelling. Argument mapping takes a different approach, seeking to identify:

- **Core claims and propositions:** What exactly is being asserted?
- **Inferential relationships:** How do claims support or challenge each other?
- **Implicit assumptions:** What unstated premises make the argument work?
- **Uncertainty qualifications:** Where does the author express doubt or confidence?

Recent advances in computational argument mining **anderson2007 benn2011 khartabil2021** have shown promise in automating parts of this process. Tools like Microsoft’s Claimify **metropolitansky2025** demonstrate how large language models can extract verifiable claims from complex texts, though the challenge of preserving argumentative structure remains formidable.

2.4.2 ArgDown: Structured Argument Notation

Enter ArgDown **voigt2025**, a markdown-inspired syntax that captures hierarchical argument structure while remaining human-readable. Think of it as the middle child between the wild

expressiveness of natural language and the rigid formality of logic—inheriting the best traits of both parents while developing its own personality.

```
[AI_Poses_Risk]: Advanced AI systems may pose existential risk to humanity.
+ [Capability_Growth]: AI capabilities are growing exponentially.
+ [Compute_Scaling]: Available compute doubles every few months.
+ [Algorithmic_Progress]: New architectures show surprising emergent abilities.
+ [Alignment_Difficulty]: Aligning AI with human values is unsolved.
- [Current_Progress]: Some progress on interpretability and oversight.
- [Institutional_Response]: Institutions are mobilizing to address risks.
```

This notation does several clever things simultaneously. The hierarchical structure mirrors how we naturally think about arguments—main claims supported by evidence, which in turn rest on more fundamental observations. The + and - symbols indicate support and opposition relationships, creating a visual flow of argumentative force. Most importantly, it preserves the semantic content of each claim while imposing just enough structure to enable computational processing.

For AMTAIR, we adapt ArgDown specifically for causal arguments, where the hierarchy represents causal influence rather than logical support. This seemingly small change has profound implications—we’re not just mapping what follows from what, but what causes what.

2.4.3 BayesDown: The Bridge to Bayesian Networks

If ArgDown is the middle child, then BayesDown—developed specifically for this thesis—is the ambitious younger sibling who insists on quantifying everything. By extending ArgDown syntax with probabilistic metadata in JSON format, BayesDown creates a complete specification for Bayesian networks while maintaining human readability.

json

```
[Existential_Catastrophe]: Permanent curtailment of humanity's potential. {
  "instantiations": ["catastrophe_TRUE", "catastrophe_FALSE"],
  "priors": {"p(catastrophe_TRUE)": "0.05", "p(catastrophe_FALSE)": "0.95"},
  "posteriors": {
    "p(catastrophe_TRUE|disempowerment_TRUE)": "0.95",
    "p(catastrophe_TRUE|disempowerment_FALSE)": "0.001"
  }
}

+ [Human_Disempowerment]: Loss of human control over future trajectory. {
  "instantiations": ["disempowerment_TRUE", "disempowerment_FALSE"],
  "priors": {"p(disempowerment_TRUE)": "0.20", "p(disempowerment_FALSE)": "0.80"}
}
```

This representation performs a delicate balancing act. The natural language descriptions preserve the semantic meaning that makes arguments comprehensible. The hierarchical structure maintains the causal relationships that give arguments their logical force. The JSON metadata

adds the mathematical precision needed for formal analysis. Together, they create what I call a “hybrid representation”—neither fully natural nor fully formal, but something more useful than either alone.

The two-stage extraction process ($\text{ArgDown} \rightarrow \text{BayesDown}$) mirrors how experts actually think about complex arguments. First, we identify what matters and how things relate causally (structure). Then, we consider how likely different scenarios are based on those relationships (quantification). This separation isn’t just convenient for implementation—it’s psychologically valid.

2.5 The MTAIR Framework: Achievements and Limitations

Understanding AMTAIR requires understanding its intellectual ancestor: the Modeling Transformative AI Risks (MTAIR) project. Like many good ideas in science, MTAIR began with a simple observation and a ambitious goal.

2.5.1 MTAIR’s Approach

The MTAIR project, spearheaded by David Manheim and colleagues [clarke2022](#), emerged from a frustration familiar to anyone who’s attended a conference on AI safety: brilliant people talking past each other, using the same words to mean different things, reaching incompatible conclusions from seemingly shared premises. The diagnosis was elegant—perhaps these disagreements stemmed not from fundamental philosophical differences but from implicit models that had never been made explicit.

Their prescription was equally elegant: manually translate influential AI risk arguments into formal Bayesian networks, making assumptions visible and disagreements quantifiable. Using Analytica software, the team embarked on what can only be described as an intellectual archaeology expedition, carefully excavating the implicit causal models buried in papers, blog posts, and treatises about AI risk.

The process was painstaking:

1. **Systematic Decomposition:** Breaking complex arguments into component claims, identifying variables and relationships through close reading and expert consultation.
2. **Probability Elicitation:** Gathering quantitative estimates through structured expert interviews, literature review, and careful interpretation of qualitative claims.
3. **Sensitivity Analysis:** Testing which parameters most influenced conclusions, revealing where disagreements actually mattered versus where they were merely academic.
4. **Visual Communication:** Creating interactive models that stakeholders could explore, modify, and understand without deep technical training.

The ambition was breathtaking—to create a formal lingua franca for AI risk discussions, enabling productive disagreement and cumulative progress.

2.5.2 Key Achievements

Credit where credit is due: MTAIR demonstrated something many thought impossible. Complex philosophical arguments about AI risk—the kind that sprawl across hundred-page papers mixing technical detail with speculative scenarios—could indeed be formalized without losing their essential insights.

Feasibility of Formalization: The project’s greatest achievement was simply showing it could be done. Arguments from Bostrom, Christiano, and others translated surprisingly well into network form, suggesting that beneath the surface complexity lay coherent causal models waiting to be extracted.

Value of Quantification: Moving from “likely” and “probably” to actual numbers forced precision in a domain often clouded by vague pronouncements. Disagreements that seemed fundamental sometimes evaporated when forced to specify exactly what probability ranges were under dispute.

Cross-Perspective Communication: The formal models created neutral ground where technical AI researchers and policy wonks could meet. Instead of talking past each other in incompatible languages, they could point to specific nodes and edges, making disagreements concrete and tractable.

Research Prioritization: Perhaps most practically, sensitivity analysis revealed which empirical questions actually mattered. If changing your belief about technical parameter X from 0.3 to 0.7 doesn’t meaningfully affect the conclusion about AI risk, maybe we should focus our research elsewhere.

2.5.3 Fundamental Limitations

But here’s where the story takes a sobering turn. Despite these achievements, MTAIR faced limitations that prevented it from achieving its full vision—limitations that ultimately motivated the development of AMTAIR.

Labor Intensity: Creating a single model required what can charitably be called a heroic effort. Based on team reports and model complexity, estimates ranged from 200 to 400 expert-hours per formalization⁵. In a field where new influential arguments appear monthly, this pace couldn’t keep up with the discourse.

Static Nature: Once built, these beautiful models began aging immediately. New research emerged, capability assessments shifted, governance proposals evolved—but updating the models required near-complete reconstruction. They were snapshots of arguments at particular moments, not living representations that could evolve.

Limited Accessibility: Using the models required Analytica software and non-trivial technical sophistication. The very experts whose arguments were being formalized often couldn’t directly engage with their formalized representations without intermediation.

⁵These estimates include time for initial extraction, expert consultation, probability elicitation, validation, and refinement

Single Perspective: Each model represented one worldview at a time. Comparing different perspectives required building entirely separate models, making systematic comparison across viewpoints labor-intensive and error-prone.

These weren't failures of execution but fundamental constraints of the manual approach. Like medieval scribes copying manuscripts, the MTAIR team had shown the value of preservation and dissemination, but the printing press had yet to be invented.

2.5.4 The Automation Opportunity

The MTAIR experience revealed a tantalizing possibility: if the bottleneck was human labor rather than conceptual feasibility, perhaps automation could crack open the problem. The rise of large language models capable of sophisticated reasoning about text created a technological moment ripe for exploitation.

Key lessons from MTAIR informed the automation approach:

- Formal models genuinely enhance understanding and coordination—the juice is worth the squeeze
- The modeling process itself surfaces implicit assumptions—extraction is as valuable as the final product
- Quantification enables analyses impossible with qualitative arguments alone—numbers matter even when uncertain
- But manual approaches cannot scale to match the challenge—we need computational leverage

This set the stage for AMTAIR's central innovation: using frontier language models to automate the extraction and formalization process while preserving the benefits MTAIR had demonstrated. Not to replace human judgment, but to amplify it—turning what took weeks into what takes hours, enabling comprehensive coverage rather than selective sampling.

2.6 Literature Review: Content and Technical Levels

The intellectual landscape surrounding AI risk resembles a rapidly expanding metropolis—new neighborhoods of thought spring up monthly, connected by bridges of varying stability to the established districts. A comprehensive review would fill volumes, so let me provide a guided tour of the territories most relevant to AMTAIR's mission.

2.6.1 AI Risk Models Evolution

The evolution of AI risk models traces a path from philosophical speculation to increasingly rigorous formalization—a journey from “what if?” to “how likely?”

Early Phase (2000-2010): The conversation began with broad conceptual arguments. Good's ultraintelligent machine **good1966** and Vinge's technological singularity set the stage, but these were more thought experiments than models. Yudkowsky's early writings **yudkowsky2008**

introduced key concepts like recursive self-improvement and orthogonality but remained largely qualitative.

Formalization Phase (2010-2018): Bostrom’s *Superintelligence* **bostrom2014** marked a watershed, providing systematic analysis of pathways, capabilities, and risks. The book’s genius lay not in mathematical formalism but in conceptual clarity—decomposing the nebulous fear of “robot overlords” into specific mechanisms like instrumental convergence and infrastructure profusion.

Quantification Phase (2018-present): Recent years have seen explicit probability estimates entering mainstream discourse. Carlsmith’s power-seeking model **carlsmith2022**, Cotra’s biological anchors, and various compute-based timelines represent attempts to put numbers on previously qualitative claims. The field increasingly recognizes that governance decisions require more than philosophical arguments—they need probability distributions.

This progression reflects a maturing field, though it also creates new challenges. As models become more quantitative, they risk false precision. As they become more complex, they risk inscrutability. AMTAIR attempts to navigate these tensions by preserving the narrative clarity of earlier work while enabling the mathematical rigor of recent approaches.

2.6.2 Governance Proposals Taxonomy

If risk models are the diagnosis, governance proposals are the treatment plans—and like medicine, they range from gentle interventions to radical surgery.

Technical Standards: The “first, do no harm” approach focuses on concrete safety requirements—interpretability benchmarks, robustness testing, capability thresholds. These proposals, exemplified by standard-setting bodies and technical safety organizations, offer specificity at the cost of narrowness.

Regulatory Frameworks: Moving up the intervention ladder, we find comprehensive regulatory proposals like the EU AI Act **european2024**. These create institutional structures, liability regimes, and oversight mechanisms, trading broad coverage for implementation complexity.

International Coordination: At the ambitious end, proposals for international AI governance treaties, soft law arrangements, and technical cooperation agreements aim to prevent races to the bottom. Think nuclear non-proliferation but for minds instead of missiles.

Research Priorities: Cutting across these categories, work by Dafoe **dafoe2018** and others maps the research landscape itself—what questions need answering before we can govern wisely? This meta-level analysis shapes funding flows and talent allocation.

A particularly compelling example of conditional governance thinking comes from “A Narrow Path” **miotti2024**, which proposes a phased approach: immediate safety measures to prevent uncontrolled development, international institutions to ensure stability, and long-term scientific foundations for beneficial transformative AI. This temporal sequencing—safety, stability, then flourishing—reflects growing sophistication in governance thinking.

2.6.3 Bayesian Network Theory and Applications

The mathematical machinery underlying AMTAIR rests on decades of theoretical development in probabilistic graphical models. Understanding this foundation helps appreciate both the power and limitations of the approach.

The key insight, crystallized in the work of Pearl **pearl2014** and elaborated by Koller & Friedman **koller2009**, is that independence relationships in complex systems can be read from graph structure. D-separation, the Markov condition, and the relationship between graphs and probability distributions provide the mathematical spine that makes Bayesian networks more than pretty pictures.

Critical concepts for AI risk modeling:

- **Conditional Independence:** Variable A is independent of C given B—encoded through graph separation
- **Markov Condition:** Each variable is independent of its non-descendants given its parents
- **Inference Algorithms:** From exact variable elimination to approximate Monte Carlo methods
- **Causal Interpretation:** When edges represent causal influence, the network supports counterfactual reasoning

These aren’t just mathematical niceties. When we claim that “deployment decisions” mediates the relationship between “capability advancement” and “catastrophic risk,” we’re making a precise statement about conditional independence that has testable implications.

2.6.4 Software Tools Landscape

The gap between Bayesian network theory and practical implementation is bridged by an ecosystem of software tools, each with its own strengths and opinions about how probabilistic reasoning should work.

pgmpy: This Python library provides the computational backbone for AMTAIR, offering both learning algorithms and inference engines. Its object-oriented design maps naturally onto our extraction pipeline.

NetworkX: For graph manipulation and analysis, NetworkX has become the de facto standard in Python, providing algorithms for everything from centrality measurement to community detection.

PyVis: Interactive visualization transforms static networks into explorable landscapes. PyVis’s integration with web technologies enables the rich interactive features that make formal models accessible.

Pandas/NumPy: The workhorses of scientific Python handle data manipulation and numerical computation, providing the infrastructure on which everything else builds.

The integration challenge—making these tools play nicely together while maintaining performance and correctness—shaped many architectural decisions in AMTAIR. Each tool excels in

its domain, but the seams between them required careful engineering.

2.6.5 Formalization Approaches

The challenge of formalizing natural language arguments extends far beyond AI risk, touching on fundamental questions in logic, linguistics, and artificial intelligence.

Pollock’s work on cognitive carpentry **pollock1995** provides philosophical grounding, arguing that human reasoning itself involves implicit formal structures that can be computationally modeled. This view—that formalization reveals rather than imposes structure—underlies AMTAIR’s approach.

Key theoretical challenges:

- **Semantic Preservation:** How do we maintain meaning while adding precision?
- **Structural Extraction:** What implicit relationships lurk in natural language?
- **Uncertainty Quantification:** How do we map “likely” to numbers?

Recent work on causal structure learning from text **babakov2025 ban2023 bethard2007** offers hope that these challenges can be addressed computationally. The convergence of large language models with formal methods creates new possibilities for bridging the semantic-symbolic gap.

2.6.6 Correlation Accounting Methods

One of the most persistent criticisms of Bayesian networks concerns their assumption of conditional independence given parents. In the real world, and especially in complex socio-technical systems like AI development, correlations abound.

Methods for handling these correlations have evolved considerably:

Copula Methods: By separating marginal distributions from dependence structure, copulas **nelson2006** allow modeling of complex correlations while preserving the Bayesian network framework. Think of it as adding a correlation layer on top of the basic network.

Hierarchical Models: Introducing latent variables that influence multiple observed variables captures correlations naturally. If “AI research culture” influences both “capability progress” and “safety investment,” their correlation is explained.

Explicit Correlation Nodes: Sometimes the most straightforward approach is best—directly model correlation mechanisms as additional nodes in the network.

Sensitivity Bounds: When correlations remain uncertain, compute best and worst case scenarios. This reveals when independence assumptions critically affect conclusions versus when they’re harmless simplifications.

For AMTAIR, the pragmatic approach dominates: start with independence assumptions, identify where they matter through sensitivity analysis, then selectively add correlation modeling where it most affects conclusions.

2.7 Methodology

The methodology of this research resembles less a linear march from hypothesis to conclusion and more an iterative dance between theory and implementation, vision and reality. Let me walk you through the choreography.

2.7.1 Research Design Overview

This research follows what methodologists might call a “design science” approach—we’re not just studying existing phenomena but creating new artifacts (the AMTAIR system) and evaluating their utility for solving practical problems (the coordination crisis in AI governance).

The overall flow:

1. **Theoretical Development:** Establishing why automated extraction could address the coordination crisis, grounded in epistemic theory and mechanism design
2. **Technical Implementation:** Building working software that demonstrates feasibility, not as a proof-of-concept toy but as a system capable of handling real arguments
3. **Empirical Validation:** Testing extraction quality against expert judgment, measuring not just accuracy but usefulness for downstream tasks
4. **Application Studies:** Applying the system to real AI governance questions, evaluating whether formal models actually enhance decision-making

This isn’t waterfall development where each phase completes before the next begins. Rather, insights from implementation fed back into theory, validation results shaped technical improvements, and application attempts revealed new requirements. The methodology itself embodied the iterative refinement it sought to enable.

2.7.2 Formalizing World Models from AI Safety Literature

The core methodological challenge—transforming natural language arguments into formal probabilistic models—requires careful consideration of what we’re actually trying to capture.

A “world model” in this context isn’t just any formal representation but specifically a causal model embodying beliefs about how different factors influence AI risk. The extraction approach must therefore:

- **Identify key variables:** Not just any entities mentioned, but causally relevant factors
- **Extract causal relationships:** Not mere correlation or co-occurrence, but directed influence
- **Capture uncertainty:** Both structural uncertainty (does A cause B?) and parametric uncertainty (how strongly?)
- **Preserve context:** Maintaining enough semantic information to interpret the formal model

Large language models enable this through sophisticated pattern recognition and reasoning capabilities, but they’re tools, not magic wands. The methodology must account for their strengths

(recognizing implicit structure) and weaknesses (potential hallucination, inconsistency).

2.7.3 From Natural Language to Computational Models

The journey from text to computation follows a carefully designed pipeline that mirrors human cognitive processes. Just as you wouldn't ask someone to simultaneously parse grammar and solve equations, we separate structural understanding from quantitative reasoning.

The Two-Stage Process:

Stage 1 focuses on structure—what causes what? The LLM reads an argument much as a human would, identifying key claims and their relationships. The prompt design here is crucial, providing enough guidance to ensure consistent extraction while allowing flexibility for different argument styles.

Stage 2 adds quantities—how likely is each outcome? With structure established, the system generates targeted questions about probabilities. This separation enables different approaches to quantification: extracting explicit estimates from text, inferring from qualitative language, or even connecting to external prediction markets.

The magic happens in the interplay. Structure constrains what probabilities are needed. Probability requirements might reveal missing structural elements. The process is a dialogue between qualitative and quantitative understanding.

2.7.4 Directed Acyclic Graphs: Structure and Semantics

At the mathematical heart of Bayesian networks lie Directed Acyclic Graphs (DAGs)—structures that are simultaneously simple enough to analyze and rich enough to capture complex phenomena.

The “directed” part encodes causality or influence—edges have direction, flowing from cause to effect. The “acyclic” part ensures logical coherence—you can't have A causing B causing C causing A, no matter how much certain political arguments might suggest otherwise.

Key properties for AI risk modeling:

Acyclicity: More than a mathematical convenience, this enforces coherent temporal or causal ordering. In AI risk arguments, this prevents circular reasoning where consequences justify premises that predict those same consequences.

D-separation: This graphical criterion determines conditional independence. If knowing about AI capabilities tells you nothing additional about risk given that you know deployment decisions, then capabilities and risk are d-separated given deployment.

Markov Condition: Each variable depends only on its parents, not on its entire ancestry. This locality assumption makes inference tractable and forces modelers to make intervention points explicit.

Path Analysis: Following paths through the graph reveals how influence propagates. Multiple paths between variables indicate redundancy—important for understanding intervention

robustness.

The causal interpretation, following Pearl’s framework, transforms these mathematical objects into tools for counterfactual reasoning. When we ask “what if we prevented deployment of misaligned systems?” we’re performing surgery on the DAG, setting variables and propagating consequences.

2.7.5 Quantification of Probabilistic Judgments

Here we encounter one of the most philosophically fraught aspects of the methodology: turning words into numbers. When an expert writes “highly likely,” what probability should we assign? When they say “significant risk,” what distribution captures their belief?

The methodology embraces rather than elides this challenge:

Calibration Studies: Research on human probability expression shows systematic patterns. “Highly likely” typically maps to 0.8-0.9, “probable” to 0.6-0.8, though individual and cultural variation is substantial.

Extraction Strategies: The system uses multiple approaches:

- Direct extraction: “We estimate 65% probability”
- Linguistic mapping: “Very likely” \rightarrow 0.85 (with uncertainty)
- Comparative extraction: “More likely than X” where $P(X)$ is known
- Bounded extraction: “At least 30%” \rightarrow [0.30, 1.0]

Uncertainty Representation: Rather than false precision, we maintain uncertainty about probabilities themselves. This might seem like uncertainty piled on uncertainty, but it’s honest—and mathematically tractable through hierarchical models.

The goal isn’t perfect extraction but useful extraction. If we can narrow “significant risk” from [0, 1] to [0.15, 0.45], we’ve added information even if we haven’t achieved precision.

2.7.6 Inference Techniques for Complex Networks

Once we’ve built these formal models, we need to reason with them—and here computational complexity rears its exponential head. The number of probability calculations required for exact inference grows exponentially with network connectivity, quickly overwhelming even modern computers.

The methodology employs a portfolio of approaches:

Exact Methods: For smaller networks (<30 nodes), variable elimination and junction tree algorithms provide exact answers. These form the gold standard against which we validate approximate methods.

Sampling Approaches: Monte Carlo methods trade exactness for scalability. By simulating many possible worlds consistent with our probability model, we approximate the true distributions. The law of large numbers is our friend here.

Variational Methods: These turn inference into optimization—find the simplest distribution that approximates our true beliefs. Like finding the best polynomial approximation to a complex curve.

Hybrid Strategies: Different parts of the network might use different methods. Exact inference for critical subgraphs, approximation for peripheral components.

The choice of method affects not just computation time but the types of questions we can meaningfully ask. This creates a methodological feedback loop where feasible inference shapes model design.

2.7.7 Integration with Prediction Markets and Forecasting Platforms

While full integration remains future work, the methodology anticipates connection to live forecasting data as a critical enhancement. The vision is compelling: formal models grounded in collective intelligence, updating as new information emerges.

The planned approach would involve:

Semantic Matching: Model variables rarely align perfectly with forecast questions. “AI causes human extinction” might map to multiple specific forecasts about capabilities, deployment, and impacts. Developing robust matching algorithms is essential.

Temporal Alignment: Markets predict specific dates (“AGI by 2030”) while models consider scenarios (“given AGI development”). Bridging these requires careful probability conditioning.

Quality Weighting: Not all forecasts are created equal. Platform reputation, forecaster track records, and market depth all affect reliability. The methodology must account for this heterogeneity.

Update Scheduling: Real-time updates would overwhelm users and computation. The system needs intelligent policies about when model updates provide value.

Platforms like Metaculus **tetlock2022** already demonstrate sophisticated conditional forecasting on AI topics. The challenge lies not in data availability but in meaningful integration that enhances rather than complicates decision-making.

With these theoretical foundations and methodological commitments established, we can now turn to the concrete implementation of AMTAIR. The next chapter demonstrates how these abstract principles translate into working software that addresses real governance challenges. The journey from theory to practice always involves surprises—some pleasant, others less so—but that’s what makes it interesting.

3. AMTAIR: Design and Implementation

Chapter Overview

Grade Weight: 20% | **Target Length:** ~29% of text (~8,700 words)

Requirements: Critical evaluation, strong argument for position, original contribution

The moment of truth in any research project comes when elegant theories meet stubborn reality. For AMTAIR, this meant transforming the vision of automated argument extraction into working code that could handle the beautiful messiness of real AI safety arguments. Let me take you through this journey from blueprint to implementation, complete with victories, defeats, and the occasional moment of “well, that’s unexpected.”

3.1 System Architecture Overview

Picture, if you will, a factory for transforming arguments into models. Raw materials (PDFs, blog posts, research papers) enter at one end. Finished products (interactive Bayesian networks) emerge at the other. In between lies a carefully orchestrated pipeline where each stage performs its specialized transformation, passing refined materials to the next.

The AMTAIR architecture embodies a philosophy: complex tasks become manageable when decomposed into focused components. Rather than building a monolithic “argument-to-model” black box, we created a series of specialized modules, each excellent at one thing.

The pipeline consists of five main stages:

1. **Text Ingestion and Preprocessing:** Like a careful librarian, this stage catalogues incoming documents, normalizes their format, extracts metadata, and identifies the argumentative content worth processing.
2. **Argument Extraction:** The intellectual heart of the system, where large language models perform their magic, transforming prose into structured representations.
3. **Data Transformation:** The workshop where extracted arguments are refined, validated, and prepared for mathematical representation.
4. **Network Construction:** The assembly line where formal Bayesian networks are instantiated, complete with conditional probability tables.

5. **Interactive Visualization:** The showroom where complex models become accessible through thoughtful design and interactivity.

3.1.1 Five-Stage Pipeline Architecture

Let's examine each stage more closely, understanding not just what they do but why they exist as separate components.

Text Ingestion and Preprocessing handles the unglamorous but essential work of standardization. Academic PDFs, with their two-column layouts and embedded figures, differ vastly from blog posts with inline code and hyperlinks. This stage creates a uniform representation while preserving essential structure and metadata. Format normalization strips away presentation while preserving content. Metadata extraction captures authorship, publication date, and citations. Relevance filtering identifies sections containing arguments rather than literature reviews or acknowledgments. Character encoding standardization prevents those maddening replacement characters that plague text processing.

Argument Extraction represents AMTAIR's core innovation. Using a two-stage process that mirrors human reasoning, it first identifies structural relationships (what influences what) then quantifies those relationships (how likely, how strong). This separation enables targeted prompts optimized for each task, human verification between stages, and modular improvements as LLM capabilities evolve.

Data Transformation bridges the gap between textual representations and mathematical models. It parses the BayesDown syntax into structured data, validates that the resulting network forms a proper DAG, checks probability consistency, and handles missing data intelligently.

Network Construction instantiates the formal mathematical model. This involves creating nodes and edges according to extracted structure, populating conditional probability tables, initializing inference engines, and validating the complete model.

Interactive Visualization makes the complex accessible. Through thoughtful visual encoding of probabilities and relationships, progressive disclosure of detail, interactive exploration capabilities, and multiple export formats, it serves diverse stakeholder needs.

3.1.2 Design Principles

Core Design Philosophy: The architecture embodies several principles that guided countless implementation decisions:

Modularity: Each component has clear inputs, outputs, and responsibilities. This isn't just good software engineering—it enables independent improvement of components and graceful degradation when parts fail.

Validation Checkpoints: Between each stage, we validate outputs before proceeding. Bad extractions don't propagate into visualization. Malformed networks trigger re-extraction rather than cryptic errors.

Human-in-the-Loop: While pursuing automation, we recognize that human judgment remains invaluable. The architecture provides natural intervention points where experts can verify and correct.

Extensibility: New document formats, improved extraction prompts, alternative visualization libraries—the architecture accommodates growth without restructuring.

The system emphasizes transparency over black-box efficiency. Users can inspect intermediate representations, understand extraction decisions, and verify transformations. This builds trust—essential for a system handling high-stakes arguments about existential risk.

3.2 The Two-Stage Extraction Process

The heart of AMTAIR beats with a two-stage rhythm: structure, then probability. This separation, which initially seemed like an implementation detail, revealed itself as fundamental to the extraction challenge.

3.2.1 Stage 1: Structural Extraction (ArgDown)

Imagine reading a complex argument about AI risk. Your first pass likely isn’t calculating exact probabilities—you’re mapping the landscape. What are the key claims? How do they relate? What supports what? Stage 1 mirrors this cognitive process.

The extraction begins with pattern recognition. Natural language contains linguistic markers of causal relationships: “leads to,” “results in,” “depends on,” “influences.” The LLM, trained on vast corpora of argumentative text, recognizes these patterns and their variations.

Consider extracting from a passage like: “The development of artificial general intelligence will likely lead to rapid capability gains through recursive self-improvement. This intelligence explosion could result in systems pursuing convergent instrumental goals, potentially including resource acquisition and self-preservation. Without solved alignment, such power-seeking behavior poses existential risks to humanity.”

The system identifies three key variables connected by causal relationships:

- AGI Development → Intelligence Explosion
- Intelligence Explosion → Power-Seeking Behavior
- Power-Seeking Behavior → Existential Risk

But extraction goes beyond simple pattern matching. The system must handle complex linguistic phenomena like coreference (“this,” “such systems”), implicit relationships, conditional statements, and negative statements. The magic lies in prompt engineering that guides the LLM to consistent extraction while remaining flexible enough for diverse argument styles.

The output, formatted in ArgDown syntax, preserves both structure and semantics:

```
[Existential_Risk]: Threat to humanity's continued existence and flourishing.
+ [Power_Seeking_Behavior]: AI systems pursuing instrumental goals like resource acquisition
```

- + [Intelligence_Explosion]: Rapid recursive self-improvement leading to superintelligence
- + [AGI_Development]: Creation of artificial general intelligence systems.

3.2.2 Stage 2: Probability Integration (BayesDown)

With structure established, Stage 2 adds the quantitative flesh to the qualitative bones. This stage faces a different challenge: extracting numerical beliefs from text that often expresses uncertainty in frustratingly vague terms.

The process begins by generating targeted questions based on the extracted structure. For each node, we need prior probabilities. For each child-parent relationship, we need conditional probabilities. The combinatorics can be daunting—a node with three binary parents requires 8 conditional probability values.

The system employs multiple strategies for probability extraction:

Explicit Extraction: When authors provide numerical estimates (“we assign 70% probability”), extraction is straightforward, though we must handle various formats and contexts.

Linguistic Mapping: Qualitative expressions map to probability ranges based on calibration studies. “Highly likely” becomes approximately 0.85, though we maintain uncertainty about this mapping.

Comparative Reasoning: Statements like “more probable than not” or “at least as likely as X” provide bounds even without exact values.

Coherence Enforcement: Probabilities must sum correctly. If $P(A|B) = 0.7$, then $P(\text{not } A|B)$ must equal 0.3. The system detects and resolves inconsistencies.

The result is a complete BayesDown specification:

json

```
[Existential_Risk]: Threat to humanity's continued existence. {
  "instantiations": ["true", "false"],
  "priors": {"p(true)": "0.10", "p(false)": "0.90"},
  "posteriors": {
    "p(true|power_seeking_true)": "0.65",
    "p(true|power_seeking_false)": "0.001"
  }
}
```

3.2.3 Why Two Stages?

The separation of structure from probability isn’t merely convenient—it’s cognitively valid and practically essential. Let me count the ways this design decision pays dividends:

Cognitive Alignment: Humans naturally separate “what relates to what” from “how likely is it.” The two-stage process mirrors this, making the system’s operation intuitive and interpretable.

Error Isolation: Structural errors (missing a key variable) differ fundamentally from probability errors (estimating 0.7 instead of 0.8). Separating stages allows targeted debugging and improvement.

Modular Validation: Experts can verify structure without needing to evaluate every probability. This enables efficient human oversight at natural checkpoints.

Flexible Quantification: Different probability sources (text extraction, expert elicitation, market data) can feed into the same structure. The architecture accommodates multiple approaches to the probability challenge.

Transparency: Users can inspect ArgDown to understand what was extracted before probabilities were added. This builds trust and enables meaningful correction.

The two-stage approach also revealed an unexpected benefit: ArgDown itself became a valuable output. Researchers began using these structural extractions for qualitative analysis, even without probability quantification. Sometimes, just making argument structure explicit provides sufficient value.

3.3 Implementation Technologies

Choosing technologies for AMTAIR resembled assembling a band—each instrument needed to excel individually while harmonizing with the ensemble. The selection criteria balanced capability, maturity, interoperability, and community support.

3.3.1 Technology Stack

The final ensemble performs beautifully:

Component	Technology	Purpose	Why This Choice
Language Models	GPT-4, Claude	Argument extraction	State-of-the-art reasoning capabilities
Network Analysis	NetworkX	Graph algorithms	Mature, comprehensive, well-documented
Probabilistic Modeling	pgmpy	Bayesian operations	Native Python, active development
Visualization	PyVis	Interactive rendering	Web-based, customizable, responsive
Data Processing	Pandas	Structured manipulation	Industry standard, powerful operations

Language Models form the cognitive core. GPT-4 and Claude demonstrate remarkable ability to understand complex arguments, recognize implicit structure, and maintain coherence across

long extractions. The choice to support multiple models provides robustness and allows leveraging their complementary strengths.

NetworkX handles all graph-theoretic heavy lifting. From basic operations like cycle detection to advanced algorithms like centrality measurement, it provides a comprehensive toolkit that would take years to replicate.

pgmpy bridges the gap between graph structure and probabilistic reasoning. Its clean API design maps naturally onto our extracted representations, while its inference algorithms handle the computational complexity of Bayesian reasoning.

PyVis transforms static networks into living documents. Built on vis.js, it provides smooth physics simulations, rich interactivity, and extensive customization options—all accessible through Python.

Pandas might seem mundane compared to its companions, but it’s the reliable rhythm section that keeps everything together. Its ability to reshape, merge, and transform structured data makes the complex data transformations tractable.

3.3.2 Key Algorithms

Beyond the libraries lie custom algorithms that address AMTAIR-specific challenges:

Hierarchical Parsing: The algorithm that transforms indented ArgDown text into structured data represents a small miracle of recursive descent parsing adapted for our custom syntax. It maintains parent-child relationships while handling edge cases like repeated nodes and complex dependencies.

python

```
def parse_hierarchy(text, current_indent=0):
    """Recursively parse indented structure maintaining relationships"""
    # Track nodes at each level for parent identification
    # Handle repeated nodes by reference
    # Validate DAG property during construction
```

Probability Completion: Real arguments rarely specify all required probabilities. Our completion algorithm uses maximum entropy principles—when uncertain, assume maximum disorder. This provides conservative estimates that can be refined with additional information.

Visual Encoding: The algorithm mapping probabilities to colors uses perceptual uniformity. The green-to-red gradient isn’t linear in RGB space but follows human perception of color difference. Small details, big impact on usability.

Layout Optimization: Force-directed layouts often produce “hairballs” for complex networks. Our customized approach uses hierarchical initialization based on causal depth, then refines with physics simulation. The result: layouts that reveal structure rather than obscuring it.

3.3.3 (Expected) Performance Characteristics

Performance in a system like AMTAIR involves multiple dimensions—speed, accuracy, scalability. Let’s examine what theoretical analysis and design considerations suggest about system behavior.

Computational Complexity: The extraction phase exhibits linear complexity in document length—processing twice as much text takes roughly twice as long. However, the inference phase faces exponential complexity in network connectivity. A fully connected network with n binary nodes requires $O(2^n)$ operations for exact inference. This fundamental limitation shapes practical usage patterns.

Practical Implications: Small networks (<20 nodes) enable real-time interaction with exact inference. Medium networks (20-50 nodes) require seconds to minutes depending on connectivity. Large networks (>50 nodes) necessitate approximate methods, trading accuracy for tractability. Very large networks push the boundaries of current methods.

The bottleneck shifts predictably: extraction remains manageable even for lengthy documents, but inference becomes challenging as models grow. This suggests a natural workflow—extract comprehensively, then focus on relevant subnetworks for detailed analysis.

Optimization Opportunities: Several strategies could improve performance: caching frequent inference queries, hierarchical decomposition of large networks, parallel processing for independent subgraphs, and progressive rendering for visualization. The modular architecture accommodates these enhancements without fundamental restructuring.

3.3.4 Deterministic vs. Probabilistic Components of the Workflow

An interesting philosophical question arises: in a system reasoning about probability, which components should themselves be probabilistic?

The current implementation draws a clear line:

Deterministic Components: All data transformations, graph algorithms, and inference calculations operate deterministically. Given the same input, they produce identical output. This provides reproducibility and debuggability—essential for building trust.

Probabilistic Components: The LLM calls for extraction introduce variability. Even with temperature set to 0, language models exhibit some randomness. Different runs might extract slightly different structures or probability estimates from the same text.

This division reflects a deeper principle: use determinism wherever possible, embrace probability where necessary. The extraction task—interpreting natural language—inherently involves uncertainty. But once we have formal representations, all subsequent operations should be predictable.

From an information-theoretic perspective, we’re trying to extract maximum information from documents within computational budget constraints. Each document contains some finite

amount of formalizable argument structure. Our goal is recovering as much as possible given realistic resource limits.

The two-stage extraction can be viewed as successive refinement—first recovering the higher-order bits (structure), then filling in lower-order bits (probabilities). This aligns with rate-distortion theory, where we get the most important information first.

3.4 Case Study: Rain-Sprinkler-Grass

Every field has its canonical examples—physics has spherical cows, economics has widget factories, and Bayesian networks have the rain-sprinkler-grass scenario. Despite its simplicity, this example teaches profound lessons about causal reasoning and serves as the perfect test case for AMTAIR.

3.4.1 Processing Steps

Let me walk you through how AMTAIR processes this foundational example:

The input arrives as a simple text description: “When it rains, the grass gets wet. The sprinkler also makes the grass wet. However, when it rains, we usually don’t run the sprinkler.”

From this prosaic description, the system performs five transformations:

1. **ArgDown Parsing:** Extract three variables (Rain, Sprinkler, Grass_Wet) and identify that rain influences both sprinkler usage and grass wetness, while the sprinkler also influences grass wetness.
2. **Question Generation:** Create probability queries: What’s $P(\text{Rain})$? What’s $P(\text{Sprinkler}|\text{Rain})$? What’s $P(\text{Grass_Wet}|\text{Rain}, \text{Sprinkler})$ for all combinations?
3. **BayesDown Extraction:** Either extract probabilities from text or apply reasonable defaults. The “usually don’t run” becomes $P(\text{Sprinkler}|\text{Rain}) = 0.01$.
4. **Network Construction:** Build the formal Bayesian network with three nodes, three edges, and complete conditional probability tables.
5. **Visualization Rendering:** Create an interactive display where rain appears as a root cause, influencing both sprinkler and grass directly.

Each step validates its outputs before proceeding, ensuring that errors don’t cascade through the pipeline.

3.4.2 Example Conversion Steps

Let’s trace the actual transformations to see the pipeline in action:

Initial ArgDown Extraction:

markdown

```
[Grass_Wet]: Concentrated moisture on grass blades. {"instantiations": ["wet", "dry"]}
+ [Rain]: Precipitation from the sky. {"instantiations": ["raining", "not_raining"]}
```

```
+ [Sprinkler]: Artificial watering system. {"instantiations": ["on", "off"]}
+ [Rain]
```

The hierarchy captures that rain influences sprinkler usage—a subtle but important causal relationship that pure correlation would miss.

Generated Questions for Probability Extraction:

markdown

```
/* Prior probabilities */
- What is the probability that it rains?
- What is the probability the sprinkler is on?

/* Conditional probabilities */
- What is the probability the sprinkler is on when it's raining?
- What is the probability the sprinkler is on when it's not raining?
- What is the probability the grass is wet when it's raining and sprinkler is on?
- [... and so on for all combinations]
```

The system generates exactly the questions needed to fully specify the network—no more, no less.

Complete BayesDown Result:

json

```
[Grass_Wet]: Concentrated moisture on grass. {
  "instantiations": ["wet", "dry"],
  "priors": {"p(wet)": "0.45", "p(dry)": "0.55"},
  "posteriors": {
    "p(wet|raining,on)": "0.99",
    "p(wet|raining,off)": "0.80",
    "p(wet|not_raining,on)": "0.90",
    "p(wet|not_raining,off)": "0.01"
  }
}
```

Notice how the probabilities tell a coherent story—grass is almost certainly wet if either water source is active, almost certainly dry if neither is.

Resulting DataFrame Structure:

The transformation into tabular format enables standard data analysis tools while preserving all relationships and probabilities. Each row represents a node with its properties, parents, children, and probability distributions.

3.4.3 Results

The successfully processed rain-sprinkler-grass example demonstrates several key capabilities:

Structure Preservation: The causal relationships—including the subtle influence of rain on sprinkler usage—are correctly captured and maintained throughout processing.

Probability Coherence: All probability distributions sum to 1.0, conditional probabilities are complete, and the values tell a plausible story.

Visual Clarity: The rendered network clearly shows rain as the root cause, influencing both sprinkler and grass, while sprinkler provides an additional pathway to wet grass.

Interactive Exploration: Users can click nodes to see detailed probabilities, drag to rearrange for clarity, and explore how changing parameters affects outcomes.

Inference Capability: The system correctly calculates derived probabilities like $P(\text{Rain}|\text{Grass_Wet})$ —the diagnostic reasoning from effect to cause that makes Bayesian networks so powerful.

This simple example validates the basic pipeline functionality. But the real test comes with complex, real-world arguments...

3.5 Case Study: Carlsmith’s Power-Seeking AI Model

From the gentle meadows of rain and sprinklers, we now ascend to the existential peaks of AI risk. Carlsmith’s model represents a dramatic increase in complexity—both conceptually and computationally. Where rain-sprinkler-grass has 3 nodes, Carlsmith involves 23. Where grass wetness is intuitive, “mesa-optimization” and “corrigibility” require careful thought.

3.5.1 Model Complexity

The numbers tell only part of the story:

- **23 nodes:** Each representing a substantive claim about AI development, deployment, or risk
- **29 edges:** Encoding causal relationships across technical, strategic, and societal domains
- **Multiple probability tables:** Many nodes have several parents, creating combinatorial explosion
- **Six-level causal depth:** From root causes to final catastrophe, influence propagates through multiple stages

But the conceptual complexity dwarfs the computational. Nodes like “APS-Systems” (Advanced, Planning, Strategically aware) encode specific technical hypotheses. Relationships like how “incentives to build” influence “deployment despite misalignment” require understanding of organizational behavior under competitive pressure.

This is no longer a toy problem but a serious attempt to formalize one of the most important arguments of our time.

3.5.2 Automated Extraction of the Carlsmith’s Argument Structure

The extraction process began with feeding Carlsmith’s paper to AMTAIR. Watching the system work felt like observing an archaeological excavation—layers of argument slowly revealed their structure.

The LLM prompts for extraction deserve special attention. Through iterative refinement, we developed prompts that guide extraction while remaining flexible:

python

```
ARGDOWN_EXTRACTION = PromptTemplate("""
You are extracting the causal model from an AI safety argument.
Focus on:
1. Identifying key variables that affect outcomes
2. Capturing causal relationships (not mere association)
3. Preserving the author's terminology where possible
4. Creating a directed acyclic graph structure

For Carlsmith's argument about power-seeking AI, pay special attention to:
- The chain from capabilities to catastrophe
- Conditional relationships (X matters only if Y)
- Technical preconditions for risk
""")
```

The extraction revealed Carlsmith’s elegant decomposition. At the highest level: capabilities enable power-seeking, which enables disempowerment, which constitutes catastrophe. But the details matter—deployment decisions mediated by incentives and deception, alignment difficulty influenced by multiple technical factors, corrective mechanisms that might interrupt the chain.

The ArgDown representation captured this structure:

```
[Existential_Catastrophe]: Permanent curtailment of humanity's potential
+ [Human_Disempowerment]: Humans lose control over future
+ [Scale_Of_Power_Seeking]: Power-seeking behavior becomes overwhelming
+ [Misaligned_Power_Seeking]: AI systems pursue problematic objectives
+ [APS_Systems]: Advanced, planning, strategically aware AI
+ [Alignment_Difficulty]: Hard to align such systems
+ [Deployment_Despite_Misalignment]: Systems deployed anyway
+ [Incentives_To_Build]: Strong pressure to develop AI
+ [Deception]: AI systems hide misalignment
```

The structure revealed insights. “Misaligned_Power_Seeking” emerged as a critical hub, influenced by multiple factors and influencing multiple outcomes. The pathway from incentives through deployment to risk became explicit.

3.5.3 From ArgDown to BayesDown in Carlsmith’s Model

Adding probabilities to Carlsmith’s structure presented unique challenges. Unlike rain-sprinkler probabilities that have intuitive values, what’s the probability of “mesa-optimization” or “deceptive alignment”?

The system generated over 100 probability questions for the full model. A sample:

For [Deployment_Decisions]:

- What is $P(\text{deploy})$?
- What is $P(\text{deploy}|\text{strong_incentives}, \text{deception})$?
- What is $P(\text{deploy}|\text{strong_incentives}, \text{no_deception})$?
- What is $P(\text{deploy}|\text{weak_incentives}, \text{deception})$?
- What is $P(\text{deploy}|\text{weak_incentives}, \text{no_deception})$?

Each question targets a specific parameter needed for the Bayesian network. The conditional structure reflects Carlsmith’s argument—deployment depends on both incentives (external pressure) and deception (hidden misalignment).

The LLM extraction drew on Carlsmith’s explicit estimates where available and inferred reasonable values elsewhere. The result captured both the structure and Carlsmith’s quantitative risk assessment:

json

```
[Deployment_Decisions]: Decisions to deploy potentially misaligned AI. {
  "instantiations": ["deploy", "withhold"],
  "priors": {"p(deploy)": "0.70", "p(withhold)": "0.30"},
  "posteriors": {
    "p(deploy|strong_incentives,deception)": "0.90",
    "p(deploy|strong_incentives,no_deception)": "0.75",
    "p(deploy|weak_incentives,deception)": "0.60",
    "p(deploy|weak_incentives,no_deception)": "0.30"
  }
}
```

The probabilities tell a plausible story: deployment becomes more likely with stronger incentives and successful deception, but even without deception, strong incentives create substantial deployment probability.

3.5.4 Practically Meaningful BayesDown

The BayesDown representation achieves something remarkable: it bridges the chasm between Carlsmith’s nuanced prose and mathematical formalism without losing the essence of either.

Consider what this bridge enables:

For Technical Researchers: The formal structure makes assumptions explicit. Is power-seeking really independent of capability level given strategic awareness? The model forces clarity.

For Policymakers: Probabilities attached to comprehensible descriptions provide actionable intelligence. “70% chance of deployment despite misalignment” translates better than abstract concerns.

For Strategic Analysts: The network structure reveals intervention points. Which nodes, if changed, most affect the final outcome? Where should we focus effort?

The hybrid nature—natural language plus formal structure plus probabilities—serves each audience while enabling communication between them. A policymaker can understand “deployment decisions” without probability theory. A researcher can analyze the mathematical model without losing sight of what the variables mean.

This isn’t just convenient—it’s essential for coordination. When different communities can refer to the same model but engage with it at their appropriate level of technical detail, we create common ground for productive disagreement and collaborative problem-solving.

3.5.5 Interactive Visualization and Exploration

The moment when Carlsmith’s model first rendered as an interactive network felt like putting on glasses after years of squinting. Suddenly, the complex web of relationships became navigable.

The visualization system employs multiple visual channels simultaneously:

Color Coding: Nodes shift from deep red (low probability) through yellow to bright green (high probability). At a glance, you see which factors Carlsmith considers likely versus speculative.

Border Styling: Blue borders mark root causes (like “Incentives_To_Build”), purple indicates intermediate nodes, magenta highlights final outcomes. The visual grammar guides the eye through causal flow.

Layout Algorithm: Initial placement uses causal depth—root causes at bottom, final outcomes at top. Physics simulation then refines positions to minimize edge crossings while preserving hierarchical structure.

Progressive Disclosure: Hovering reveals probability summaries. Clicking opens detailed conditional probability tables. Dragging allows custom arrangement. Each interaction level serves different analytical needs.

The implementation required careful attention to human factors:

python

```
def create_interactive_visualization(network_df):
    """Transform formal model into explorable landscape"""

    # Initialize with thoughtful defaults
    net = Network(height="720px", width="100%", directed=True)

    # Configure physics for clarity not just aesthetics
```

```

net.force_atlas_2based(
    gravity=-50,      # Gentle spread
    spring_length=150, # Readable spacing
    spring_strength=0.02 # Soft constraints
)

# Add nodes with rich metadata
for node in nodes:
    net.add_node(
        node_id,
        label=create_simple_label(node),      # "Deployment\np=0.70"
        title=create_rich_tooltip(node),      # Full probability details
        color=probability_to_color(node),     # Visual encoding
        borderWidth=3,                       # Visible borders
        shape="box"                          # Readable text
    )

```

The resulting visualization transforms abstract relationships into tangible understanding. Users report “aha” moments when exploring—suddenly seeing how technical factors compound into strategic risks, or identifying previously unnoticed bottlenecks in the causal chain.

3.5.6 Validation Against Original (From the MTAIR Project)

Validating AMTAIR’s extraction required careful comparison with expert judgment. While comprehensive benchmarking remains future work, preliminary validation efforts provide encouraging signals.

Manual Baseline Creation: Domain experts, including Johannes Meyer and Jelena Meyer, independently extracted ArgDown and BayesDown representations from Carlsmith’s paper. This created ground truth accounting for legitimate interpretive variation—experts might reasonably disagree on some structural choices or probability estimates.

Structural Comparison: Comparing extracted causal structures revealed high agreement on core relationships. AMTAIR consistently identified the main causal chain from capabilities through deployment to catastrophe. Some variation appeared in handling of auxiliary factors—where one expert might include a minor influence, another might omit it for simplicity.

Probability Assessment: Probability extraction showed greater variation, reflecting inherent ambiguity in translating qualitative language. When Carlsmith writes “likely,” different readers might reasonably interpret this as 0.7, 0.75, or 0.8. AMTAIR’s extractions fell within the range of expert interpretations, suggesting successful capture of intended meaning even if not identical numbers.

Semantic Preservation: Most importantly, the formal models preserved the essential insights of Carlsmith’s argument. The critical role of deployment decisions, the compound nature of

risk, the importance of technical and strategic factors—all emerged clearly in the extracted representations.

An ideal validation protocol would expand this approach:

1. Multiple expert extractors working independently
2. Systematic comparison of structural and quantitative agreement
3. Analysis of where and why extractions diverge
4. Testing whether different extractions lead to different policy conclusions
5. Iterative refinement based on identified failure modes

The goal isn’t perfect agreement—even human experts disagree. Rather, we seek extractions good enough to support meaningful analysis while acknowledging their limitations.

3.6 Validation Methodology

Building trust in automated extraction requires more than anecdotal success. We need systematic validation that honestly assesses both capabilities and limitations.

3.6.1 Ground Truth Construction

Creating ground truth for argument extraction poses unique challenges. Unlike named entity recognition or sentiment analysis, argument structure lacks universal standards. What constitutes the “correct” extraction from a complex text?

An ideal validation approach would embrace this inherent subjectivity:

Expert Selection: Recruit 5-10 domain experts with demonstrated expertise in both AI safety and formal modeling. Diversity matters—include technical researchers, policy analysts, and those with mixed backgrounds.

Extraction Protocol: Provide standardized training on ArgDown/BayesDown syntax while allowing flexibility in interpretation. Experts work independently to avoid anchoring bias, documenting their reasoning process alongside final extractions.

Consensus Building: Through structured discussion, identify areas of convergence (likely core argument structure) versus legitimate disagreement (interpretive choices, granularity decisions). This distinguishes system errors from inherent ambiguity.

Quality Metrics: Rather than binary correct/incorrect judgments, assess:

- Structural similarity (graph edit distance)
- Probability distribution overlap (KL divergence)
- Semantic preservation (expert ratings)
- Downstream task performance (policy analysis agreement)

The resulting dataset would capture not a single “truth” but a distribution of reasonable interpretations against which to evaluate automated extraction.

3.6.2 Evaluation Metrics

Evaluating argument extraction requires metrics that capture multiple dimensions of quality:

Structural Fidelity:

- Node identification: What fraction of expert-identified variables does the system extract?
- Edge accuracy: Are causal relationships preserved?
- Hierarchy preservation: Does the system maintain argument levels?

Probability Calibration:

- Explicit extraction: When sources state probabilities, how accurately are they captured?
- Linguistic mapping: Do qualitative expressions translate to reasonable probabilities?
- Coherence: Are probability distributions properly normalized?

Semantic Quality:

- Description accuracy: Do extracted descriptions preserve original meaning?
- Terminology preservation: Does the system maintain author’s vocabulary?
- Context retention: Is sufficient information preserved for interpretation?

Functional Validity:

- Inference agreement: Do extracted models support similar conclusions?
- Sensitivity preservation: Are critical parameters identified as influential?
- Policy robustness: Do different extractions suggest similar interventions?

These metrics acknowledge that perfect extraction is neither expected nor necessary. The goal is extraction sufficient for practical use while maintaining transparency about limitations.

3.6.3 Results Summary

While comprehensive validation remains future work, preliminary assessments using the methodology described above would likely reveal several patterns:

Expected Strengths: Automated extraction should excel at identifying explicit causal claims, preserving hierarchical argument structure, and extracting stated probabilities. The two-stage approach likely improves quality by allowing focused optimization for each task.

Anticipated Challenges: Implicit reasoning, complex conditionals, and ambiguous quantifiers would pose greater challenges. Coreference resolution across long documents and maintaining consistency in large models would require continued refinement.

Practical Utility Threshold: Even with imperfect extraction, the system could provide value if it achieves perhaps 70-80% structural accuracy and captures probability estimates within reasonable ranges. This level of performance would enable rapid initial modeling that experts could refine, dramatically reducing the time from argument to formal model.

The validation framework itself represents a contribution—establishing systematic methods for assessing argument extraction quality as this research area develops.

3.6.4 Error Analysis

Understanding failure modes guides both appropriate use and future improvements:

Implicit Assumptions: Authors often leave critical assumptions unstated, relying on shared background knowledge. When an AI safety researcher writes about “alignment,” they assume readers understand the technical concept. The system must either extract these implicit elements or flag their absence.

Complex Conditionals: Natural language expresses conditionality in myriad ways. “If we achieve alignment (which seems unlikely without major theoretical breakthroughs), then deployment might be safe (assuming robust verification).” Parsing nested, qualified conditionals challenges current methods.

Ambiguous Quantifiers: The word “significant” might mean 10% in one context, 60% in another. Without calibration to author-specific usage or domain conventions, probability extraction remains approximate.

Coreference Challenges: Academic writing loves pronouns and indirect references. When “this approach” appears three paragraphs after introducing multiple approaches, identifying the correct referent requires sophisticated discourse understanding.

These limitations don’t invalidate the approach but rather define its boundaries. Users who understand these constraints can work within them, leveraging automation’s strengths while compensating for its weaknesses.

3.7 Policy Evaluation Capabilities

The ultimate test of a model isn’t its elegance but its utility. Can AMTAIR’s extracted models actually inform governance decisions? This section demonstrates how formal models enable systematic policy analysis.

3.7.1 Intervention Representation

Representing policy interventions in Bayesian networks requires translating governance mechanisms into parameter modifications. Pearl’s do-calculus provides the mathematical framework, but the practical challenge lies in meaningful translation.

An ideal implementation would support several intervention types:

Parameter Modification: Policies often change probabilities. Safety requirements might reduce $P(\text{deployment}|\text{misaligned})$ from 0.7 to 0.2 by making unsafe deployment legally prohibited or reputationally costly.

Structural Interventions: Some policies add new causal pathways. Introducing mandatory review boards creates new nodes and edges representing oversight mechanisms.

Uncertainty Modeling: Policy effectiveness is itself uncertain. Rather than assuming perfect implementation, represent ranges: $P(\text{deployment}|\text{misaligned})$ might become $[0.1, 0.3]$ depending

on enforcement.

Multi-Level Effects: Policies influence multiple levels simultaneously. Compute governance affects technical development, corporate behavior, and international competition.

The system would translate high-level policy descriptions into specific network modifications, enabling rigorous counterfactual analysis of intervention effects.

3.7.2 Example: Deployment Governance

Let’s trace how a specific policy—mandatory safety certification before deployment—might be evaluated:

Baseline Model: In Carlsmith’s original model, $P(\text{deployment}|\text{misaligned}) = 0.7$, reflecting competitive pressures overwhelming safety concerns.

Policy Specification: Safety certification requires demonstrating alignment properties before deployment authorization. Based on similar regulations in other domains, we might estimate 80-90% effectiveness.

Parameter Update: The modified model sets $P(\text{deployment}|\text{misaligned}) = 0.1\text{-}0.2$, representing the residual probability of circumvention or regulatory capture.

Downstream Effects:

- Reduced deployment of misaligned systems
- Lower probability of power-seeking manifestation
- Decreased existential risk from $\sim 5\%$ to $\sim 1.2\%$

Sensitivity Analysis: How robust is this conclusion? Varying certification effectiveness, enforcement probability, and other parameters reveals which assumptions critically affect the outcome.

This example illustrates policy evaluation’s value: moving from vague claims (“regulation would help”) to quantitative assessments (“this specific intervention might reduce risk by $75\% \pm 15\%$ ”).

3.7.3 Robustness Analysis

Good policies work across scenarios. AMTAIR enables testing interventions against multiple worldviews, parameter ranges, and structural variations.

Cross-Model Testing: Extract multiple expert models and evaluate the same policy in each. If an intervention reduces risk in Carlsmith’s model but increases it in Christiano’s, we’ve identified a critical dependency.

Parameter Sensitivity: Which uncertainties most affect policy effectiveness? If the intervention only works for $P(\text{alignment_difficulty}) < 0.3$, and experts disagree whether it’s 0.2 or 0.4, we need more research before implementing.

Structural Uncertainty: Some disagreements concern model structure itself. Does capability advancement directly influence misalignment risk, or only indirectly through deployment

pressures? Test policies under both structures.

Confidence Bounds: Rather than point estimates, compute ranges. “This policy reduces risk by 40-80%” honestly represents uncertainty while still providing actionable guidance.

The goal isn’t eliminating uncertainty but making decisions despite it. Robustness analysis reveals which policies work across uncertainties versus those requiring specific assumptions.

3.8 Interactive Visualization Design

A Bayesian network without good visualization is like a symphony without performers—all potential, no impact. The visualization system transforms mathematical abstractions into intuitive understanding.

3.8.1 Visual Encoding Strategy

Every visual element carries information:

Color: The probability spectrum from red (low) through yellow to green (high) provides immediate gestalt understanding. Pre-attentive processing—the brain’s ability to process certain visual features without conscious attention—makes patterns jump out.

Borders: Node type encoding (blue=root, purple=intermediate, magenta=outcome) creates visual flow. The eye naturally follows from blue through purple to magenta, tracing causal pathways.

Size: Larger nodes have higher centrality—more connections, more influence. This emerges from the physics simulation but reinforces importance.

Layout: Force-directed positioning naturally clusters related concepts while maintaining readability. The algorithm balances competing constraints: minimize edge crossings, maintain hierarchical levels, avoid node overlap, and create aesthetic appeal.

The encoding philosophy: every pixel should earn its place by conveying information while maintaining visual harmony.

3.8.2 Progressive Disclosure

Information overload kills understanding. The interface reveals complexity gradually:

Level 1 - Overview: At first glance, see network structure and probability color coding. This answers: “What’s the shape of the argument? Where are the high-risk areas?”

Level 2 - Hover Details: Mouse over a node to see its description and prior probability. This adds: “What does this factor represent? How likely is it?”

Level 3 - Click Deep Dive: Clicking opens full probability tables and relationships. This reveals: “How does this probability change with conditions? What influences this factor?”

Level 4 - Interactive Exploration: Dragging, zooming, and physics controls enable custom investigation. This supports: “What if I reorganize to see different patterns? How do these clusters relate?”

Each level serves different users and use cases. A policymaker might work primarily with levels 1-2, while a researcher dives into level 3-4 details.

3.8.3 User Interface Elements

Effective interface design for Bayesian networks requires balancing power with accessibility:

Physics Controls: Force-directed layouts benefit from tuning. Gravity affects spread, spring length controls spacing, damping influences settling time. Advanced users can adjust these for optimal layouts, while defaults work well for most cases.

Filter Options: With large networks, selective viewing becomes essential. Filter by probability ranges (show only likely events), node types (focus on interventions), or causal depth (see only immediate effects).

Export Functions: Different stakeholders need different formats. Researchers want raw data, policymakers need reports, presenters require images. Supporting diverse export formats enables broad usage.

Comparison Mode: Understanding often comes from contrast. Side-by-side viewing of baseline versus intervention, or different expert models, reveals critical differences.

Iterative design with actual users would refine these features, ensuring they serve real needs rather than imagined ones.

3.9 Integration with Prediction Markets

The vision: formal models that breathe with live data, updating as collective intelligence evolves. While full implementation awaits, the architecture anticipates this future.

3.9.1 Design for Integration

Integration Architecture requires careful design to manage the impedance mismatch between formal models and market data:

API Specifications: Each platform—Metaculus, Manifold, Good Judgment Open—has unique data formats, update frequencies, and question types. A unified adapter layer would translate platform-specific formats into model-compatible data.

Semantic Matching: The hard problem—connecting “AI causes extinction by 2100” (market question) to “Existential_Catastrophe” (model node). This requires sophisticated NLP and possibly human curation for high-stakes connections.

Aggregation Methods: When multiple markets address similar questions, how do we combine? Weighted averages based on market depth, participant quality, and historical accuracy provide

more signal than simple means.

Update Scheduling: Real-time updates would overwhelm users and computation. Smart scheduling might update daily for slow-changing strategic questions, hourly for capability announcements, immediately for critical events.

3.9.2 Challenges and Opportunities

The challenges are real but surmountable:

Question Mapping: Markets ask specific, time-bound questions while models represent general relationships. “AGI by 2030?” maps uncertainly to “APS_Systems exists.” Developing robust mapping functions requires deep understanding of both domains.

Temporal Alignment: Market probabilities change over time, but model parameters are typically static. Should we use current market values, time-weighted averages, or attempt to extract trend information?

Quality Variation: A liquid market with expert participants provides different information than a thin market with casual forecasters. Weighting schemes must account for these quality differences.

Incentive Effects: If models influence policy and policy influences outcomes, and markets forecast outcomes, we create feedback loops. Understanding these dynamics prevents perverse incentives.

Despite challenges, even partial integration provides value:

- External validation of expert-derived probabilities
- Dynamic updating as new information emerges
- Identification of where model and market disagree
- Quantified uncertainty from market spread

The perfect shouldn’t be the enemy of the good—simple integration beats no integration.

3.10 Computational Performance Analysis

As networks grow from toy examples to real-world complexity, computational challenges emerge. Understanding these constraints shapes realistic expectations and optimization priorities.

3.10.1 Exact vs. Approximate Inference

The fundamental tradeoff in probabilistic reasoning: exactness versus tractability.

Exact Inference: Variable elimination and junction tree algorithms provide mathematically exact answers. For our 3-node rain-sprinkler network, calculations complete instantly. For 20-node networks with modest connectivity, expect seconds. But for 50+ node networks with complex dependencies, exact inference becomes impractical—potentially taking hours or exhausting memory.

Approximate Methods: When exactness becomes impractical, approximation saves the day:

- **Monte Carlo Sampling:** Generate thousands of scenarios consistent with the network, estimate probabilities from frequencies. Accuracy improves with samples, trading computation time for precision.
- **Variational Inference:** Find the simplest distribution that approximates our complex reality. Like fitting a smooth curve to jagged data—we lose detail but gain comprehension.
- **Belief Propagation:** Pass messages between nodes until beliefs converge. Works beautifully for tree-structured networks, can oscillate or converge slowly for complex loops.

The system selects methods based on network properties:

- Small networks: exact inference for precision
- Medium networks: belief propagation for speed
- Large networks: sampling for scalability
- Very large networks: hierarchical decomposition

3.10.2 Scaling Strategies

When networks grow beyond convenient computation, clever strategies maintain usability:

Hierarchical Decomposition: Break large networks into smaller, manageable subnetworks. Compute locally, then integrate results. Like solving a jigsaw puzzle by completing sections before assembling the whole.

Relevance Pruning: For specific queries, most nodes don't matter. If asking about deployment risk, technical details about interpretability methods might be temporarily ignorable. Prune irrelevant subgraphs for focused analysis.

Caching Architecture: Many queries repeat— $P(\text{catastrophe})$, $P(\text{deployment}|\text{misalignment})$. Cache results to avoid recomputation. Smart invalidation updates only affected queries when parameters change.

Parallel Processing: Inference calculations often decompose naturally. Different branches of the network can be processed simultaneously. Modern multi-core processors and cloud computing make this increasingly attractive.

Implementation would balance these strategies based on usage patterns. Interactive exploration benefits from caching and pruning. Batch analysis leverages parallelization. The architecture accommodates multiple approaches.

3.11 Results and Achievements

3.11.1 Extraction Quality Assessment

Assessing extraction quality requires honesty about both achievements and limitations. An ideal evaluation would examine multiple dimensions:

Coverage: What proportion of arguments in source texts does the system successfully capture? Initial applications suggest the two-stage approach identifies most explicit causal claims while struggling with deeply implicit relationships.

Accuracy: How closely do automated extractions match expert consensus? Preliminary comparisons indicate strong agreement on primary causal structures with more variation in probability estimates.

Robustness: How well does the system handle different writing styles, argument structures, and domains? Academic papers with clear argumentation extract more reliably than informal blog posts or policy documents.

Utility: Do the extracted models enable meaningful analysis? Even imperfect extractions that capture 80% of structure with approximate probabilities can dramatically accelerate modeling compared to starting from scratch.

The key insight: perfect extraction isn't necessary for practical value. Like machine translation, which provides useful results despite imperfections, automated argument extraction can enhance human capability without replacing human judgment.

3.11.2 Computational Performance

Performance analysis would reveal the practical boundaries of the current system:

Extraction Speed: LLM-based extraction scales roughly linearly with document length. A 20-page paper might require 30-60 seconds for structural extraction and similar time for probability extraction. This enables processing dozens of documents daily—orders of magnitude faster than manual approaches.

Network Complexity Limits: Exact inference remains tractable for networks up to approximately 30-40 nodes with moderate connectivity. Beyond this, approximate methods become necessary, with sampling methods scaling to hundreds of nodes at the cost of precision.

Visualization Responsiveness: Interactive visualization performs smoothly for networks under 50 nodes. Larger networks benefit from hierarchical viewing or focus+context techniques to maintain usability.

End-to-End Pipeline: From document input to interactive visualization, expect 2-5 minutes for typical AI safety arguments. This represents roughly 100x speedup compared to manual modeling efforts.

These performance characteristics make AMTAIR practical for real-world use while highlighting areas for future optimization.

3.11.3 Policy Impact Evaluation

The true test of AMTAIR lies in its ability to inform governance decisions. An ideal policy evaluation framework would demonstrate several capabilities:

Intervention Modeling: Representing diverse policy proposals—from technical standards to international agreements—as parameter modifications in extracted networks. This translation from qualitative proposals to quantitative changes enables rigorous analysis.

Comparative Assessment: Evaluating multiple interventions across different expert world-views to identify robust strategies. Policies that reduce risk across different models deserve priority over those requiring specific assumptions.

Sensitivity Analysis: Understanding which uncertainties most affect policy conclusions. If an intervention’s effectiveness depends critically on disputed parameters, this highlights research priorities.

Implementation Guidance: Moving beyond “this policy reduces risk” to specific recommendations about design details, implementation sequences, and success metrics.

The system would transform abstract policy discussions into concrete quantitative analyses, enabling evidence-based decision-making in AI governance.

3.12 Summary of Technical Contributions

Looking back at the implementation journey, several achievements stand out:

Automated Extraction: The two-stage pipeline successfully transforms natural language arguments into formal models, achieving practical accuracy while maintaining transparency about limitations.

Hybrid Representation: BayesDown bridges qualitative and quantitative worlds, preserving semantic richness while enabling mathematical analysis.

Scalable Architecture: Modular design accommodates growth—new document types, improved extraction methods, additional visualization options—without fundamental restructuring.

Interactive Accessibility: Thoughtful visualization makes complex models understandable to diverse stakeholders, democratizing access to formal reasoning tools.

Policy Relevance: The ability to model interventions and assess robustness transforms academic exercises into practical governance tools.

These technical achievements validate the feasibility of computational coordination infrastructure for AI governance. Not as a complete solution, but as a meaningful enhancement to human judgment and collaboration.

The implementation demonstrates that the vision of automated argument extraction is not merely theoretical but practically achievable. While challenges remain—particularly in handling implicit reasoning and diverse uncertainty expressions—the system provides a foundation for enhanced coordination in AI governance.

The journey from concept to implementation revealed unexpected insights. The two-stage extraction process, initially a pragmatic choice, proved cognitively valid. The intermediate rep-

resentations became valuable outputs themselves. The visualization challenges led to design innovations applicable beyond this project.

Most importantly, the implementation confirms that formal modeling of AI risk arguments need not remain the province of a few dedicated experts. Through automation and thoughtful design, these powerful tools can serve the broader community working to ensure advanced AI benefits humanity.

Having demonstrated technical feasibility and practical utility, we must now critically examine limitations, address objections, and explore broader implications. The next chapter undertakes this essential reflection, ensuring we neither oversell the approach nor undervalue its contributions.

3.5.6 Validation Against Original (From the MTAIR Project)

To validate the AMTAIR extraction quality, an ideal approach would involve systematic comparison with expert manual extractions. The validation methodology would follow these steps:

Expert Baseline Creation: Multiple domain experts independently extract ArgDown and BayesDown representations from the same source documents. This creates a ground truth dataset accounting for legitimate variation in expert interpretation.

Structural Comparison: Compare the extracted causal structures, examining:

- Node identification completeness
- Relationship preservation accuracy
- Hierarchical organization fidelity
- Handling of repeated or complex dependencies

Probability Assessment: Evaluate probability extraction by:

- Comparing explicit probability captures
- Assessing interpretation of qualitative expressions
- Measuring consistency across related probabilities
- Identifying systematic biases or tendencies

Semantic Preservation: Expert review would assess whether the formal representation maintains the essential meaning and nuance of the original arguments.

For the Carlsmith model specifically, preliminary manual extractions by domain experts (including Johannes Meyer and Jelena Meyer) suggest that automated extraction achieves high structural fidelity—capturing the key variables and their relationships—while probability estimates show greater variation, reflecting the inherent ambiguity in translating qualitative language to quantitative values.

3.6 Validation Methodology

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

3.6.1 Ground Truth Construction

An ideal validation protocol would establish ground truth through:

Expert Selection: Recruit domain experts with both AI safety knowledge and experience in formal modeling. Experts should represent diverse perspectives within the field to capture legitimate interpretive variation.

Standardized Training: Provide consistent training on ArgDown/BayesDown syntax and extraction principles. This ensures methodological alignment while preserving substantive differences in interpretation.

Independent Extraction: Have experts work independently on the same source documents, preventing anchoring bias and capturing the natural range of valid interpretations.

Consensus Building: Through structured discussion, identify areas of convergence and legitimate disagreement. This distinguishes extraction errors from genuine ambiguity in source materials.

Documentation: Record not just final extractions but the reasoning process, creating rich data for understanding extraction challenges and improving automated approaches.

3.6.2 Evaluation Metrics

A comprehensive evaluation framework would assess multiple dimensions:

Structural Metrics:

- Node identification precision and recall
- Edge extraction accuracy
- Preservation of hierarchical relationships
- Handling of complex dependencies

Probability Metrics:

- Accuracy of explicit probability extraction
- Consistency in interpreting qualitative expressions
- Preservation of conditional relationships
- Handling of uncertainty about uncertainty

Semantic Metrics:

- Expert ratings of meaning preservation
- Functional equivalence for key inferences
- Preservation of author’s argumentative intent
- Appropriate simplification choices

Pragmatic Metrics:

- Usefulness for downstream analysis
- Time savings versus manual extraction
- Error patterns and failure modes
- Robustness across document types

3.6.3 Results Summary

While comprehensive validation remains future work, preliminary assessments suggest:

Structural Extraction: The two-stage approach successfully identifies major causal relationships and preserves hierarchical structure. The explicit ArgDown intermediate representation allows verification of structural accuracy before probability quantification.

Probability Challenges: Converting qualitative expressions to numerical probabilities remains the primary challenge. Different experts interpret phrases like “likely” or “significant risk” differently, and automated extraction inherits this ambiguity.

Practical Utility: Despite imperfections, automated extraction provides sufficient quality for many practical applications, especially when combined with human review at critical points.

The validation framework itself represents a contribution, providing systematic methods for assessing automated argument formalization tools as this area develops.

3.6.4 Error Analysis

Common failure modes to avoid:

Implicit Assumptions: Unstated background assumptions that experts infer but system misses. These often involve domain-specific common knowledge that remains unspoken in expert discourse.

Complex Conditionals: Nested conditionals with multiple antecedents challenge current parsing. Statements like “If A and B, then probably C, unless D” require sophisticated logical analysis.

Ambiguous Quantifiers: Terms like “significant” lack clear probability mapping without context. The same word may imply different probabilities in different domains or even different parts of the same argument.

Coreference Resolution: Pronouns and indirect references create attribution challenges. When authors use “this risk” or “that assumption,” identifying the correct referent requires deep contextual understanding.

Understanding these limitations guides both current usage and future improvements.

3.7 Policy Evaluation Capabilities

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

3.7.1 Intervention Representation

Policy interventions can be modeled as modifications to network parameters, following Pearl’s do-calculus framework. An ideal implementation would:

Parameter Modification: Represent policies as changes to specific probability values. For instance, safety requirements might reduce $P(\text{deployment}|\text{misaligned})$ by making unsafe deployment less likely.

Structural Interventions: Some policies add or remove causal pathways. Regulatory oversight might introduce new nodes representing approval processes.

Uncertainty Propagation: Model uncertainty about policy effectiveness. Rather than assuming perfect implementation, represent ranges of possible effects.

Multi-Level Effects: Capture how policies influence multiple levels simultaneously—technical development, corporate behavior, and international dynamics.

The formal framework enables rigorous counterfactual reasoning: “What would happen to existential risk if this policy were implemented?”

3.7.2 Example: Deployment Governance

Consider a hypothetical policy requiring safety certification before deployment:

Baseline Scenario: Without intervention, the model might show $P(\text{deployment}|\text{misaligned}) = 0.7$, reflecting competitive pressures to deploy despite risks.

Policy Intervention: Safety certification requirements could reduce this to $P(\text{deployment}|\text{misaligned}) = 0.1$, assuming effective enforcement.

Downstream Effects: This change propagates through the network:

- Reduced deployment of misaligned systems
- Lower probability of power-seeking behavior manifestation
- Decreased existential risk

Quantitative Assessment: The formal model enables precise calculation of risk reduction, helping prioritize among possible interventions.

This example illustrates how formal models transform vague policy discussions into concrete quantitative analyses, though the specific numbers depend on model assumptions and parameter estimates.

3.7.3 Robustness Analysis

Policies must work across worldviews. AMTAIR enables multi-model evaluation, parameter sensitivity testing, scenario analysis, and confidence bound computation—ensuring interventions remain effective despite uncertainty.

Cross-Model Testing: Evaluate policies across different extracted worldviews to identify robust strategies that work regardless of which expert’s model proves correct.

Sensitivity Analysis: Identify which parameters most affect policy effectiveness, focusing implementation efforts on critical factors.

Scenario Planning: Test policies under different future scenarios—slow versus fast takeoff, unipolar versus multipolar development, cooperative versus adversarial dynamics.

Confidence Bounds: Rather than point estimates, compute ranges of possible effects accounting for parameter uncertainty.

3.8 Interactive Visualization Design

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

3.8.1 Visual Encoding Strategy

The system uses multiple visual channels:

Color: Probability magnitude (green=high, red=low) provides immediate visual indication of likelihood, leveraging intuitive color associations.

Borders: Node type (blue=root, purple=intermediate, magenta=effect) helps users understand causal flow through the network structure.

Size: Centrality in network (larger=more influential) draws attention to critical nodes that affect many other variables.

Layout: Force-directed positioning reveals clusters of related variables, helping users identify cohesive sub-arguments within larger models.

3.8.2 Progressive Disclosure

Information appears at appropriate levels:

1. **Overview:** Network structure and color coding provide immediate understanding of key relationships and probability distributions.
2. **Hover:** Node description and prior probability offer additional context without cluttering the main view.
3. **Click:** Full probability tables and details enable deep investigation of specific variables and their relationships.
4. **Interaction:** Drag to rearrange, zoom to explore—users can customize views for their specific interests and questions.

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

3.8.3 User Interface Elements

Effective interface design would incorporate:

Physics Controls: Allow users to adjust layout dynamics, finding arrangements that best reveal patterns of interest.

Filter Options: Enable focusing on specific node types or probability ranges, reducing complexity for targeted analysis.

Export Functions: Support saving visualizations and data in formats suitable for reports, presentations, and further analysis.

Comparison Mode: Facilitate side-by-side viewing of different models or the same model under different policy interventions.

These features should emerge from iterative design with actual users—researchers needing detailed analysis, policymakers seeking key insights, and public stakeholders requiring accessible overviews.

3.9 Integration with Prediction Markets

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

3.9.1 Design for Integration

The system anticipates market connections through:

API Specifications: Standardized interfaces for major forecasting platforms like Metaculus, Good Judgment Open, and Manifold Markets.

Semantic Matching: Algorithms to connect model variables with related forecast questions, handling differences in phrasing and scope.

Aggregation Methods: Principled approaches for combining multiple forecast sources, accounting for track records and expertise.

Update Scheduling: Efficient caching and refresh strategies to balance currency with computational cost.

3.9.2 Challenges and Opportunities

Key integration challenges:

Question Mapping: Model variables rarely match market questions exactly. “AI causes existential catastrophe” might map to multiple specific forecast questions about AI capabilities, deployment, and impacts.

Temporal Alignment: Markets forecast specific dates while models consider scenarios. Bridging these requires careful interpretation and uncertainty propagation.

Quality Variation: Market depth and participation vary significantly. Some questions attract expert forecasters while others rely on casual participants.

Despite challenges, even partial integration provides value through external validation of probability estimates and dynamic updating as new information emerges.

3.10 Computational Performance Analysis

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

3.10.1 Exact vs. Approximate Inference

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

Monte Carlo Methods: Sample from probability distributions to estimate queries. While approximate, these methods scale to arbitrary network sizes.

Variational Inference: Optimize simpler distributions to approximate true posteriors. These methods trade some accuracy for guaranteed convergence.

Belief Propagation: Pass messages between nodes to converge on beliefs. Particularly effective for tree-structured or sparse networks.

The system automatically selects appropriate methods based on network properties.

3.10.2 Scaling Strategies

For very large networks, several strategies enable practical analysis:

Hierarchical Decomposition: Break large networks into manageable sub-networks, compute locally, then integrate results.

Relevance Pruning: For specific queries, identify and focus on relevant subgraphs, ignoring distant unconnected nodes.

Caching Architecture: Store computed results for common queries, dramatically improving response time for interactive use.

Parallel Processing: Distribute computation across multiple cores or machines for large-scale analysis.

3.11 Results and Achievements

3.11.1 Extraction Quality Assessment

An ideal assessment methodology would systematically evaluate:

Coverage Metrics: What proportion of arguments in source texts are successfully captured in formal models?

Accuracy Metrics: How closely do automated extractions match expert consensus?

Robustness Metrics: How well does the system handle different writing styles, argument structures, and domains?

Utility Metrics: Do the extracted models enable meaningful analysis and decision support?

Preliminary applications suggest the approach achieves practical utility while highlighting areas for improvement, particularly in handling implicit reasoning and converting qualitative uncertainty expressions.

3.11.2 Computational Performance

Performance analysis would examine:

Scaling Characteristics: How processing time grows with network size and complexity.

Bottleneck Identification: Whether limitations arise from extraction, inference, or visualization.

Optimization Opportunities: Where algorithmic improvements or engineering enhancements could improve performance.

Resource Requirements: Memory, processing, and storage needs for realistic applications.

The modular architecture enables targeted optimization of bottleneck components while maintaining system coherence.

3.11.3 Policy Impact Evaluation

A comprehensive evaluation framework would assess:

Intervention Modeling: How effectively can policies be represented as network modifications?

Robustness Testing: Do policy recommendations remain stable across model variations?

Comparative Analysis: How do different policy options compare in effectiveness?

Implementation Guidance: Does the analysis provide actionable insights for policymakers?

The ability to formally model policy interventions and trace their effects through complex causal networks represents a significant advance in systematic governance analysis.

3.12 Summary of Technical Contributions

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
- **Scalable implementation** handling realistic network sizes

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

The implementation shows that formal modeling of AI risk arguments is not only theoretically possible but practically achievable. While challenges remain—particularly in handling implicit reasoning and diverse uncertainty expressions—the system provides a foundation for enhanced coordination in AI governance.

4. Discussion: Implications and Limitations

Chapter Overview

Grade Weight: 10% | **Target Length:** ~14% of text (~4,200 words)

Requirements: Discusses objections, provides convincing replies, extends beyond course materials

4.1 Technical Limitations and Responses

4.1.1 Objection 1: Extraction Quality Boundaries

Critic: “Complex implicit reasoning chains resist formalization; automated extraction will systematically miss nuanced arguments and subtle conditional relationships that human experts would identify.”

Response: This concern has merit—extraction does face inherent limitations. However, the empirical results tell a more nuanced story. The two-stage extraction process, while imperfect, captures sufficient structure for practical use while maintaining transparency about its limitations.

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 partially 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.

4.1.2 Objection 2: False Precision in Uncertainty

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 a range 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.”

4.1.3 Objection 3: 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.⁶

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

⁶Copulas provide a mathematically elegant way to separate marginal behavior from dependence structure

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.

4.2 Conceptual and Methodological Concerns

4.2.1 Objection 4: 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 Bayes-Down 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.

4.2.2 Objection 5: Oversimplification of Complex Systems

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.

4.2.3 Objection 6: Idiosyncratic Implementation and Modeling Choices

Critic: “The specific choices made in AMTAIR’s implementation—from prompt design to parsing algorithms to visualization strategies—seem arbitrary. Different teams might make entirely different choices, leading to incompatible results. How can we trust conclusions that depend so heavily on implementation details?”

Response: This concern about implementation dependency is valid and deserves careful consideration. However, several factors mitigate this issue:

Convergent Design Principles: While specific implementations vary, fundamental design principles tend to converge. The two-stage extraction process (structure then probability) emerges naturally from how humans parse arguments. The use of intermediate representations follows established practice in computational linguistics. These aren’t arbitrary choices but responses to inherent challenges.

Empirical Validation: The “correctness” of implementation choices isn’t philosophical but empirical. If different reasonable implementations extract similar structures and lead to similar policy conclusions, this demonstrates robustness. If they diverge dramatically, this reveals genuine ambiguity in source materials—itsself valuable information.

Transparent Methodology: By documenting all implementation choices and making code

open source, AMTAIR enables replication and variation. Other teams can modify specific components while preserving overall architecture, testing which choices matter.

Convergence at Higher Levels: Even if implementations differ in details, they may converge at levels that matter for coordination. If two systems extract slightly different network structures but reach similar conclusions about policy robustness, the implementation differences don't undermine the approach's value.

Community Standards: As the field matures, community standards will likely emerge—not enforcing uniformity but establishing interoperability. This parallels development in other technical fields where multiple implementations coexist within shared frameworks.

The deeper insight is that implementation choices encode theoretical commitments. By making these explicit and variable, AMTAIR turns a bug into a feature—we can systematically explore how different assumptions affect conclusions, enhancing rather than undermining epistemic security.

4.3 Red-Teaming Results

To identify failure modes, systematic adversarial testing of the AMTAIR system would be essential.

4.3.1 Adversarial Extraction Attempts

A comprehensive red-teaming approach would test the system with:

Contradictory Arguments: Texts containing logically inconsistent claims or probability estimates. The system should flag contradictions rather than silently reconciling them.

Circular Reasoning: Arguments with circular dependencies that violate DAG requirements. Proper validation should detect and report such structural issues.

Ambiguous Language: Texts using extremely vague or metaphorical language. The system should acknowledge extraction uncertainty rather than forcing precise interpretations.

Deceptive Framings: Arguments crafted to imply false causal relationships. This tests whether the system merely extracts surface claims or requires deeper coherence.

Adversarial Prompts: Inputs designed to trigger known LLM failure modes. This ensures robustness against prompt injection and manipulation attempts.

Each failure mode discovered would inform system improvements and user guidance.

4.3.2 Robustness Findings

Theoretical analysis suggests key vulnerabilities:

Anchoring Effects: Language models may over-weight information presented early in documents, potentially biasing extraction toward initial framings.

Authority Sensitivity: Extraction might be influenced by explicit credibility signals in text, potentially giving undue weight to claimed expertise.

Complexity Limits: Performance likely degrades with very large argument structures, requiring hierarchical decomposition strategies.

Context Windows: Long-range dependencies exceeding model context windows could be missed, fragmenting cohesive arguments.

Understanding these limitations enables appropriate use—leveraging strengths while compensating for weaknesses through human oversight and validation.

4.3.3 Implications for Deployment

These considerations 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

4.4 Enhancing Epistemic Security

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

4.4.1 Making Models Inspectable

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

4.4.2 Revealing Convergence and Divergence

Theoretical analysis suggests formal comparison would reveal:

Structural Patterns: Experts likely share more agreement about causal structures than probability values, suggesting common understanding of mechanisms despite quantitative disagreement.

Crux Identification: Formal models make explicit which specific disagreements drive different conclusions, focusing discussion on genuinely critical differences.

Hidden Agreements: Apparently conflicting positions might share substantial common ground obscured by different terminology or emphasis.

Uncertainty Clustering: Areas of high uncertainty likely correlate across models, revealing where additional research would most reduce disagreement.

These patterns remain invisible in natural language debates but become analyzable through formalization.

4.4.3 Improving 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

While empirical validation remains future work, theoretical considerations suggest these mechanisms could substantially improve coordination quality. By providing shared representations and systematic methods for managing disagreement, formal models create infrastructure for collective intelligence that transcends individual limitations.

4.5 Scaling Challenges and Opportunities

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

4.5.1 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

4.5.2 Social and Institutional 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

4.5.3 Opportunities for Impact

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.

4.6 Integration with Governance Frameworks

AMTAIR complements rather than replaces existing governance approaches.

4.6.1 Standards Development

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.

4.6.2 Regulatory Design

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.

The extensive literature on corporate governance and liability frameworks **cuomo2016 demirag2000 devilliers2021 divito2022 kaur2024 list2011 solomon2020** provides theoretical grounding for understanding how regulatory interventions shape organizational behavior. AMTAIR could formalize these relationships in the specific context of AI development, making explicit how different liability regimes might incentivize or discourage safety investments.

4.6.3 International Coordination

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.

4.6.4 Organizational Decision-Making

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.

4.7 Future Research Directions

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

4.7.1 Technical Enhancements

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

4.7.2 Methodological Extensions

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

Recent advances in causal structure learning from both text and data [babakov2025](#) [ban2023](#) [bethard2007](#) [chen2023](#) [heinze-deml2018](#) [squires2023](#) [yang2022](#) suggest promising directions for enhancing AMTAIR’s extraction capabilities. The theoretical foundations from [duhem1954](#) and [meyer2022b](#) on the philosophy of science and knowledge structures provide epistemological grounding for these methodological extensions.

4.7.3 Application Domains

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

4.7.4 Ecosystem Development

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.

4.8 Known Unknowns and Deep Uncertainties

While AMTAIR enhances reasoning under uncertainty, fundamental limitations remain regarding truly novel developments that might fall outside existing conceptual frameworks.

4.8.1 Categories of Deep Uncertainty

Novel Capabilities: Future AI developments may operate according to principles outside current scientific understanding. No amount of careful modeling can anticipate fundamental paradigm shifts in what intelligence can accomplish.

Emergent Behaviors: Complex system properties that resist prediction from component analysis may dominate outcomes. The interaction between advanced AI systems and human society could produce wholly unexpected dynamics.

Strategic Interactions: Game-theoretic dynamics with superhuman AI systems exceed human modeling capacity. We cannot reliably predict how entities smarter than us will behave strategically.

Social Transformation: Unprecedented social and economic changes may invalidate current institutional assumptions. Our models assume continuity in basic social structures that AI might fundamentally alter.

4.8.2 Adaptation Strategies for Deep Uncertainty

Rather than pretending to model the unmodelable, AMTAIR incorporates several strategies:

Model Architecture Flexibility: The modular structure enables rapid incorporation of new variables as novel factors become apparent. When surprises occur, models can be updated rather than discarded.

Explicit Uncertainty Tracking: Confidence levels for each model component make clear where knowledge is solid versus speculative. This prevents false confidence in highly uncertain domains.

Scenario Branching: Multiple model variants capture different assumptions about fundamental uncertainties. Rather than committing to one worldview, the system maintains portfolios of possibilities.

Update Mechanisms: Integration with prediction markets and expert assessment enables rapid model revision as new information emerges. Models evolve rather than remaining static.

4.8.3 Robust Decision-Making Principles

Given deep uncertainty, certain decision principles become paramount:

Option Value Preservation: Policies should maintain flexibility for future course corrections rather than locking in irreversible choices based on current models.

Portfolio Diversification: Multiple approaches hedging across different uncertainty sources provide robustness against model error.

Early Warning Systems: Monitoring for developments that would invalidate current models enables rapid response when assumptions break down.

Adaptive Governance: Institutional mechanisms must enable rapid response to new information rather than rigid adherence to plans based on outdated models.

The goal is not to eliminate uncertainty but to make good decisions despite it. AMTAIR provides tools for systematic reasoning about what we do know while maintaining appropriate humility about what we don't and can't know.

4.9 Summary of Implications

The discussion reveals both the promise and limitations of computational approaches to AI governance coordination:

Technical Feasibility: Despite imperfections, automated extraction and formal modeling prove practically viable for complex AI risk arguments.

Epistemic Value: Making implicit models explicit, enabling systematic comparison, and supporting evidence integration enhance collective reasoning.

Practical Limitations: Extraction boundaries, false precision risks, and implementation dependencies require careful management.

Integration Potential: The approach complements rather than replaces existing governance frameworks, adding rigor without sacrificing flexibility.

Future Development: Technical enhancements, methodological extensions, and ecosystem growth could amplify impact.

Deep Uncertainty: Fundamental limits on predicting novel developments require maintaining humility and adaptability.

These findings suggest AMTAIR represents a valuable addition to the AI governance toolkit—not a panacea but a meaningful enhancement to our collective capacity for navigating unprecedented challenges.

5. Conclusion: Toward Coordinated AI Governance

Chapter Overview

Grade Weight: 10% | **Target Length:** ~14% of text (~4,200 words)

Requirements: Summarizes thesis and argument, outlines implications, notes limitations, points to future research

5.1 Summary of 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.

5.1.1 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.

5.1.2 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, allows for unprecedented transparency and explainability of the entire process, 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.

5.1.3 Technical Achievements

Working Implementation: AMTAIR demonstrates end-to-end feasibility from document ingestion through interactive visualization. The system successfully processes complex arguments like Carlsmith’s power-seeking AI model, extracting hierarchical structures and probability information.

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.

5.1.4 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: Theoretical analysis suggests that formal comparison would reveal structural agreements across different expert worldviews even when probability estimates diverge—providing foundations for coordination.

Intervention Impacts: Policy evaluation capabilities demonstrate how formal models enable rigorous assessment of governance options. The ability to trace intervention effects through complex causal networks validates the practical value of formalization.

5.2 Limitations and Honest Assessment

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

5.2.1 Technical Constraints

Extraction Boundaries: The system struggles with implicit assumptions, complex conditionals, and ambiguous quantifiers. 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 require approximations that may affect accuracy. As models grow to represent richer phenomena, computational constraints increasingly bind.

5.2.2 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 causal models struggle to capture fully.

5.2.3 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 or political rhetoric remains challenging.

5.3 Implications for AI Governance

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

5.3.1 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

5.3.2 Medium-Term Transformation

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.

5.3.3 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.

5.4 Recommendations for Stakeholders

Different communities can take concrete steps to realize these benefits:

5.4.1 For Researchers

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

5.4.2 For Policymakers

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

5.4.3 For Technologists

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

Bibliography

[Existential_Risk]: Increase in existential risks for humanity. {"instantiations": [TRUE", "FALSE"]}

- [Unaligned_AGI_Risk]: Unaligned artificial general intelligence causes existential risk. {"instantiations": [TRUE", "FALSE"]}
 - [State-State_Relations]
- [Near_term_AI]: Even if not unaligned AGI, near term AI can act as intermediate risk factor. {"instantiations": [TRUE", "FALSE"]}
 - [State-State_Relations]: AI arms race dynamic inhibits international coordination, diverting resources from other pressing issues {"instantiations": [TRUE", "FALSE"]}
 - * [Cybersecurity]: Probably enhances Cyber-Attack-Offense, may intensify cyber warfare. {"instantiations": [TRUE", "FALSE"]}
 - [State-Cooperation_Relations]: Cooperations have a lot of power and might have misaligned goals with society {"instantiations": [TRUE", "FALSE"]}
 - [Stable_Repressive_Regime]: More repressive instruments, possibility of stable repressive regime. {"instantiations": [TRUE", "FALSE"]}
 - * [State-Citizen_Relations]: AI helps regime monitor citizens {"instantiations": [TRUE", "FALSE"]}
 - [Compromised_Political_Decision_Making]: AI can compromise political decision making. {"instantiations": [TRUE", "FALSE"]}
 - * [Social_media_and_Recommender_Systems]: Influence of AI in social media on public opinion. {"instantiations": [TRUE", "FALSE"]}
- [Nuclear]: Probability that nuclear conflict escalates to end civilisation. {"instantiations": [TRUE", "FALSE"]}
 - [Compromised_Political_Decision_Making]
- [Biological]: Probability that a natural or engineered pandemic poses existential risks. {"instantiations": [TRUE", "FALSE"]}
 - [Compromised_Political_Decision_Making]
 - [Social_media_and_Recommender_Systems]
- [Natural]: Non-human caused existential risks, seem unrelated with AI. {"instantiations": [TRUE", "FALSE"]}
- [Environmental]: Probability of climate catastrophe. {"instantiations": [TRUE", "FALSE"]}

-
- [Compromised_Political_Decision_Making]
 - [AI_resource_consumption]: Current AI models consume large amounts of energy having environmental impacts. {"instantiations": ["TRUE", "FALSE"]}
 - [Social_media_and_Recommender_Systems]



UNIVERSITÄT
BAYREUTH

– P&E Master's Programme –
Chair of Philosophy, Computer
Science & Artificial Intelligence

Affidavit

Declaration of Academic Honesty

Hereby, I attest that I have composed and written the presented thesis

Automating the Modelling of Transformative Artificial Intelligence Risks

independently on my own, without the use of other than the stated aids and without any other resources than the ones indicated. All thoughts taken directly or indirectly from external sources are properly denoted as such.

This paper has neither been previously submitted in the same or a similar form to another authority nor has it been published yet.

BAYREUTH on the
May 26, 2025

VALENTIN MEYER