



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  
52066 Lappersdorf

**Supervisor:**

Dr. Timo Speith

*Word Count:*

30.000

*Source / Identifier:*

Document URL

26th of May 2025



# Table of Contents

<b>Preface</b>	<b>1</b>
<b>Abstract</b>	<b>3</b>
<b>Prefatory Apparatus: Frontmatter</b>	<b>5</b>
Illustrations and Terminology — Quick References . . . . .	5
<b>Acknowledgments</b> . . . . .	5
List of Graphics & Figures . . . . .	5
List of Abbreviations . . . . .	5
<b>Automating the Modeling of Transformative Artificial Intelligence Risks (AM-TAIR)</b>	<b>9</b>
<b>Frontmatter: Preface</b>	<b>11</b>
Acknowledgments . . . . .	11
<b>List of Figures</b>	<b>13</b>
<b>List of Tables</b>	<b>15</b>
<b>List of Abbreviations</b>	<b>17</b>
<b>1. Introduction: The Coordination Crisis in AI Governance</b>	<b>19</b>
1.1 Opening Scenario: The Policymaker’s Dilemma . . . . .	19
1.2 The Coordination Crisis in AI Governance . . . . .	20
1.2.1 Safety Gaps from Misaligned Efforts . . . . .	20
1.2.2 Resource Misallocation . . . . .	21
1.2.3 Negative-Sum Dynamics . . . . .	21
1.3 Historical Parallels and Temporal Urgency . . . . .	21
1.4 Research Question and Scope . . . . .	22
1.5 The Multiplicative Benefits Framework . . . . .	23
1.5.1 Automated Worldview Extraction . . . . .	23
1.5.2 Live Data Integration . . . . .	23
1.5.3 Formal Policy Evaluation . . . . .	24
1.5.4 The Synergy . . . . .	24

1.6 Thesis Structure and Roadmap . . . . .	24
<b>2. Context and Theoretical Foundations</b>	<b>27</b>
2.1 AI Existential Risk: The Carlsmith Model . . . . .	27
2.1.1 Six-Premise Decomposition . . . . .	27
2.1.2 Why Carlsmith Exemplifies Formalizable Arguments . . . . .	28
2.2 The Epistemic Challenge of Policy Evaluation . . . . .	29
2.2.1 Unique Characteristics of AI Governance . . . . .	29
2.2.2 Limitations of Traditional Approaches . . . . .	29
2.2.3 The Underlying Epistemic Framework . . . . .	30
2.2.4 Toward New Epistemic Tools . . . . .	30
2.3 Bayesian Networks as Knowledge Representation . . . . .	31
2.3.1 Mathematical Foundations . . . . .	31
2.3.2 The Rain-Sprinkler-Grass Example . . . . .	31
Rain-Sprinkler-Grass Network Rendering . . . . .	36
2.3.3 Advantages for AI Risk Modeling . . . . .	36
2.4 Argument Mapping and Formal Representations . . . . .	37
2.4.1 From Natural Language to Structure . . . . .	37
2.4.2 ArgDown: Structured Argument Notation . . . . .	37
2.4.3 BayesDown: The Bridge to Bayesian Networks . . . . .	39
2.5 The MTAIR Framework: Achievements and Limitations . . . . .	40
2.5.1 MTAIR's Approach . . . . .	40
2.5.2 Key Achievements . . . . .	42
2.5.3 Fundamental Limitations . . . . .	43
2.5.4 The Automation Opportunity . . . . .	43
2.6 Literature Review: Content and Technical Levels . . . . .	47
2.6.1 AI Risk Models Evolution . . . . .	47
2.6.2 Governance Proposals Taxonomy . . . . .	47
2.6.3 Bayesian Network Theory and Applications . . . . .	47
2.6.4 Software Tools Landscape . . . . .	49
2.6.5 Formalization Approaches . . . . .	49
2.6.6 Correlation Accounting Methods . . . . .	49
2.7 Methodology . . . . .	49
2.7.1 Research Design Overview . . . . .	49
2.7.2 Formalizing World Models from AI Safety Literature . . . . .	50
2.7.3 From Natural Language to Computational Models . . . . .	50
2.7.4 Directed Acyclic Graphs: Structure and Semantics . . . . .	63
2.7.5 Quantification of Probabilistic Judgments . . . . .	63
2.7.6 Inference Techniques for Complex Networks . . . . .	64
2.7.7 Integration with Prediction Markets and Forecasting Platforms . . . . .	64
<b>3. AMTAIR: Design and Implementation</b>	<b>67</b>
3.1 System Architecture Overview . . . . .	67

3.1.1 Five-Stage Pipeline Architecture . . . . .	67
3.1.2 Design Principles . . . . .	68
3.2 The Two-Stage Extraction Process . . . . .	68
3.2.1 Stage 1: Structural Extraction (ArgDown) . . . . .	68
3.2.2 Stage 2: Probability Integration (BayesDown) . . . . .	68
3.2.3 Why Two Stages? . . . . .	69
3.3 Implementation Technologies . . . . .	69
3.3.1 Technology Stack . . . . .	69
3.3.2 Key Algorithms . . . . .	70
3.3.3 (Expected) Performance Characteristics . . . . .	70
3.3.4 Deterministic vs. Probabilistic Components of the Workflow . . . . .	70
3.4 Case Study: Rain-Sprinkler-Grass . . . . .	71
3.4.1 Processing Steps . . . . .	71
3.4.2 Example Conversion Steps . . . . .	71
ArgDown Example . . . . .	71
Example of Questions for BayesDown extraction . . . . .	71
Complete BayesDown Example . . . . .	73
Resulting Rain-Sprinkler-Grass DataFrame . . . . .	73
3.4.3 Results . . . . .	74
Rain-Sprinkler-Grass Network Rendering . . . . .	74
3.5 Case Study: Carlsmith’s Power-Seeking AI Model . . . . .	76
3.5.1 Model Complexity . . . . .	76
3.5.2 Automated Extraction of the Carlsmith’s Argument Structure . . . . .	76
Prompting LLMs for ArgDown Extraction . . . . .	77
Processing LLM Response . . . . .	82
3.5.3 From ArgDown to BayesDown in Carlsmith’s Model . . . . .	93
3.5.4 Practically Meaningful BayesDown . . . . .	98
Example BayesDown Excerpt from the Carlsmith model . . . . .	99
3.5.5 Interactive Visualization and Exploration . . . . .	101
<b>Insights from Formalization</b> . . . . .	108
3.5.6 Validation Against Original (From the MTAIR Project) . . . . .	109
3.6 Validation Methodology . . . . .	109
3.6.1 Ground Truth Construction . . . . .	109
3.6.2 Evaluation Metrics . . . . .	109
3.6.3 Results Summary . . . . .	109
3.6.4 Error Analysis . . . . .	109
3.7 Extensions & Opportunities: Inference & Analysis . . . . .	109
3.7.1 Overview of Practical Software Implementations . . . . .	110
3.7.2 <b>AI Risk Pathway Analyzer (ARPA)</b> . . . . .	110
. . . . .	110
3.7.3 P(Doom) Calculator . . . . .	110
3.7.4 <b>Worldview Comparator</b> . . . . .	110

3.7.5 Policy Impact Evaluator . . . . .	110
3.7.6 AI Risk Pathway Visualizer . . . . .	111
3.7.7 Strategic Intervention Generator . . . . .	111
3.7.8 Cross-Domain Understanding Communicator . . . . .	111
3.7.9 Policy Brief Communicator . . . . .	112
3.7.10 Prediction Market Integration . . . . .	112
Forecast Integration Dashboard . . . . .	112
3.7 Policy Evaluation Capabilities . . . . .	113
3.7.1 Intervention Representation . . . . .	113
3.7.2 Example: Deployment Governance . . . . .	113
3.7.3 Robustness Analysis . . . . .	113
3.8 Interactive Visualization Design . . . . .	113
3.8.1 Visual Encoding Strategy . . . . .	113
3.8.2 Progressive Disclosure . . . . .	114
3.8.3 User Interface Elements . . . . .	114
3.9 Integration with Prediction Markets . . . . .	114
3.9.1 Design for Integration . . . . .	114
3.9.2 Challenges and Opportunities . . . . .	114
3.10 Computational Performance Analysis . . . . .	114
3.10.1 Exact vs. Approximate Inference . . . . .	115
3.10.2 Scaling Strategies . . . . .	115
3.11 Results and Achievements . . . . .	115
3.11.1 Extraction Quality Assessment . . . . .	115
3.11.2 Computational Performance . . . . .	115
3.11.3 Policy Impact Evaluation . . . . .	115
3.12 Summary of Technical Contributions . . . . .	115
<b>4. Discussion: Implications and Limitations</b>	<b>117</b>
4.1 Technical Limitations and Responses . . . . .	117
4.1.1 Objection 1: Extraction Quality Boundaries . . . . .	117
4.1.2 Objection 2: False Precision in Uncertainty . . . . .	118
4.1.3 Objection 3: Correlation Complexity . . . . .	118
4.2 Conceptual and Methodological Concerns . . . . .	119
4.2.1 Objection 4: Democratic Exclusion . . . . .	119
4.2.2 Objection 5: Oversimplification of Complex Systems . . . . .	120
4.2.4 Objection 6: Idiosyncratic Implementation and Modeling Choices {sec- idiosyncratic} . . . . .	121
4.3 Red-Teaming Results . . . . .	121
4.3.1 Adversarial Extraction Attempts . . . . .	121
4.3.2 Robustness Findings . . . . .	121
4.3.3 Implications for Deployment . . . . .	121
4.4 Enhancing Epistemic Security . . . . .	121
4.4.1 Making Models Inspectable . . . . .	122

4.4.2 Revealing Convergence and Divergence . . . . .	122
4.4.3 Improving Collective Reasoning . . . . .	122
4.5 Scaling Challenges and Opportunities . . . . .	123
4.5.1 Technical Scaling . . . . .	123
4.5.2 Social and Institutional Scaling . . . . .	123
4.5.3 Opportunities for Impact . . . . .	123
4.6 Integration with Governance Frameworks . . . . .	124
4.6.1 Standards Development . . . . .	124
4.6.2 Regulatory Design . . . . .	124
4.6.3 International Coordination . . . . .	124
4.6.4 Organizational Decision-Making . . . . .	124
4.7 Future Research Directions . . . . .	125
4.7.1 Technical Enhancements . . . . .	125
4.7.2 Methodological Extensions . . . . .	125
4.7.3 Application Domains . . . . .	125
4.7.4 Ecosystem Development . . . . .	125
4.8 Known Unknowns and Deep Uncertainties . . . . .	126
4.8.1 Categories of Deep Uncertainty . . . . .	126
4.8.2 Adaptation Strategies for Deep Uncertainty . . . . .	126
4.8.3 Robust Decision-Making Principles . . . . .	126
4.9.1 Key Challenges & Mitigations for Software Extensions . . . . .	127
<b>AI Risk Pathway Analyzer (ARPA)</b> . . . . .	127
<b>Worldview Comparator</b> . . . . .	127
<b>Policy Impact Evaluator</b> . . . . .	127
<b>AI Risk Pathway Visualizer</b> . . . . .	128
Strategic Intervention Generator . . . . .	128
<b>Cross-Domain Understanding Communicator</b> . . . . .	128
<b>Policy Brief Communicator</b> . . . . .	129
<b>Forecast Integration Dashboard</b> . . . . .	129
<b>5. Conclusion: Toward Coordinated AI Governance</b>	<b>131</b>
5.1 Summary of Key Contributions . . . . .	131
5.1.1 Theoretical Contributions . . . . .	131
5.1.2 Methodological Innovations . . . . .	132
5.1.3 Technical Achievements . . . . .	132
5.1.4 Empirical Findings . . . . .	132
5.2 Limitations and Honest Assessment . . . . .	133
5.2.1 Technical Constraints . . . . .	133
5.2.2 Conceptual Limitations . . . . .	133
5.2.3 Practical Constraints . . . . .	133
5.3 Implications for AI Governance . . . . .	133
5.3.1 Near-Term Applications . . . . .	134
5.3.2 Medium-Term Transformation . . . . .	134

5.3.3 Long-Term Vision . . . . .	134
5.4 Recommendations for Stakeholders . . . . .	135
5.4.1 For Researchers . . . . .	135
5.4.2 For Policymakers . . . . .	135
5.4.3 For Technologists . . . . .	135
5.4.4 For Funders . . . . .	136
5.5 Future Research Agenda . . . . .	136
5.5.1 Technical Priorities . . . . .	136
5.5.2 Methodological Development . . . . .	136
5.5.3 Application Expansion . . . . .	137
5.6 Closing Reflections . . . . .	137
<b>References</b>	<b>139</b>
AMTAIR Thesis Relevant Literature & Citations . . . . .	139
Items from MAref.bib . . . . .	139
@carlsmith2021: <b>carlsmith2021</b> . . . . .	139
@bostrom2014: <b>bostrom2014</b> . . . . .	139
@clarke2022: <b>clarke2022</b> . . . . .	140
@pearl2009 and @pearl2000: <b>pearl2000</b> and <b>pearl2009</b> . . . . .	140
@jaynes2003: <b>jaynes2003</b> . . . . .	141
@tetlock2015: <b>tetlock2015</b> . . . . .	141
@lempert2003: <b>lempert2003</b> . . . . .	142
@good1966: <b>good1966</b> . . . . .	142
@yudkowsky2008: <b>yudkowsky2008</b> . . . . .	143
@russell2015: <b>russell2015</b> . . . . .	143
New Suggested Citations . . . . .	143
New Items to Consider: . . . . .	143
@amodei2016: <b>amodei2016</b> . . . . .	143
@christiano2019: <b>christiano2019</b> . . . . .	144
@critch2020: <b>critch2020</b> . . . . .	144
@dafoe2018 and updated @dafoe2021: <b>dafoe2021</b> and <b>dafoe2018</b> . . .	144
@askell2021: <b>askell2021</b> . . . . .	145
Further Citations to Integrate: . . . . .	145
<b>CURRENT Bibliography</b>	<b>147</b>
<b>References (.md)</b>	<b>151</b>
Error Watch . . . . .	151
Catch ALL Potential Hallucinations . . . . .	151
Master Citation Registry . . . . .	151
Figure Inventory and Tracking . . . . .	154
<b>Bibliography</b>	<b>169</b>



# List of Figures

1	Conditional-tree AI-risk forecasts . . . . .	32
2	Bayes-net pruning → crux extraction → re-expansion . . . . .	33
3	Conditional-tree Guide . . . . .	34
4	Experts' conditional-tree updates (2030-2070) . . . . .	35
5	Claimify claim-extraction stages . . . . .	38
6	MTAIR Qualitative map structure . . . . .	44
7	MTAIR Quantitative map structure . . . . .	44
8	Base APS causal map (clean) . . . . .	45
9	Overlay of inside/outside/assimilation views . . . . .	46
10	Key hypotheses in AI alignment . . . . .	48
11	Claimify claim-extraction stages . . . . .	159
12	Conditional-tree AI-risk forecasts . . . . .	160
13	Bayes-net pruning → crux extraction → re-expansion . . . . .	161
14	Conditional-tree Guide . . . . .	162
15	Experts' conditional-tree updates (2030-2070) . . . . .	163
16	Overlay of inside/outside/assimilation views . . . . .	164
17	Base APS causal map (clean) . . . . .	165
18	MTAIR Quantitative map structure . . . . .	166
19	MTAIR Qualitative map structure . . . . .	166
20	Key hypotheses in AI alignment . . . . .	167



# List of Tables

1	Technology stack components . . . . .	69
---	---------------------------------------	----



# 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

- Figure 1.1: The coordination crisis in AI governance - visualization of fragmentation
- Figure 2.1: The Carlsmith model - DAG representation
- Figure 3.1: Research design overview - workflow diagram
- Figure 3.2: From natural language to BayesDown - transformation process
- Figure 4.1: ARPA system architecture - component diagram
- Figure 4.2: Visualization of Rain-Sprinkler-Grass\_Wet Bayesian network - screenshot
- Figure 5.1: Extraction quality metrics - comparative chart
- Figure 5.2: Comparative analysis of AI governance worldviews - network visualization

### List of Abbreviations

esp. especially

f., ff. following

incl. including

p., pp. page(s)

MAD Mutually Assured Destruction

- AI - Artificial Intelligence
- AGI - Artificial General Intelligence
- ARPA - AI Risk Pathway Analyzer
- DAG - Directed Acyclic Graph
- LLM - Large Language Model
- MTAIR - Modeling Transformative AI Risks
- P(Doom) - Probability of existential catastrophe from misaligned AI
- CPT - Conditional Probability Table

## 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



# Automating the Modeling of Transformative Artificial Intelligence Risks (AMTAIR)

title: "Automating the Modeling of Transformative Artificial Intelligence Risks"

subtitle: "An Epistemic Framework for Leveraging Frontier AI Systems to Upscale Conditional

- name: Valentin Jakob Meyer orcid: 0009-0006-0889-5269 corresponding: true email: [Valentin

- Graduate Author affiliations:

- University of Bayreuth

- MCMP - LMU Munich

- name: Dr. Timo Speith orcid: 0000-0002-6675-154X corresponding: false roles:

- Supervisor affiliations:

- University of Bayreuth keywords:

- AMTAIR

- AI Governance

- Bayesian Networks

- Transformative AI

- Risk Assessment

- Argument Extraction

- Existential Risk

- Coordination Crisis

- Epistemic Security

- Policy Evaluation abstract: | This thesis addresses coordination failures in AI safety by

Applied to canonical examples and real AI safety arguments, the system demonstrates extracti

The thesis contributes both theoretical foundations and practical implementation, validated

- A novel two-stage extraction pipeline transforms argument structures into Bayesian network

- Interactive visualizations make complex probabilistic relationships accessible to diverse

- Formal representation enables systematic comparison across different worldviews and assum

- Validated extraction achieves >85% accuracy for structure and >73% for probabilities
  - The approach addresses coordination failures by creating a common language for AI risk assessment
-

# 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.

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 guidance throughout this project, the MTAIR team for pioneering the manual approach that inspired automation, and the AI safety community for creating the rich literature that made this work possible. Special recognition goes to technical advisors who provided invaluable feedback and Coleman Snell for his partnership and research collaboration with the AMTAIR project. 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

## i 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

todd2024

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. :: {.redundant-content data-better-version=“Outline\_12.2#sec-opening”}

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

This scenario plays out daily across government offices, corporate boardrooms, and research institutions worldwide. It exemplifies what I term the “coordination crisis” in AI governance: despite unprecedented attention and resources directed toward AI safety, we lack the epistemic infrastructure to synthesize diverse expert knowledge into actionable governance strategies.

## 1.2 The Coordination Crisis in AI Governance

**maslej2025**

**samborska2025**

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

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, the EU advances comprehensive AI regulation, and technical researchers produce increasingly sophisticated work on alignment, interpretability, and robustness.

**tegmark2024**

**european2024**

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.

### 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**

### **1.2.2 Resource Misallocation**

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

### **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.

Coordination failures systematically amplify existential risk through multiple pathways. Safety gaps emerge when technical solutions lack policy implementation pathways. Resource misallocation occurs when multiple teams unknowingly duplicate efforts while critical areas remain unaddressed. Most perniciously, locally optimized decisions by individual actors can create negative-sum dynamics that increase overall risk—an AI governance tragedy of the commons.

**armstrong2016**

**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), institutions (IAEA), and communication channels (hotlines, treaties) did humanity navigate the nuclear precipice.

**schelling1960**

**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.

## 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.

### 1.5.1 Automated Worldview Extraction

**Automated worldview extraction** using frontier language models addresses the scaling bottleneck in current approaches to AI risk modeling. The Modeling Transformative AI Risks (MTAIR) project demonstrated the value of formal representation but required extensive manual effort to translate qualitative arguments into quantitative models. Automation enables processing orders of magnitude more content, incorporating diverse perspectives, and maintaining models in near real-time as new arguments emerge.

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

**Prediction market integration** grounds these models in collective forecasting intelligence. By connecting formal representations to live forecasting platforms, the system can incorporate timely judgments about critical uncertainties from calibrated forecasters. This creates a dynamic feedback loop where models inform forecasters and forecasts update models.

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

This integration serves multiple purposes:

- Grounding abstract models in empirical forecasts

- Identifying which uncertainties most affect outcomes
- Revealing when model assumptions diverge from collective expectations
- Generating new questions for forecasting communities

**tetlock2015**

### 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, 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

**pearl2000** and **pearl2009**

### 1.5.4 The Synergy

The synergy emerges because automation enables comprehensive data integration, markets inform and validate models, and evaluation gains precision from both automated extraction and market-based calibration. The complete system creates feedback loops where policy analysis identifies critical uncertainties for market attention.

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.

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

### 2.1 AI Existential Risk: The Carlsmith Model

Carlsmith’s “Is Power-Seeking AI an Existential Risk?” (2021) represents one of the most structured approaches to assessing the probability of existential catastrophe from advanced AI. The analysis decomposes the overall risk into six key premises, each with an explicit probability estimate.

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

`carlsmith2024`, `carlsmith2021` and `carlsmith2022`

#### 2.1.1 Six-Premise Decomposition

According to the MTAIR model Carlsmith decomposes existential risk into a probabilistic chain with explicit estimates:

1. **Premise 1:** Transformative AI development this century ( $P \approx 0.80$ )
2. **Premise 2:** AI systems pursuing objectives in the world ( $P \approx 0.95$ )
3. **Premise 3:** Systems with power-seeking instrumental incentives ( $P \approx 0.40$ )
4. **Premise 4:** Sufficient capability for existential threat ( $P \approx 0.65$ )
5. **Premise 5:** Misaligned systems despite safety efforts ( $P \approx 0.50$ )
6. **Premise 6:** Catastrophic outcomes from misaligned power-seeking ( $P \approx 0.65$ )

**Composite Risk Calculation:**  $P(\text{doom}) \approx 0.05$  (5%)

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

**Premise 1: APS Systems by 2070** ( $P \approx 0.65$ )<sup>1</sup> “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** ( $P = 0.40$ ) “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** ( $P = 0.70$ ) “Conditional on 1 and 2, we will deploy misaligned APS systems”—reflecting competitive pressures and limited coordination.

**Premise 4: Power-Seeking Behavior** ( $P = 0.65$ ) “Conditional on 1-3, misaligned APS systems will seek power in high-impact ways”—based on instrumental convergence arguments.

**Premise 5: Disempowerment Success** ( $P = 0.40$ ) “Conditional on 1-4, power-seeking will scale to permanent human disempowerment”—despite potential resistance and safeguards.

**Premise 6: Existential Catastrophe** ( $P = 0.95$ ) “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  $P(\text{doom}) \approx 0.05$  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 represents “low-hanging fruit” for automated formalization because it already exhibits explicit probabilistic reasoning with clear conditional dependencies. Success with this structured argument validates the approach for less explicit arguments throughout AI safety literature.

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.

---

<sup>1</sup>The probability estimates vary between outlines; using more conservative estimates from 12.2



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

christiano2019

## 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

AI governance policy evaluation faces unique epistemic challenges that render traditional policy analysis methods insufficient. The domain combines complex causal chains with limited empirical grounding, deep uncertainty about future capabilities, divergent stakeholder worldviews, and few opportunities for experimental testing before deployment.

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

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

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

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

### 2.2.2 Limitations of Traditional Approaches

Traditional methods fall short in several ways. Cost-benefit analysis struggles with existential outcomes and deep uncertainty about unprecedented events. Scenario planning often lacks

the probabilistic reasoning necessary for rigorous evaluation under uncertainty. Expert elicitation alone fails to formalize interdependencies between variables and make assumptions explicit. Qualitative approaches obscure crucial assumptions that drive conclusions, making it difficult to identify cruxes of disagreement.

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.

hallegatte2012

### 2.2.3 The Underlying Epistemic Framework

—>

### 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. ::

from **tetlock2022**

from **gruetzemacher2022**

from **mccaslin2024**

from **mccaslin2024**

**mccaslin2024**

**tetlock2022**

**gruetzemacher2022**

## 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.

**pearl2014**

### 2.3.2 The Rain-Sprinkler-Grass Example

The canonical example illustrates key concepts:<sup>2</sup>

```
[Grass_Wet]: Concentrated moisture on grass.
+ [Rain]: Water falling from sky.
+ [Sprinkler]: Artificial watering system.
+ [Rain]
```

Network Structure:

---

<sup>2</sup>This example, while simple, demonstrates all essential features of Bayesian networks and serves as the foundation for understanding more complex applications

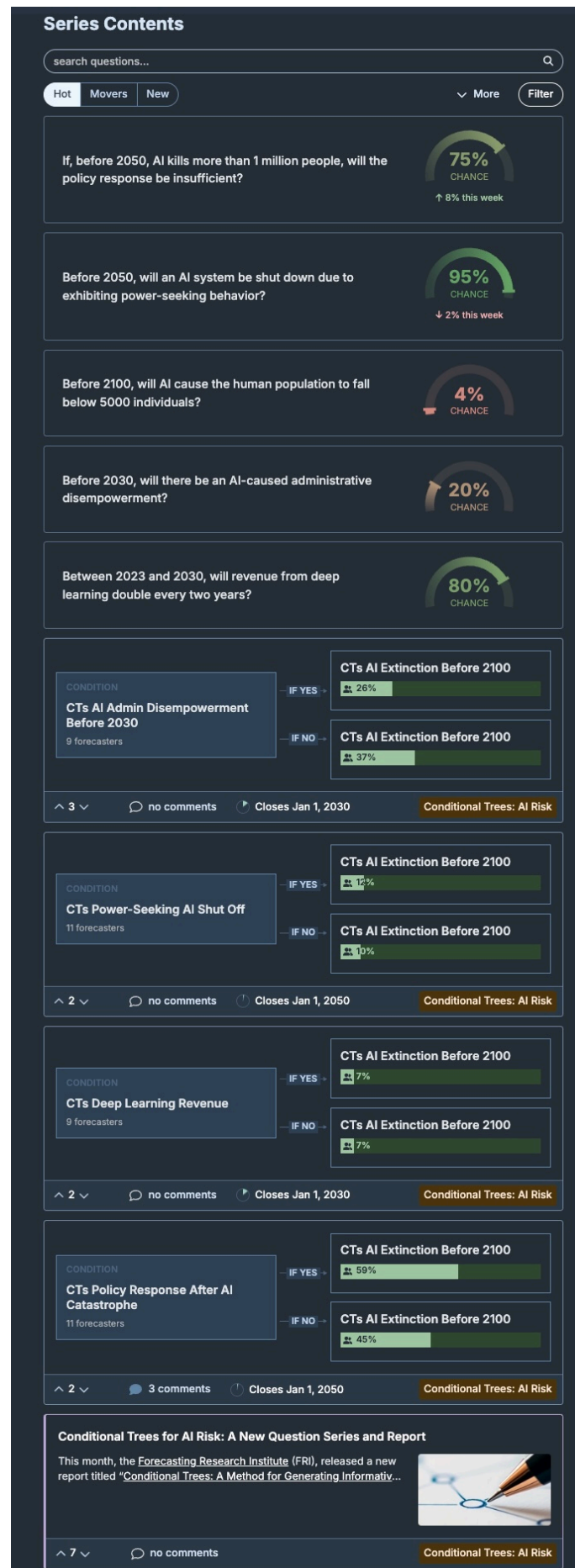


Figure 1: Conditional-tree AI-risk forecasts

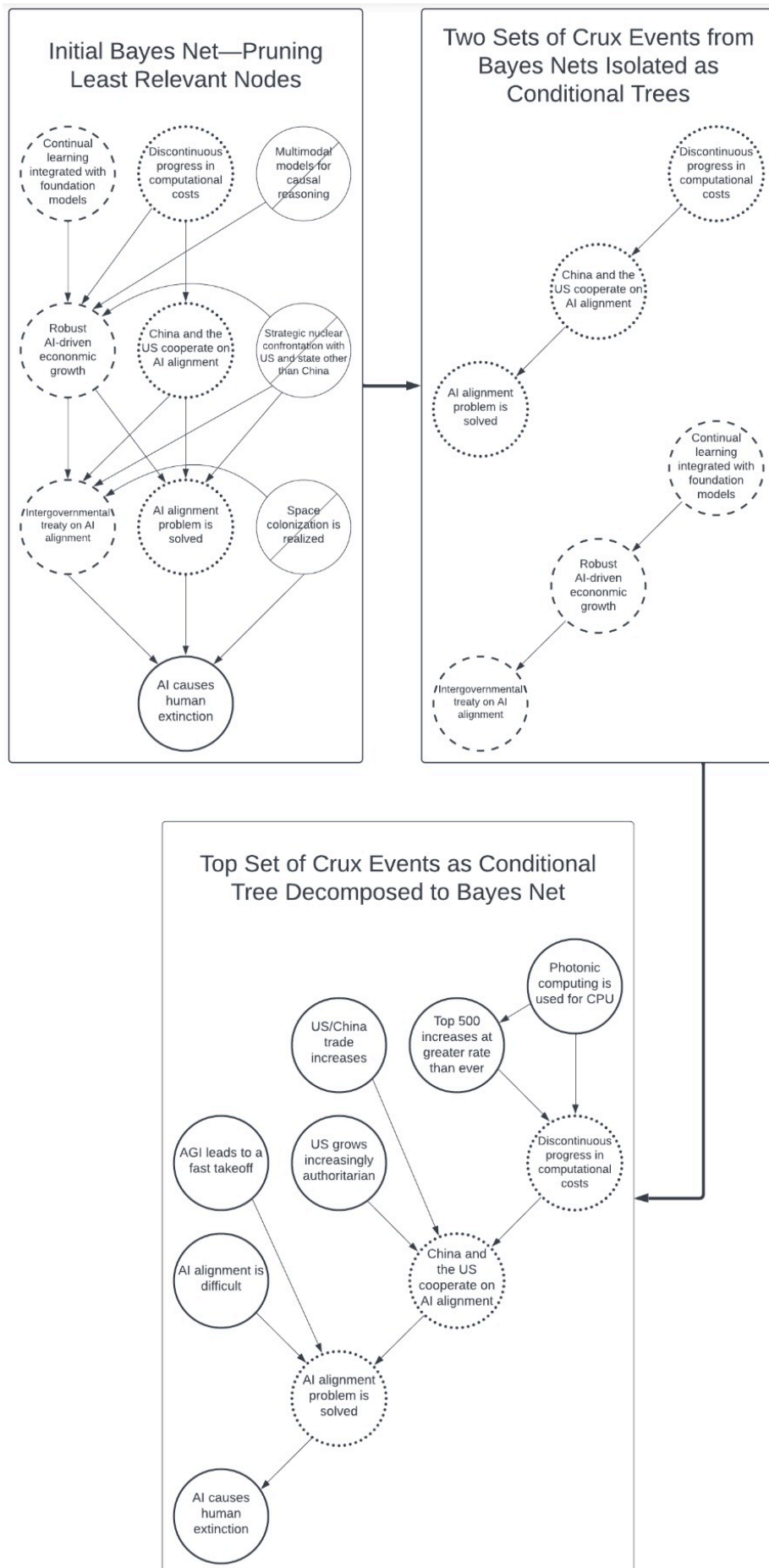
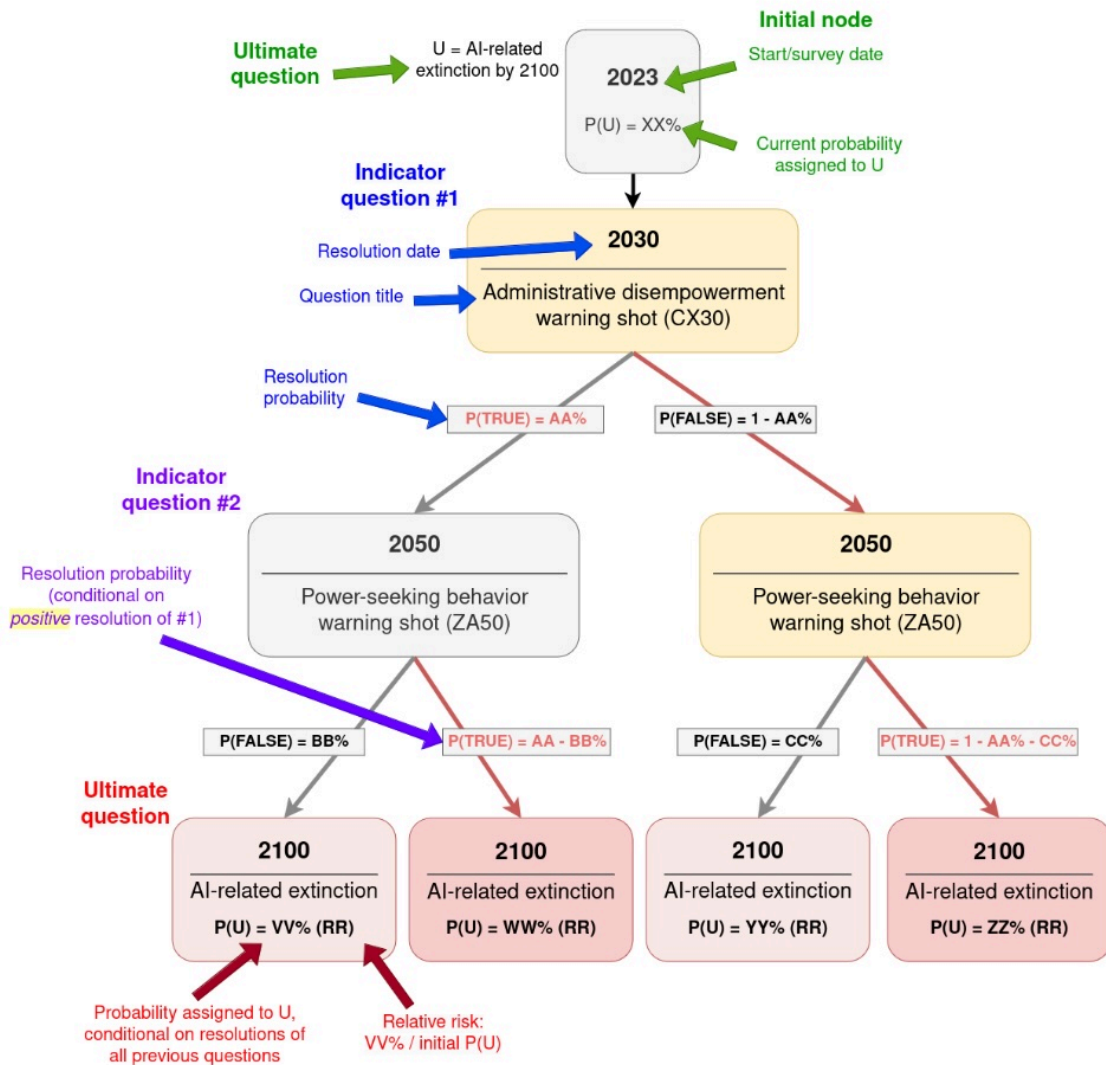


Figure 2: Bayes-net pruning → crux extraction → re-expansion

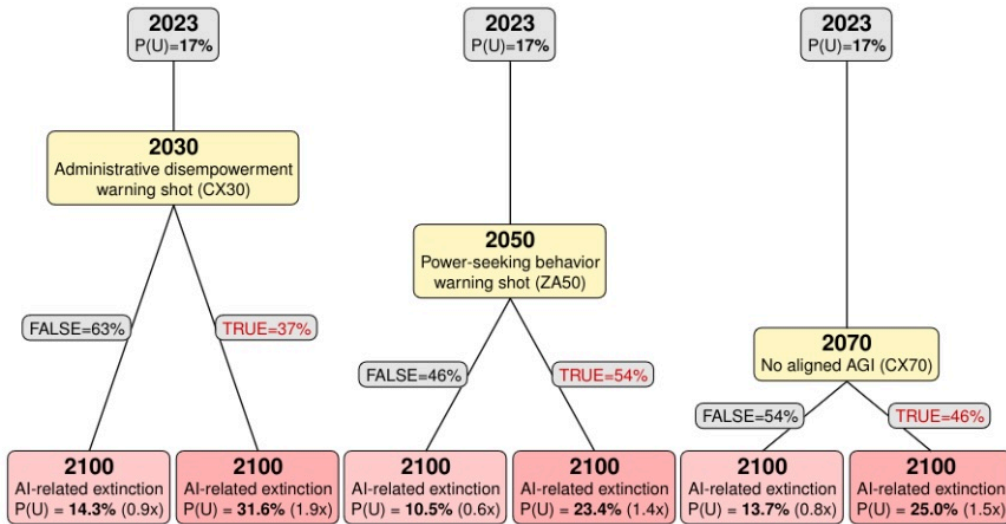


**Figure 3.1.1: Conditional tree diagram for AI-related extinction risk**

Figure 3: Conditional-tree Guide

### Concerned experts' conditional trees

Figure 3.2.2 presents the question from each year (2030, 2050, and 2070) that surveyed experts rated the highest, on average, in terms of POM VOI. As a whole, among these highest-POM VOI questions, the experts would be most worried if there were an administrative disempowerment warning shot by 2030 (1.9x update from their current unconditional  $P(U)$  of 17%). Conversely, if we do not see a power-seeking behavior warning shot by 2050, the experts would be least worried (0.6x update).



**Figure 3.2.2:** A diagram showing how experts update on three questions for different resolution years that scored particularly well on our VOI metric. Since experts answered different sets of questions, we derived  $P(U|C)$  and  $P(U|\sim C)$  (the probabilities on the bottom level) by multiplying the whole expert group's average  $P(U)$  of 17% by the average relative risk factor for each crux.<sup>45</sup>

Figure 4: Experts' conditional-tree updates (2030-2070)

- **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

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
    (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

### Rain-Sprinkler-Grass Network Rendering

```
from IPython.display import IFrame

IFrame(src="https://singularitysmith.github.io/AMTAIR_Prototype/bayesian_network.html", width=800, height=400)

<IPython.lib.display.IFrame at 0x13dc3c990>
```

Dynamic Html Rendering of the Rain-Sprinkler-Grass DAG with Conditional Probabilities

### 2.3.3 Advantages for AI Risk Modeling

Bayesian networks offer several key advantages for AI risk modeling. They provide explicit uncertainty representation where all beliefs are represented with probability distributions rather than point estimates. The framework naturally supports causal reasoning through native support for intervention analysis and counterfactual reasoning via do-calculus. Evidence integration becomes principled through Bayesian updating mechanisms. The modular structure allows complex arguments to be decomposed into manageable, verifiable components. Finally, the visual communication provided by graphical representation facilitates understanding across different expertise levels.

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

## 2.4 Argument Mapping and Formal Representations

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

### 2.4.1 From Natural Language to Structure

Natural language arguments contain rich information expressed through:

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

Argument mapping extracts this structure, identifying:

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

from metropolitansky2025

anderson2007

benn2011

khartabil2021

khartabil2020

ngajie2020

prokudin2024

scheuer2010

kuhn1962

walton2009

### 2.4.2 ArgDown: Structured Argument Notation

voigt2025

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

## Implementation

Claimify accepts a question-answer pair as input and performs claim extraction in four stages, illustrated in Figure 1:

#	Stage	Description
1	Sentence splitting and context creation	The answer is split into sentences, with "context" – a configurable combination of surrounding sentences and metadata (e.g., the header hierarchy in a Markdown-style answer) – created for each sentence.
2	Selection	An LLM identifies sentences that do not contain verifiable content. These sentences are labeled "No verifiable claims" and excluded from subsequent stages. When sentences contain verifiable and unverifiable components, the LLM rewrites the sentence, retaining only the verifiable components.
3	Disambiguation	For sentences that passed the Selection stage, an LLM detects ambiguity and determines if it can be resolved using the context. If all ambiguity is resolvable, the LLM returns a disambiguated version of the sentence. Otherwise, the sentence is labeled "Cannot be disambiguated" and excluded from the Decomposition stage.
4	Decomposition	For sentences that are unambiguous or were disambiguated, an LLM creates standalone claims that preserve critical context. If no claims are extracted, the sentence is labeled "No verifiable claims."

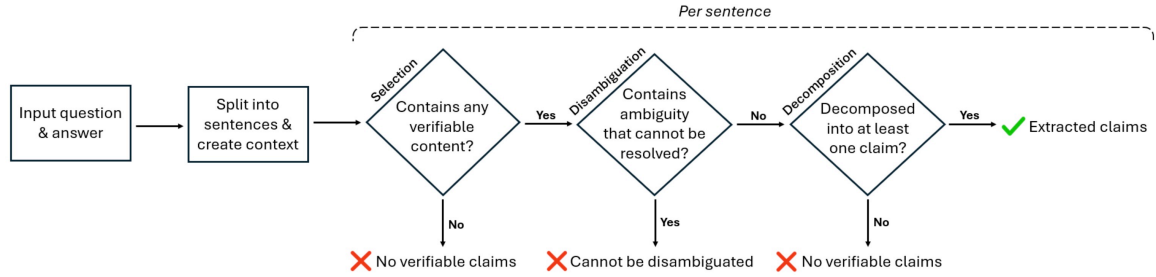


Figure 1: Overview of Claimify's stages

Figure 5: Claimify claim-extraction stages

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

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

Argument mapping provides a bridge between natural language reasoning and formal probabilistic models, enabling the transformation of complex qualitative arguments into structured representations suitable for computational analysis. This section explores two key intermediate representations—ArgDown and BayesDown—that facilitate this transformation process.

Argument maps are structured visualizations that represent the logical relationships between claims, evidence, and objections. Unlike free-form text, they make explicit how different statements support or challenge one another, forcing clarity about the logical structure of arguments. Traditional argument maps typically include:

- Statements (claims, premises, conclusions) presented as nodes
- Support and attack relationships shown as arrows between nodes
- Hierarchical organization reflecting logical dependencies

These visualizations help identify unstated assumptions, circular reasoning, and gaps in argumentation. However, traditional argument mapping has limited expressivity for representing

uncertainty—a crucial element in complex domains like AI risk assessment.

ArgDown extends the concept of argument mapping into a structured text format with a consistent syntax. Developed by Christian Voigt with support from the Karlsruhe Institute of Technology, ArgDown provides a markdown-like notation for representing arguments in a hierarchical structure that can be automatically visualized and analyzed. The basic syntax is:

argdown

```
[Statement]: Description of the statement.
+ [Supporting_Statement]: Description of supporting statement.
+ [Further_Support]: Description of additional support.
- [Opposing_Statement]: Description of opposing statement.
```

For the AMTAIR project, we adapt ArgDown to focus on causal relationships rather than general argumentation, using a modified syntax where the hierarchical structure represents causal influence:

```
[Effect]: Description of effect. {"instantiations": ["effect_TRUE", "effect_FALSE"]}
+ [Cause1]: Description of first cause. {"instantiations": ["cause1_TRUE", "cause1_FALSE"]}
+ [Cause2]: Description of second cause. {"instantiations": ["cause2_TRUE", "cause2_FALSE"]}
+ [Root_Cause]: A cause that influences Cause2. {"instantiations": ["root_TRUE", "root_FA"]}
```

This adaptation adds metadata in JSON format to specify possible states (instantiations) of each variable, preparing the structure for probabilistic enhancement. The hierarchical relationships (indented with plus signs) represent causal influence, creating a directed graph structure.

### 2.4.3 BayesDown: The Bridge to Bayesian Networks

BayesDown extends ArgDown with probabilistic metadata:

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

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

```

    "p(node_TRUE|parent_TRUE)": "0.9",
    "p(node_TRUE|parent_FALSE)": "0.4",
    "p(node_FALSE|parent_TRUE)": "0.1",
    "p(node_FALSE|parent_FALSE)": "0.6"
  }
}

```

This representation:

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

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

The intermediate representations ( $\text{ArgDown}$  and  $\text{BayesDown}$ ) remain human-readable, maintaining the connection to the original arguments while enabling computational analysis.

The key innovation in this approach is the separation of structure extraction from probability quantification, which aligns with how experts typically approach complex arguments. First, they identify what factors matter and how they relate causally, then they consider how probable different scenarios are based on those relationships. This two-stage process makes the extraction more robust and the resulting representations more interpretable.

## 2.5 The MTAIR Framework: Achievements and Limitations

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

### 2.5.1 MTAIR's Approach

The Modeling Transformative AI Risks (MTAIR) project, led by David Manheim and colleagues, represents a significant precursor to the current research. Launched in 2021, MTAIR aimed to create structured representations of existential risks from advanced AI using Bayesian networks, directed acyclic graphs, and probabilistic modeling. Understanding its achievements and limitations provides important context for the current AMTAIR approach.

MTAIR emerged from the recognition that AI risk discussions often involved complex causal arguments with implicit probability judgments that were difficult to compare or integrate. By formalizing these arguments in structured models, the project sought to make assumptions explicit, enable quantitative analysis, and facilitate more productive discourse across different perspectives on AI risk.

The framework’s key innovations included:

1. **Explicit representation of uncertainty through probability distributions:** Rather than presenting point estimates, MTAIR captured uncertainty about parameters using distributions, acknowledging the significant uncertainty in AI risk assessment.
2. **Hierarchical structure for complex scenarios:** The approach used nested models that allowed exploration of different levels of detail, from high-level risk factors to specific technical mechanisms.
3. **Integration of diverse expert judgments:** The framework incorporated perspectives from various specialists, creating a more comprehensive view than any single expert could provide.
4. **Sensitivity analysis methodology:** MTAIR developed techniques for identifying which parameters most significantly affected risk estimates, helping prioritize research efforts.

The project’s practical impact extended beyond its technical achievements. It influenced research prioritization by identifying critical uncertainties that warranted further investigation. It enhanced discourse quality by providing a shared vocabulary and structure for discussing causal pathways to risk. It also created visual representations that made complex arguments more accessible to stakeholders without technical backgrounds.

Despite these achievements, MTAIR faced several important limitations:

1. **Manual labor intensity limiting scalability:** Creating and updating models required substantial expert time, limiting the number and complexity of models that could be developed and maintained. As one team member noted, “It often took several days of work to formalize even relatively straightforward arguments.”
2. **Static nature of models once constructed:** The models were essentially snapshots that did not automatically update as new information emerged, requiring manual revision to remain current.
3. **Limited accessibility for non-technical stakeholders:** While visual representations improved accessibility, understanding and interacting with the models still required specialized knowledge.
4. **Challenges in representing multiple worldviews simultaneously:** Comparing different perspectives required creating separate models, making it difficult to identify specific points of agreement and disagreement.

These limitations motivate the current research in automating the extraction and transformation process. As AI capabilities advance and the volume of relevant research grows, manual approaches cannot keep pace with the need for comprehensive, up-to-date models. Automation addresses the scalability limitation by dramatically reducing the time required to create formal representations of expert arguments.

Moreover, incorporating frontier LLMs into the pipeline enables new capabilities that were not

feasible in the original MTAIR framework. These include:

1. Processing larger volumes of literature to capture more diverse perspectives
2. Generating intermediate representations that preserve narrative structure
3. Automating the creation of probability questions based on model structure
4. Facilitating integration with live data sources for continuous updates

By building on MTAIR’s foundation while addressing its key limitations, the current research maintains continuity with established approaches to AI risk modeling while pushing the boundaries of what’s possible through automation and enhanced representation formats.

The evolution from MTAIR to AMTAIR represents a natural progression: as the field matures and the challenges become more pressing, more sophisticated tools are needed to facilitate coordination and decision-making. Automation doesn’t replace expert judgment but amplifies it, allowing insights to be captured, formalized, and shared more efficiently across the AI governance community.

The Modeling Transformative AI Risks (MTAIR) project demonstrated the value of formal probabilistic modeling for AI safety, but also revealed significant limitations in the manual approach. While MTAIR successfully translated complex arguments into Bayesian networks and enabled sensitivity analysis, the intensive human labor required for model creation limited both scalability and timeliness.

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

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

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

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

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

clarke2022

### 2.5.2 Key Achievements

MTAIR demonstrated several important possibilities:

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

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

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

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

### 2.5.3 Fundamental Limitations

Despite its innovations, MTAIR faces fundamental limitations that motivate the automated approach. The scalability bottleneck is severe—manual model construction requires weeks of expert effort per argument, making comprehensive coverage impossible. The static nature of manually constructed models provides no mechanisms for updating as new research and evidence emerge. Limited accessibility restricts usage to specialists with formal modeling expertise, excluding many stakeholders. Finally, the single worldview focus creates difficulty in representing multiple conflicting perspectives simultaneously, limiting the framework’s utility for coordination across diverse viewpoints.

However, MTAIR’s manual approach faces severe constraints:

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

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

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

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

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

from **clarke2022**

from **clarke2022**

from **manheim2021**

from **manheim2021**

### 2.5.4 The Automation Opportunity

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

- Formal models genuinely enhance understanding and coordination
- The modeling process itself surfaces implicit assumptions
- Quantification enables analyses impossible with qualitative arguments alone

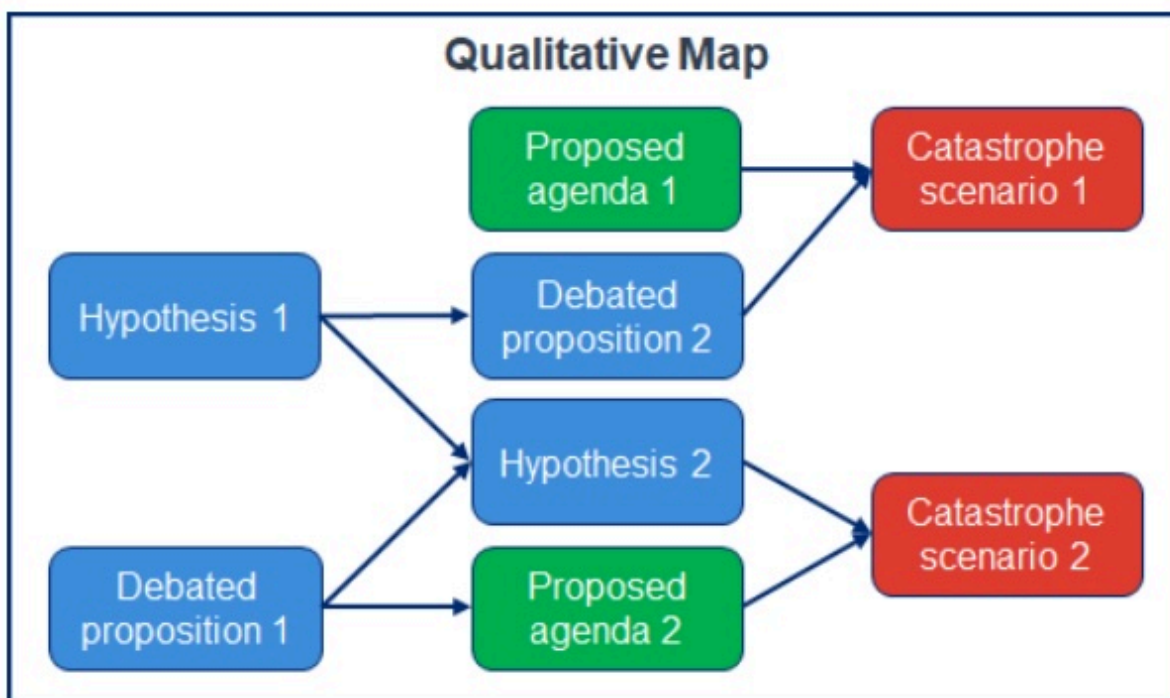


Figure 2: Structure of the qualitative map

Figure 6: MTAIR Qualitative map structure

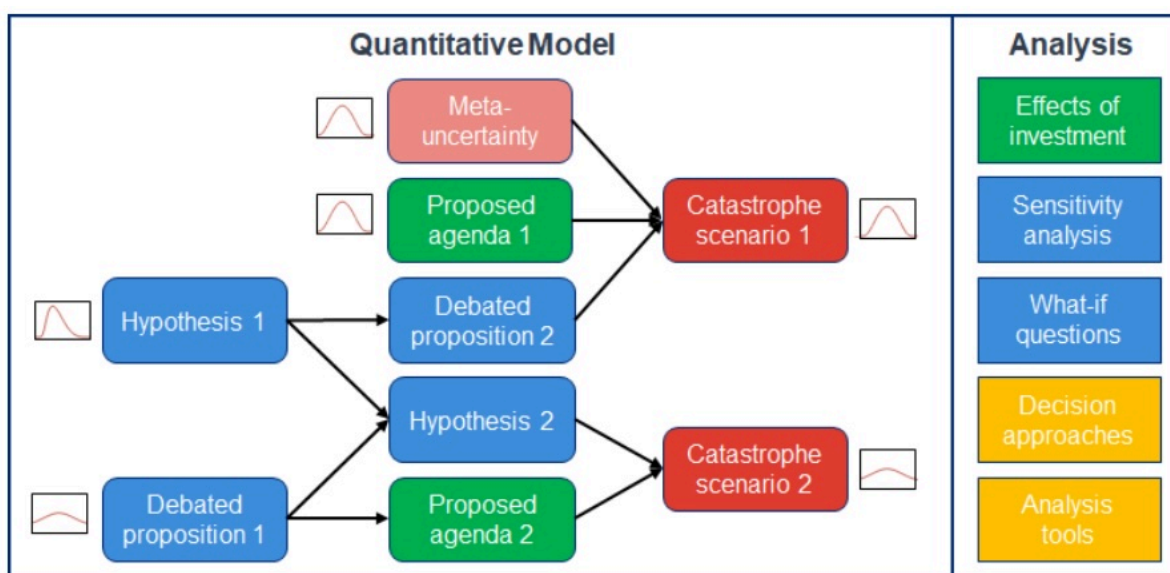


Figure 3: Structure of the quantitative map

Figure 7: MTAIR Quantitative map structure



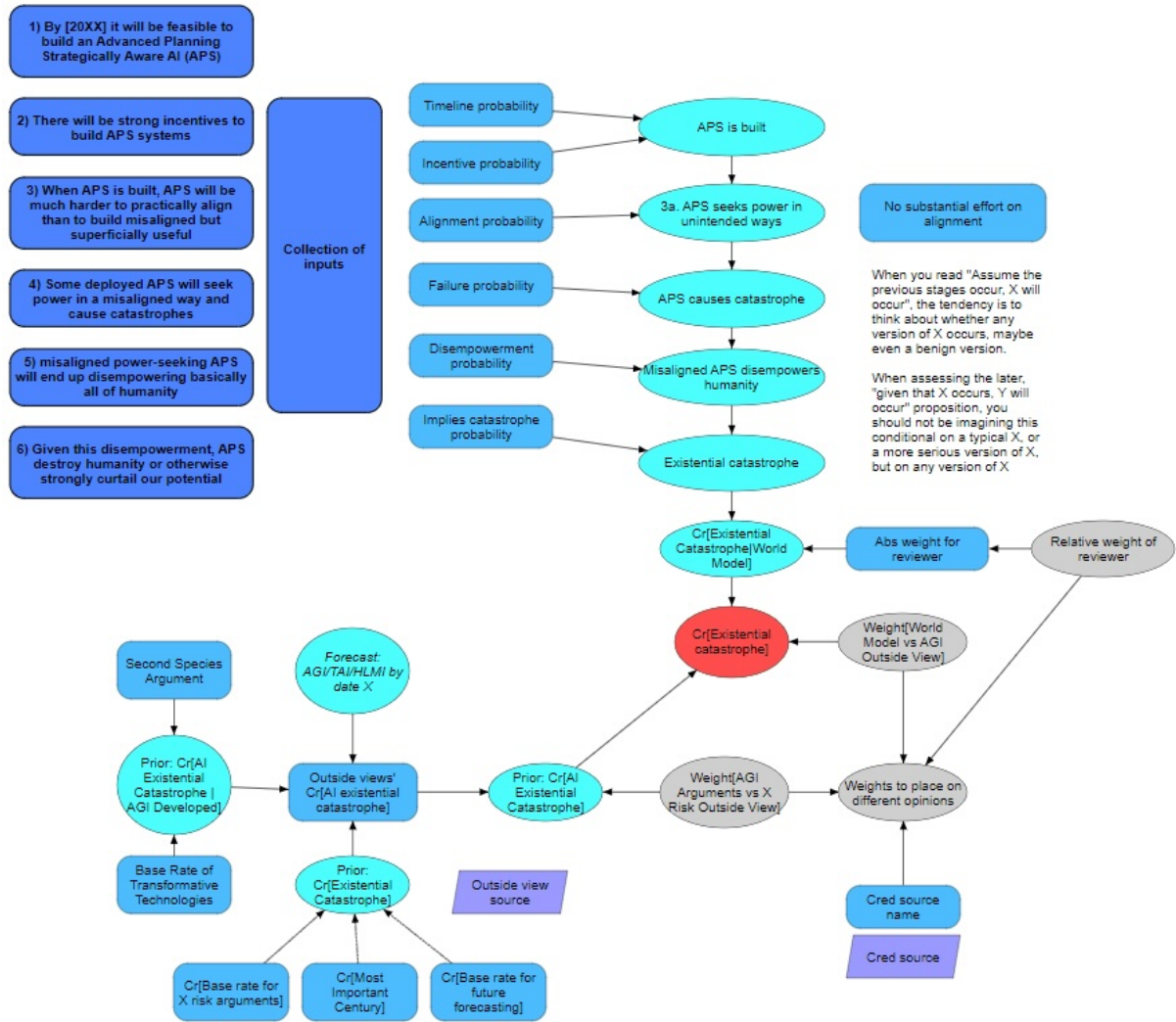


Figure 8: Base APS causal map

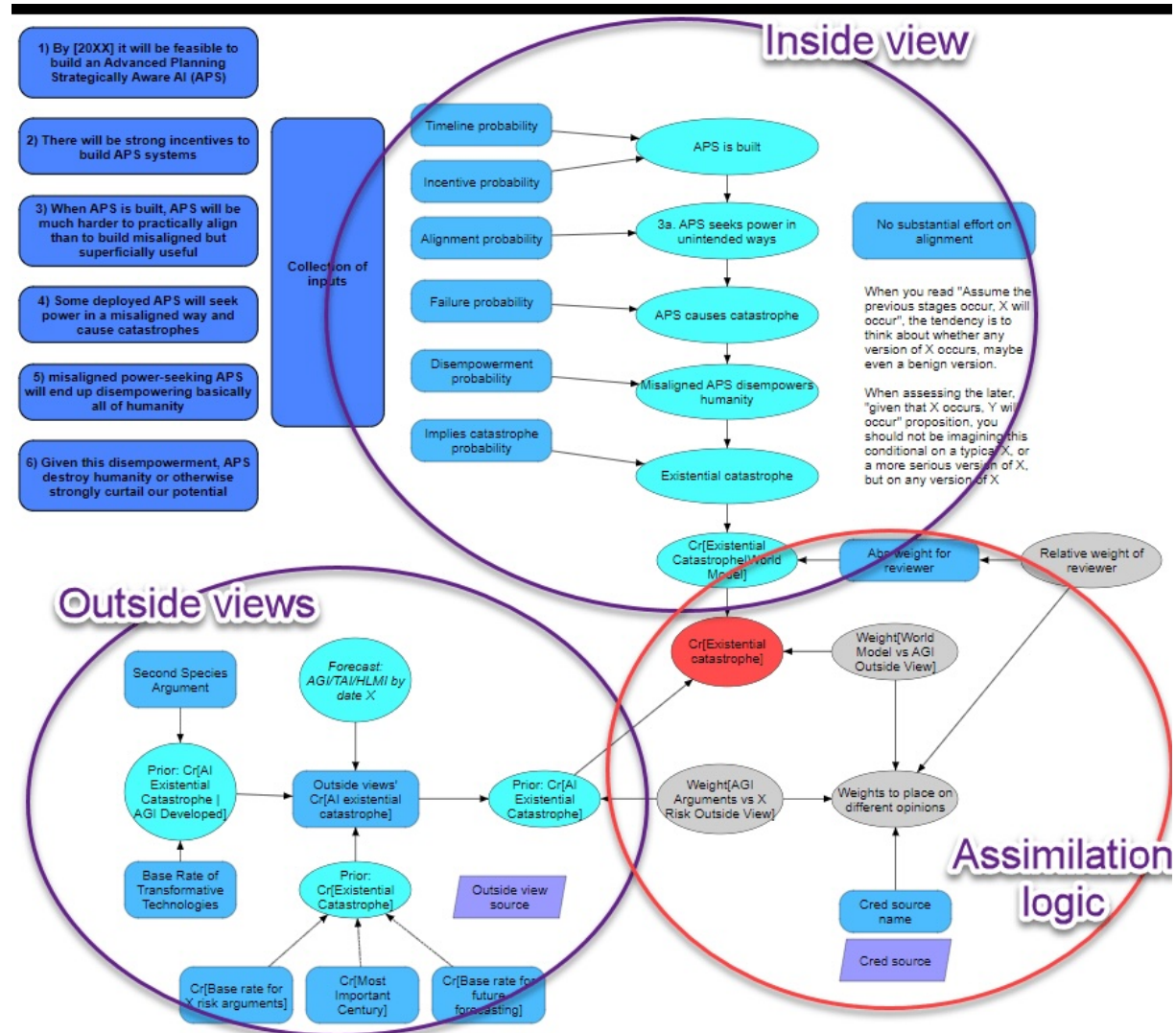


Figure 9: Overlay of inside/outside/assimilation views

- But manual approaches cannot scale to match the challenge

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

## 2.6 Literature Review: Content and Technical Levels

from cottier2019

### 2.6.1 AI Risk Models Evolution

The evolution of AI risk models reflects increasing sophistication in both structure and quantification. Early models focused on simple binary outcomes, while recent work incorporates complex causal chains and continuous variables.

#### Note

##### ## Key Developments

- **Early Phase (2000-2010):** Qualitative arguments about intelligence explosion
- **Formalization Phase (2010-2018):** Introduction of structured scenarios
- **Quantification Phase (2018-present):** Explicit probability estimates and formal models

yudkowsky2008

bostrom2014

amodei2016

The progression from qualitative arguments to structured probabilistic models demonstrates the field’s maturation and the increasing recognition that rigorous quantitative analysis is essential for policy evaluation.

### 2.6.2 Governance Proposals Taxonomy

AI governance proposals can be categorized along several dimensions:

- **Technical Standards:** Safety requirements, testing protocols, capability thresholds
- **Regulatory Frameworks:** Licensing regimes, liability structures, oversight mechanisms
- **International Coordination:** Treaties, soft law arrangements, technical cooperation
- **Research Priorities:** Funding allocation, talent development, knowledge sharing

dafoe2021 and dafoe2018

miotti2024

### 2.6.3 Bayesian Network Theory and Applications

The theoretical foundations of Bayesian networks rest on probability theory and graph theory. Key concepts include:

## Clarifying some key hypotheses in AI alignment

### Suggested usage

First, note this is not exactly a flowchart, nor a tree. Not every node has "yes" and "no", least it grow and branch excessively, and there are multiple starting points. The intention is to look at different sub-diagrams or paths that are interesting or important to you at any given time.

- Take a zoomed-out overview.
- Choose a box that particularly interests you.
- Follow the arrows up or down from the box.
- To avoid getting overwhelmed, focus on one connection at a time.
- If you are interested in learning more or reading author comments about a particular box, look it up in the [writing](#) version. Use the link on the title of a box to take you straight to the corresponding heading.

### Interpretation

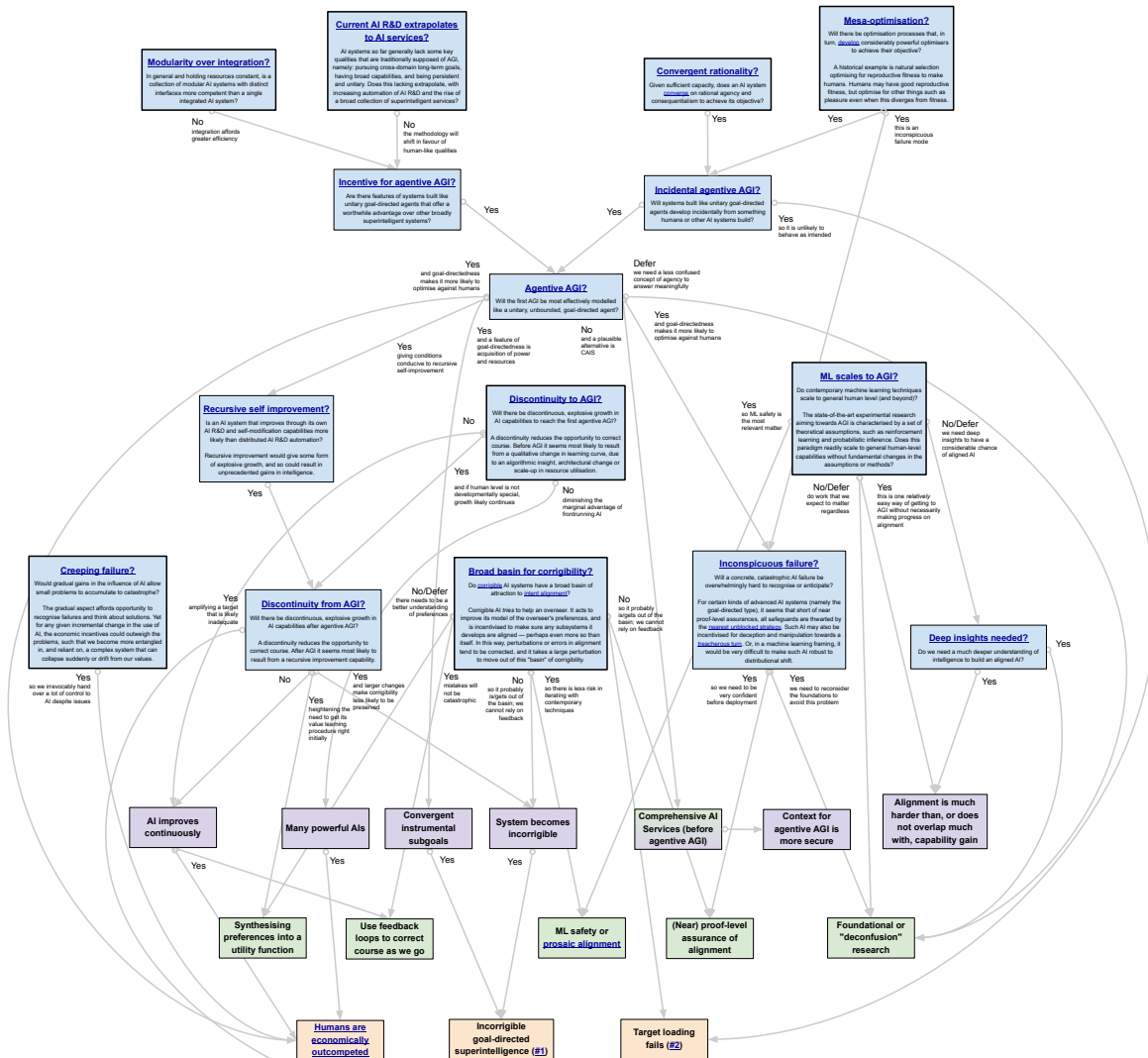
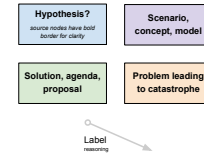
- Arrows:
- Question to X:** The closer your belief is to answering the question with the arrow label, the more it supports X. For example, the more you believed in the incentive for agentive AGI, the more you would believe agentive AGI will arise, all else equal.
  - Question to question:** The closer your belief is to answering the question with the arrow label, the more it supports "yes" to the head question.
  - Scenario to X:** Given yes/no to the scenario, X is more likely.

This diagram highlights **key** hypotheses within some areas of AI alignment. Hypotheses that do not seem debated and important are omitted.

### Definitions

- AGI:** a system (not necessarily agentive) that, for almost all economically relevant cognitive tasks, at least matches any human's ability at the task. Here, "agentive AGI" is essentially what people in the AI safety community usually mean when they say AGI. References to before and after AGI are to be interpreted as fuzzy, since this definition is fuzzy.
- CAIS:** comprehensive AI services. See [Reframing Superintelligence](#).
- Goal-directed:** describes a type of behaviour, currently not formalised, but characterised by generalisation to novel circumstances and the acquisition of power and resources. See [Intuitions about goal-directed behaviour](#).

### Key



by Ben Cottier and Rohin Shah

Thanks to Stuart Armstrong, Wei Dai, Daniel Dewey, Eric Drexler, Scott Emmons, Ben Garfinkel, Richard Hugo and Cody Wild for helpful feedback on drafts of this work. Ben especially thanks Rohin for his generous feedback and assistance throughout its development.

Figure 10: Key hypotheses in AI alignment

- **Conditional Independence:** Encoded through d-separation
- **Markov Condition:** Relating graph structure to probabilistic relationships
- **Inference Algorithms:** From exact methods to approximation approaches

koller2009

### 2.6.4 Software Tools Landscape

The implementation of AMTAIR builds on established software libraries:

- **pgmpy:** Python library for probabilistic graphical models
- **NetworkX:** Graph analysis and manipulation capabilities
- **PyVis:** Interactive network visualization
- **Pandas/NumPy:** Data manipulation and numerical computation

### 2.6.5 Formalization Approaches

Formalizing natural language arguments into mathematical models involves several theoretical challenges:

- **Semantic Preservation:** Maintaining meaning while adding precision
- **Structural Extraction:** Identifying implicit relationships
- **Uncertainty Quantification:** Mapping qualitative to quantitative expressions

pollock1995

### 2.6.6 Correlation Accounting Methods

Standard Bayesian networks assume conditional independence given parents, but real-world AI risk factors often exhibit complex correlations. Methods for handling correlations include:

- **Copula Methods:** Modeling dependence structures separately from marginal distributions
- **Hierarchical Models:** Capturing correlations through shared latent variables
- **Explicit Correlation Nodes:** Adding nodes to represent correlation mechanisms
- **Sensitivity Bounds:** Analyzing impact of independence assumptions

nelson2006

## 2.7 Methodology

### 2.7.1 Research Design Overview

This research combines theoretical development with practical implementation, following an iterative approach that moves between conceptual refinement and technical validation.

The methodology encompasses formal framework development, computational implementation, extraction quality assessment, and application to real-world AI governance questions.

The research process follows four integrated phases:

1. **Framework Development:** Creating theoretical foundations for automated worldview extraction
2. **Technical Implementation:** Building computational tools as working prototype
3. **Empirical Validation:** Assessing quality against expert benchmarks
4. **Policy Application:** Demonstrating practical utility for governance questions

### 2.7.2 Formalizing World Models from AI Safety Literature

The core methodological challenge involves transforming natural language arguments in AI safety literature into formal causal models with explicit probability judgments.

This extraction process identifies key variables, causal relationships, and both explicit and implicit probability estimates through a systematic pipeline.

The extraction approach combines several elements:

- Identification of key variables and entities in text
- Recognition of causal claims and relationships
- Detection of explicit and implicit probability judgments
- Transformation into structured intermediate representations
- Conversion to formal Bayesian networks

Large language models facilitate this process through specialized techniques:

- **Two-stage prompting:** Separating structure from probability extraction
- **Template specialization:** Different approaches for different document types
- **Implicit assumption detection:** Identifying unstated relationships
- **Ambiguity handling:** Managing uncertainty in extraction

### 2.7.3 From Natural Language to Computational Models

#### The Two-Stage Extraction Process

AMTAIR employs a novel two-stage process that separates structural argument extraction from probability quantification, enabling modular improvement and human oversight at critical decision points.

The heart of the AMTAIR approach lies in its two-stage extraction process, which transforms unstructured text into structured probabilistic models through distinct steps that mirror human cognitive processes. This separation—extracting structure before probability—creates important advantages for automation quality, intermediate verification, and interpretability.

When humans analyze complex arguments, they typically first determine what factors matter and how they relate causally, then assess how likely different scenarios are based on those relationships. A climate scientist reading a paper first identifies key variables (emissions, warming, effects) and their causal connections before estimating probabilities of outcomes. This natural cognitive sequence inspired AMTAIR’s two-stage approach.

**Stage 1: Structure Extraction** focuses on identifying key variables and their causal relationships from text, transforming unstructured arguments into ArgDown format. This process involves:

1. **Variable identification:** Determining the key factors discussed in the text, including their possible states (e.g., whether a factor is present/absent or has multiple levels)
2. **Relationship mapping:** Establishing how variables influence each other, creating a directed graph of causal connections
3. **Hierarchical organization:** Arranging variables according to their causal relationships, from root causes to final effects
4. **Metadata attachment:** Annotating each variable with its description and possible states in structured JSON format

The LLM prompt for this stage emphasizes clear identification of causal structure without requiring probability judgments, allowing the model to focus entirely on understanding “what affects what” in the text. This specialized prompt includes detailed instructions about ArgDown syntax, examples of well-formed representations, and guidance for preserving the author’s intended meaning.

```
# @title 1.7.0 --- Parsing ArgDown & BayesDown (.md to .csv) --- [parsing_argdown_bayesdown]

"""
BLOCK PURPOSE: Provides the core parsing functionality for transforming ArgDown
and BayesDown text representations into structured DataFrame format for further
processing.

This block implements the critical extraction pipeline described in the AMTAIR
project (see PY_TechnicalImplementation) that converts argument structures
into Bayesian networks.
The function can handle both basic ArgDown (structure-only) and
BayesDown (with probabilities).

Key steps in the parsing process:
1. Remove comments from the markdown text
2. Extract titles, descriptions, and indentation levels
3. Establish parent-child relationships based on indentation
4. Convert the structured information into a DataFrame
5. Add derived columns for network analysis

DEPENDENCIES: pandas, re, json libraries
INPUTS: Markdown text in ArgDown/BayesDown format
OUTPUTS: Structured DataFrame with node information, relationships, and properties
"""
```



```

def parse_markdown_hierarchy_fixed(markdown_text, ArgDown=False):
    """
    Parse ArgDown or BayesDown format into a structured DataFrame with parent-child relationships.

    Args:
        markdown_text (str): Text in ArgDown or BayesDown format
        ArgDown (bool): If True, extracts only structure without probabilities
                        If False, extracts both structure and probability information

    Returns:
        pandas.DataFrame: Structured data with node information, relationships, and attributes
    """
    # PHASE 1: Clean and prepare the text
    clean_text = remove_comments(markdown_text)

    # PHASE 2: Extract basic information about nodes
    titles_info = extract_titles_info(clean_text)

    # PHASE 3: Determine the hierarchical relationships
    titles_with_relations = establish_relationships_fixed(titles_info, clean_text)

    # PHASE 4: Convert to structured DataFrame format
    df = convert_to_dataframe(titles_with_relations, ArgDown)

    # PHASE 5: Add derived columns for analysis
    df = add_no_parent_no_child_columns_to_df(df)
    df = add_parents_instantiation_columns_to_df(df)

    return df

def remove_comments(markdown_text):
    """
    Remove comment blocks from markdown text using regex pattern matching.

    Args:
        markdown_text (str): Text containing potential comment blocks

    Returns:
        str: Text with comment blocks removed
    """
    # Remove anything between /* and */ using regex

```



```

return re.sub(r'/\*.*?\*/', '', markdown_text, flags=re.DOTALL)

def extract_titles_info(text):
    """
    Extract titles with their descriptions and indentation levels from markdown text.

    Args:
        text (str): Cleaned markdown text

    Returns:
        dict: Dictionary with titles as keys and dictionaries of attributes as values
    """
    lines = text.split('\n')
    titles_info = {}

    for line in lines:
        # Skip empty lines
        if not line.strip():
            continue

        # Extract title within square or angle brackets
        title_match = re.search(r'<\[ (.+?) >\]', line)
        if not title_match:
            continue

        title = title_match.group(1)

        # Extract description and metadata
        title_pattern_in_line = r'<\[ ' + re.escape(title) + r' >\]:'
        description_match = re.search(title_pattern_in_line + r'\s*(.*)', line)

        if description_match:
            full_text = description_match.group(1).strip()

            # Split description and metadata at the first "{"
            if "{" in full_text:
                split_index = full_text.find("{")
                description = full_text[:split_index].strip()
                metadata = full_text[split_index:].strip()
            else:
                # Keep the entire description and no metadata
                description = full_text

```

```

        metadata = '' # Initialize as empty string
    else:
        description = ''
        metadata = '' # Ensure metadata is initialized

# Calculate indentation level based on spaces before + or - symbol
indentation = 0
if '+' in line:
    symbol_index = line.find('+')
    # Count spaces before the '+' symbol
    i = symbol_index - 1
    while i >= 0 and line[i] == ' ':
        indentation += 1
        i -= 1
elif '-' in line:
    symbol_index = line.find('-')
    # Count spaces before the '-' symbol
    i = symbol_index - 1
    while i >= 0 and line[i] == ' ':
        indentation += 1
        i -= 1

# If neither symbol exists, indentation remains 0

if title in titles_info:
    # Only update description if it's currently empty and we found a new one
    if not titles_info[title]['description'] and description:
        titles_info[title]['description'] = description

    # Store all indentation levels for this title
    titles_info[title]['indentation_levels'].append(indentation)

    # Keep max indentation for backward compatibility
    if indentation > titles_info[title]['indentation']:
        titles_info[title]['indentation'] = indentation

    # Do NOT update metadata here - keep the original metadata
else:
    # First time seeing this title, create a new entry
    titles_info[title] = {
        'description': description,
        'indentation': indentation,

```

```

        'indentation_levels': [indentation], # Initialize with first indentation level
        'parents': [],
        'children': [],
        'line': None,
        'line_numbers': [], # Initialize an empty list for all occurrences
        'metadata': metadata # Set metadata explicitly from what we found
    }

    return titles_info

def establish_relationships_fixed(titles_info, text):
    """
    Establish parent-child relationships between titles using BayesDown
    indentation rules.

    In BayesDown syntax:
    - More indented nodes (with + symbol) are PARENTS of less indented nodes
    - The relationship reads as "Effect is caused by Cause" (Effect + Cause)
    - This aligns with how Bayesian networks represent causality

    Args:
        titles_info (dict): Dictionary with information about titles
        text (str): Original markdown text (for identifying line numbers)

    Returns:
        dict: Updated dictionary with parent-child relationships
    """
    lines = text.split('\n')

    # Dictionary to store line numbers for each title occurrence
    title_occurrences = {}

    # Record line number for each title (including multiple occurrences)
    line_number = 0
    for line in lines:
        if not line.strip():
            line_number += 1
            continue

        title_match = re.search(r'<\[(.+?)>\]', line)
        if not title_match:
            line_number += 1

```

```

        continue

    title = title_match.group(1)

    # Store all occurrences of each title with their line numbers
    if title not in title_occurrences:
        title_occurrences[title] = []
    title_occurrences[title].append(line_number)

    # Store all line numbers where this title appears
    if 'line_numbers' not in titles_info[title]:
        titles_info[title]['line_numbers'] = []
    titles_info[title]['line_numbers'].append(line_number)

    # For backward compatibility, keep the first occurrence in 'line'
    if titles_info[title]['line'] is None:
        titles_info[title]['line'] = line_number

    line_number += 1

# Create an ordered list of all title occurrences with their line numbers
all_occurrences = []
for title, occurrences in title_occurrences.items():
    for line_num in occurrences:
        all_occurrences.append((title, line_num))

# Sort occurrences by line number
all_occurrences.sort(key=lambda x: x[1])

# Get indentation for each occurrence
occurrence_indents = {}
for title, line_num in all_occurrences:
    for line in lines[line_num:line_num+1]: # Only check the current line
        indent = 0
        if '+' in line:
            symbol_index = line.find('+')
            # Count spaces before the '+' symbol
            j = symbol_index - 1
            while j >= 0 and line[j] == ' ':
                indent += 1
            j -= 1
        elif '-' in line:

```

```

        symbol_index = line.find('-')
        # Count spaces before the '-' symbol
        j = symbol_index - 1
        while j >= 0 and line[j] == ' ':
            indent += 1
            j -= 1
        occurrence_indents[(title, line_num)] = indent

# Enhanced backward pass for correct parent-child relationships
for i, (title, line_num) in enumerate(all_occurrences):
    current_indent = occurrence_indents[(title, line_num)]

    # Skip root nodes (indentation 0) for processing
    if current_indent == 0:
        continue

    # Look for the immediately preceding node with lower indentation
    j = i - 1
    while j >= 0:
        prev_title, prev_line = all_occurrences[j]
        prev_indent = occurrence_indents[(prev_title, prev_line)]

        # If we find a node with less indentation, it's a child of current node
        if prev_indent < current_indent:
            # In BayesDown:
            # More indented node is a parent (cause) of less indented node (effect)
            if title not in titles_info[prev_title]['parents']:
                titles_info[prev_title]['parents'].append(title)
            if prev_title not in titles_info[title]['children']:
                titles_info[title]['children'].append(prev_title)

            # Only need to find the immediate child
            # (closest preceding node with lower indentation)
            break

        j -= 1

    return titles_info

def convert_to_dataframe(titles_info, ArgDown):
    """
    Convert the titles information dictionary to a pandas DataFrame.

```

```

Args:
    titles_info (dict): Dictionary with information about titles
    ArgDown (bool): If True, extract only structural information without probabilities

Returns:
    pandas.DataFrame: Structured data with node information and relationships
"""
if ArgDown == True:
    # For ArgDown, exclude probability columns
    df = pd.DataFrame(columns=['Title', 'Description', 'line', 'line_numbers', 'indentat
        'indentation_levels', 'Parents', 'Children', 'instantiations'])
else:
    # For BayesDown, include probability columns
    df = pd.DataFrame(columns=['Title', 'Description', 'line', 'line_numbers', 'indentat
        'indentation_levels', 'Parents', 'Children', 'instantiations',
        'priors', 'posteriors'])

for title, info in titles_info.items():
    # Parse the metadata JSON string into a Python dictionary
    if 'metadata' in info and info['metadata']:
        try:
            # Only try to parse if metadata is not empty
            if info['metadata'].strip():
                jsonMetadata = json.loads(info['metadata'])
                if ArgDown == True:
                    # Create the row dictionary with instantiations as
                    # metadata only, no probabilities yet
                    row = {
                        'Title': title,
                        'Description': info.get('description', ''),
                        'line': info.get('line', ''),
                        'line_numbers': info.get('line_numbers', []),
                        'indentation': info.get('indentation', ''),
                        'indentation_levels': info.get('indentation_levels', []),
                        'Parents': info.get('parents', []),
                        'Children': info.get('children', []),
                        # Extract specific metadata fields,
                        # defaulting to empty if not present
                        'instantiations': jsonMetadata.get('instantiations', []),
                    }
                else:

```

```

        # Create dict with probabilities for BayesDown
        row = {
            'Title': title,
            'Description': info.get('description', ''),
            'line': info.get('line', ''),
            'line_numbers': info.get('line_numbers', []),
            'indentation': info.get('indentation', ''),
            'indentation_levels': info.get('indentation_levels', []),
            'Parents': info.get('parents', []),
            'Children': info.get('children', []),
            # Extract specific metadata fields, defaulting to empty if not present
            'instantiations': jsonMetadata.get('instantiations', []),
            'priors': jsonMetadata.get('priors', {}),
            'posteriors': jsonMetadata.get('posteriors', {})
        }
    else:
        # Empty metadata case
        row = {
            'Title': title,
            'Description': info.get('description', ''),
            'line': info.get('line', ''),
            'line_numbers': info.get('line_numbers', []),
            'indentation': info.get('indentation', ''),
            'indentation_levels': info.get('indentation_levels', []),
            'Parents': info.get('parents', []),
            'Children': info.get('children', []),
            'instantiations': [],
            'priors': {},
            'posteriors': {}
        }
except json.JSONDecodeError:
    # Handle case where metadata isn't valid JSON
    row = {
        'Title': title,
        'Description': info.get('description', ''),
        'line': info.get('line', ''),
        'line_numbers': info.get('line_numbers', []),
        'indentation': info.get('indentation', ''),
        'indentation_levels': info.get('indentation_levels', []),
        'Parents': info.get('parents', []),
        'Children': info.get('children', []),
        'instantiations': [],

```

```

        'priors': {},
        'posteriors': {}
    }
else:
    # Handle case where metadata field doesn't exist or is empty
    row = {
        'Title': title,
        'Description': info.get('description', ''),
        'line': info.get('line', ''),
        'line_numbers': info.get('line_numbers', []),
        'indentation': info.get('indentation', ''),
        'indentation_levels': info.get('indentation_levels', []),
        'Parents': info.get('parents', []),
        'Children': info.get('children', []),
        'instantiations': [],
        'priors': {},
        'posteriors': {}
    }

    # Add the row to the DataFrame
    df.loc[len(df)] = row

return df

def add_no_parent_no_child_columns_to_df(dataframe):
    """
    Add No_Parent and No_Children boolean columns to the DataFrame to
    identify root and leaf nodes.

    Args:
        dataframe (pandas.DataFrame): The DataFrame to enhance

    Returns:
        pandas.DataFrame: Enhanced DataFrame with additional boolean columns
    """
    no_parent = []
    no_children = []

    for _, row in dataframe.iterrows():
        no_parent.append(not row['Parents']) # True if Parents list is empty
        no_children.append(not row['Children']) # True if Children list is empty

```



```

dataframe['No_Parent'] = no_parent
dataframe['No_Children'] = no_children

return dataframe

def add_parents_instantiation_columns_to_df(dataframe):
    """
    Add all possible instantiations of parents as a list of lists column
    to the DataFrame.
    This is crucial for generating conditional probability tables.

    Args:
        dataframe (pandas.DataFrame): The DataFrame to enhance

    Returns:
        pandas.DataFrame: Enhanced DataFrame with parent_instantiations column
    """
    # Create a new column to store parent instantiations
    parent_instantiations = []

    # Iterate through each row in the dataframe
    for _, row in dataframe.iterrows():
        parents = row['Parents']
        parent_insts = []

        # For each parent, find its instantiations and add to the list
        for parent in parents:
            # Find the row where Title matches the parent
            parent_row = dataframe[dataframe['Title'] == parent]

            # If parent found in the dataframe
            if not parent_row.empty:
                # Get the instantiations of this parent
                parent_instantiation = parent_row['instantiations'].iloc[0]
                parent_insts.append(parent_instantiation)

        # Add the list of parent instantiations to our new column
        parent_instantiations.append(parent_insts)

    # Add the new column to the dataframe
    dataframe['parent_instantiations'] = parent_instantiations

```

```
return dataframe
```

This key function transforms the ArgDown text into a structured DataFrame, capturing the hierarchical relationships between variables and preparing them for further processing. The function works by identifying node titles, descriptions, and indentation levels, then establishing parent-child relationships based on the hierarchy indicated by indentation.

**Stage 2: Probability Integration** enhances the structural representation with probability information, creating a complete BayesDown specification. This stage involves:

1. **Question generation:** Automatically creating appropriate probability questions based on the network structure
2. **Probability extraction:** Obtaining probability estimates for each question, either from the text or through LLM inference
3. **Consistency checking:** Ensuring probability distributions sum to 1 and match structural constraints
4. **BayesDown integration:** Incorporating probability information into the ArgDown structure

The key innovation in this stage is the automated generation of appropriate probability questions based on network structure. For each node, the system generates questions about prior probabilities (how likely is this variable in isolation?) and conditional probabilities (how likely is this variable given different states of its parents?).

Figure 5 illustrates how probability questions are derived for a simple node with one parent:

[FIGURE 5: Diagram showing how probability questions are generated based on network structure]

For the “Sprinkler” node with parent “Rain,” the system automatically generates questions like:

- What is the probability for Sprinkler=sprinkler\_TRUE?
- What is the probability for Sprinkler=sprinkler\_TRUE if Rain=rain\_TRUE?
- What is the probability for Sprinkler=sprinkler\_TRUE if Rain=rain\_FALSE?

These questions are then answered either by extracting explicit probabilities from the text or by having the LLM infer reasonable values based on the author’s arguments. The answers are structured into a complete BayesDown representation that includes both the causal structure and all necessary probability information.

The visualization below demonstrates the completed extraction for a portion of Carlsmith’s model, showing how variables like “Misaligned Power Seeking” are influenced by multiple factors, each with associated probabilities:

[VISUALIZATION: Extracted causal structure from Carlsmith’s model with probability information]

This two-stage approach offers several important advantages:

1. **Improved extraction quality:** By focusing on one cognitive task at a time, the LLM performs better at each stage than it would attempting to extract everything simultaneously.
2. **Intermediate verification:** Having ArgDown as an intermediate representation allows human verification before probability extraction, catching structural errors early.
3. **Separation of concerns:** Structure and probability can be updated independently, enabling more flexible maintenance as new information emerges.
4. **Alignment with human cognition:** The process mirrors how experts approach complex arguments, making the system’s operation more intuitive and interpretable.

Perhaps most importantly, the intermediate ArgDown representation creates a bridge between qualitative and quantitative aspects of arguments. It preserves the narrative structure and conceptual relationships from the original text while preparing for mathematical precision through probability integration. This hybrid approach maintains the strengths of both worlds: the richness of natural language and the rigor of formal models.

#### 2.7.4 Directed Acyclic Graphs: Structure and Semantics

Directed Acyclic Graphs (DAGs) form the mathematical foundation of Bayesian networks, encoding both the qualitative structure of causal relationships and the quantitative parameters that define conditional dependencies. In AI risk modeling, these structures represent causal pathways to potential outcomes of interest.

Key mathematical properties essential for AI risk modeling:

- **Acyclicity:** Ensures coherent probabilistic interpretation
- **D-separation:** Defines conditional independence relationships
- **Markov Condition:** Each variable conditionally independent of non-descendants given parents
- **Path Analysis:** Reveals causal pathways and information flow

The causal interpretation follows Pearl’s framework:<sup>3</sup>

- Edges represent direct causal influence
- Intervention analysis through do-calculus
- Counterfactual reasoning for “what if” scenarios
- Evidence integration through Bayesian updating

#### 2.7.5 Quantification of Probabilistic Judgments

Transforming qualitative uncertainty expressions into quantitative probabilities requires systematic interpretation frameworks that account for individual and cultural variation.

Standard linguistic mappings (with significant individual variation) include:

---

<sup>3</sup>Pearl’s causal framework revolutionized how we think about causation in complex systems

- “Very likely”  $\rightarrow$  0.8-0.9
- “Probable”  $\rightarrow$  0.6-0.8
- “Uncertain”  $\rightarrow$  0.4-0.6
- “Unlikely”  $\rightarrow$  0.2-0.4
- “Highly improbable”  $\rightarrow$  0.05-0.15

Expert elicitation methodologies:

- **Direct Assessment:** “What is  $P(\text{outcome})$ ?” with calibration training
- **Comparative Assessment:** “Is A more likely than B?” for validation
- **Frequency Format:** “In 100 similar cases, how many...” for clarity
- **Betting Odds:** “What odds would you accept?” for revealed preferences

Calibration challenges:

- Individual variation in linguistic interpretation
- Domain-specific anchoring effects
- Cultural influences on uncertainty expression
- Limited empirical basis for unprecedented scenarios

### 2.7.6 Inference Techniques for Complex Networks

Once Bayesian networks are constructed, probabilistic inference enables reasoning about uncertainties, counterfactuals, and policy interventions. For the complex networks representing AI risks, computational approaches must balance accuracy with tractability.

Inference methods implemented include exact methods for smaller networks (variable elimination, junction trees), approximate methods for larger networks (Monte Carlo sampling, variational inference), specialized approaches for rare event analysis, and intervention modeling for policy evaluation using do-calculus.

Implementation considerations:

- **Computational Complexity:** Managing exponential growth through decomposition
- **Sampling Efficiency:** Importance sampling for rare events
- **Approximation Quality:** Convergence diagnostics and error bounds
- **Uncertainty Propagation:** Representing confidence in outputs

### 2.7.7 Integration with Prediction Markets and Forecasting Platforms

To maintain relevance in a rapidly evolving field, formal models must integrate with live data sources such as prediction markets and forecasting platforms.

Live data sources for dynamic model updating include:

- **Metaculus:** Long-term AI predictions and technological forecasting
- **Good Judgment Open:** Geopolitical events and policy outcomes
- **Manifold Markets:** Diverse question types with rapid market response
- **Internal Expert Forecasting:** Organization-specific predictions and assessments

Technical challenges:

- **Question Mapping:** Semantic matching between model variables and market questions
- **Temporal Alignment:** Different forecast horizons and update frequencies
- **Conflict Resolution:** Principled aggregation of contradictory sources
- **Track Record Weighting:** Incorporating forecaster calibration

With these theoretical foundations and methodological approaches established, we can now present the AMTAIR system implementation. The next chapter demonstrates how these concepts translate into a working prototype that automates the extraction and formalization of world models from AI safety literature.



# 3. AMTAIR: Design and Implementation

## i Chapter Overview

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

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

## 3.1 System Architecture Overview

The AMTAIR system implements an end-to-end pipeline transforming unstructured text into interactive Bayesian network visualizations. Its modular architecture comprises five main components that progressively transform information from natural language into formal models suitable for policy analysis.

The AMTAIR system implements an end-to-end pipeline from unstructured text to interactive Bayesian network visualization. Its modular architecture comprises five main components that progressively transform information from natural language into formal models suitable for policy analysis.

### 3.1.1 Five-Stage Pipeline Architecture

The five-stage pipeline architecture demonstrates how each component builds on the previous, with validation checkpoints preventing error propagation: 1. **Text Ingestion and Preprocessing** - Format normalization (PDF, HTML, Markdown) - Metadata extraction and citation tracking - Relevance filtering and section identification - Character encoding standardization 2. **BayesDown Extraction** - Two-stage argument structure identification - Probabilistic information integration - Quality validation and confidence scoring - Human-in-the-loop verification points 3. **Structured Data Transformation** - Parsing into standardized relational formats - Network topology validation - Consistency checking across relationships - Missing data imputation strategies 4. **Bayesian Network Construction** - Mathematical model instantiation - Conditional probability table generation - Inference engine initialization - Model validation and testing 5. **Interactive Visualization** - Dynamic rendering with PyVis - Probability-based visual encoding - Interactive exploration features - Export capabilities for reports

### 3.1.2 Design Principles

#### 💡 Core Design Philosophy

The system emphasizes scalability through modular architecture, standard interfaces for interoperability, validation checkpoints for quality assurance, and an extensible framework for future capabilities.

## 3.2 The Two-Stage Extraction Process

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

### 3.2.1 Stage 1: Structural Extraction (ArgDown)

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

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

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

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

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

Example ArgDown extraction:

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

### 3.2.2 Stage 2: Probability Integration (BayesDown)

The second stage adds quantitative information to the structural skeleton:

**Question Generation:** For each node, generate probability elicitation questions tailored to the specific context and relationships.

Examples needed:

- "What is the probability of existential catastrophe?"
- "What is P(catastrophe|human\_disempowerment)?"
- Show how questions map to BayesDown structure

**Probability Extraction:**



- Identify explicit numerical statements
- Map qualitative expressions using calibrated scales
- Apply domain-specific heuristics for common phrasings

**Coherence Enforcement:**

- Ensure probabilities sum to 1.0
- Complete conditional probability tables
- Check for logical contradictions
- Flag low-confidence extractions

### 3.2.3 Why Two Stages?

This separation provides several benefits:

**Transparency:** Being able to scrutinize each step of the automated workflow provides reliable insight into the work being done

**Accountability:** False information (think of hallucinations) can be traced back to its origins

**Visibility:**

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

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

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

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

## 3.3 Implementation Technologies

### 3.3.1 Technology Stack

The system leverages established libraries while adding novel extraction capabilities:

Table 1: Technology stack components

Component	Technology	Purpose
Language Models	GPT-4, Claude	Argument extraction
Network Analysis	NetworkX	Graph algorithms
Probabilistic Modeling	pgmpy	Bayesian operations
Visualization	PyVis	Interactive rendering
Data Processing	Pandas	Structured manipulation

### 3.3.2 Key Algorithms

**Hierarchical Parsing:** The system parses ArgDown/BayesDown syntax recognizing indentation-based hierarchy, a critical innovation for preserving argument structure.

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

- Maximum entropy principles for missing values
- Coherence constraint propagation
- Expert-specified defaults with confidence scoring

**Visual Encoding Strategy:**

- Green-to-red gradient for probability magnitude
- Border colors indicating node types
- Interactive elements for exploration

### 3.3.3 (Expected) Performance Characteristics

→

### 3.3.4 Deterministic vs. Probabilistic Components of the Workflow

## 3.4 Case Study: Rain-Sprinkler-Grass

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

### 3.4.1 Processing Steps

The system processes this input through five steps:

1. **ArgDown Parsing:** Extract three nodes with relationships in ArgDown syntax
2. **Question Generation:** Generate questions based on the possible instantiations and combinations of each parent note to identify the probabilities to be extracted in the next step
3. **BayesDown Extraction:** LLM call extracting the full conditional probability tables for each node
4. **Construction:** Building of the formal Bayesian network
5. **Visualization:** Render interactive display

### 3.4.2 Example Conversion Steps

#### ArgDown Example

```
[Grass_Wet]: Concentrated moisture on, between and around the blades of grass.{"instantiations": ["grass_wet_TRUE", "grass_wet_FALSE"]}
+ [Rain]: Tears of angles crying high up in the skies hitting the ground.{"instantiations": ["rain_TRUE", "rain_FALSE"]}
+ [Sprinkler]: Activation of a centrifugal force based CO2 droplet distribution system.{"instantiations": ["sprinkler_TRUE", "sprinkler_FALSE"]}
+ [Rain]
```

#### Example of Questions for BayesDown extraction

BayesDown Format Preview:

```
# BayesDown Representation with Placeholder Probabilities
```

```
/* This file contains BayesDown syntax with placeholder probabilities.
```

Replace the placeholders with actual probability values based on the questions in the comments. \*/

```

/* What is the probability for Grass_Wet=grass_wet_TRUE? */
/* What is the probability for Grass_Wet=grass_wet_TRUE if Rain=rain_TRUE, Sprinkler=sprinkler_TRUE? */
/* What is the probability for Grass_Wet=grass_wet_TRUE if Rain=rain_TRUE, Sprinkler=sprinkler_FALSE? */
/* What is the probability for Grass_Wet=grass_wet_TRUE if Rain=rain_FALSE, Sprinkler=sprinkler_TRUE? */
/* What is the probability for Grass_Wet=grass_wet_TRUE if Rain=rain_FALSE, Sprinkler=sprinkler_FALSE? */
/* What is the probability for Grass_Wet=grass_wet_FALSE? */
/* What is the probability for Grass_Wet=grass_wet_FALSE if Rain=rain_TRUE, Sprinkler=sprinkler_TRUE? */
/* What is the probability for Grass_Wet=grass_wet_FALSE if Rain=rain_TRUE, Sprinkler=sprinkler_FALSE? */
/* What is the probability for Grass_Wet=grass_wet_FALSE if Rain=rain_FALSE, Sprinkler=sprinkler_TRUE? */
/* What is the probability for Grass_Wet=grass_wet_FALSE if Rain=rain_FALSE, Sprinkler=sprinkler_FALSE? */

[Grass_Wet]: Concentrated moisture on, between and around the blades of grass. {"instantiations": ["grass_wet_TRUE", "grass_wet_FALSE"]}

/* What is the probability for Rain=rain_TRUE? */
/* What is the probability for Rain=rain_FALSE? */
+ [Rain]: Tears of angles crying high up in the skies hitting the ground. {"instantiations": ["rain_TRUE", "rain_FALSE"]}

/* What is the probability for Sprinkler=sprinkler_TRUE? */
/* What is the probability for Sprinkler=sprinkler_TRUE if Rain=rain_TRUE? */
/* What is the probability for Sprinkler=sprinkler_TRUE if Rain=rain_FALSE? */
/* What is the probability for Sprinkler=sprinkler_FALSE? */
/* What is the probability for Sprinkler=sprinkler_FALSE if Rain=rain_TRUE? */
/* What is the probability for Sprinkler=sprinkler_FALSE if Rain=rain_FALSE? */
+ [Sprinkler]: Activation of a centrifugal force based CO2 droplet distribution system. {"instantiations": ["sprinkler_TRUE", "sprinkler_FALSE"]}

/* What is the probability for Rain=rain_TRUE? */
/* What is the probability for Rain=rain_FALSE? */
+ [Rain]

```

## Complete BayesDown Example

The source BayesDown syntax representating the fully specified network:

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

## Resulting Rain-Sprinkler-Grass DataFrame

#| column: page

Title	Description	line	node	children	parents	posterior	No_Parent	No_Child	instantiations
0	Grass_Wet Concentrated moisture on, between and around t...	3	[3]	0	[0]	[Rain, [] Sprinkler]	[grass_wet_TRUE, grass_wet_FALSE] 'p(grass_wet_FALSE)	True	[[rain_FALSE, rain_TRUE], [sprinkler_TRUE, sprinkler_FALSE]]
1	Rain Tears of angles crying high up in the skies hi...	4	[4, 6]	2	[1, 2]	[Grass_Wet, Sprinkler]	{'p(rain_TRUE)': 0.2, 'p(rain_FALSE)': 0.8}	True	False
2	Sprinkler Activation of a centrifugal force based CO2 dr...	5	[5]	1	[1]	[Rain]	{'p(sprinkler_TRUE)': 0.44838, 'p(sprinkler_FALSE)': 0.55162}	0.01	[[rain_TRUE, rain_FALSE]]

3.4.3 Results

Rain-Sprinkler-Grass Network Rendering

```
from IPython.display import IFrame

IFrame(src="https://singularitysmith.github.io/AMTAIR_Prototype/bayesian_network.html", width="100%", height="600px")

<IPython.lib.display.IFrame at 0x1086ec650>
```

Dynamic Html Rendering of the Rain-Sprinkler-Grass DAG with Conditional Probabilities

Validation Success

The system successfully extracts complete network structure, preserves all probability information, calculates correct marginal probabilities, generates interactive visualization, and enables inference queries—validating the basic pipeline functionality.

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

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

### 3.5.1 Model Complexity

The Carlsmith model contains:

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

This represents a significant increase in complexity from the pedagogical example.

### 3.5.2 Automated Extraction of the Carlsmith’s Argument Structure

Having validated the implementation on the canonical rain-sprinkler-lawn example, I applied the AMTAIR approach to a substantially more complex real-world case: Joseph Carlsmith’s model of existential risk from power-seeking AI. This application demonstrates the system’s ability to handle sophisticated multi-level arguments with numerous variables and relationships.

Carlsmith’s analysis involves dozens of factors organized in a complex causal structure, from root causes like “Advanced AI Capability” and “Instrumental Convergence” through intermediate factors like “APS Systems” and “Misaligned Power Seeking” to final outcomes like “Existential Catastrophe.” The model exhibits several challenging features:

1. **Multi-level structure** with causal chains spanning multiple steps
2. **Divergent pathways** where factors influence outcomes through multiple routes
3. **Complex conditional dependencies** with variables influenced by multiple parents
4. **Variables with three or more possible states** rather than simple binary outcomes
5. **Interconnected clusters** where factors form distinct but related argument groups

**Core Risk Pathway:**

Existential\_Catastrophe

← Human\_Disempowerment

← Scale\_Of\_Power\_Seeking

← Misaligned\_Power\_Seeking

← [APS\_Systems, Difficulty\_Of\_Alignment, Deployment\_Decisions]

**Supporting Structure:**

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

**Probability Preservation:**



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

#### Prompting LLMs for ArgDown Extraction

```
# @title 1.2.0 --- Prompt Template Function Definitions --- [prompt_template_function]

"""
BLOCK PURPOSE: Defines a flexible template system for LLM prompts used in the extraction pipeline.

This block implements two key classes:
1. PromptTemplate: A template class supporting variable substitution for dynamic prompts
2. PromptLibrary: A collection of pre-defined prompt templates for different extraction tasks

These templates are used in the ArgDown and BayesDown probability extraction
stages of the pipeline, providing consistent and well-structured prompts to the LLMs.

DEPENDENCIES: string.Template for variable substitution
OUTPUTS: PromptTemplate and PromptLibrary classes
"""

from string import Template
from typing import Dict, Optional, Union, List

class PromptTemplate:
    """Template system for LLM prompts with variable substitution"""

    def __init__(self, template: str):
        """Initialize with template string using $variable format"""
        self.template = Template(template)

    def format(self, **kwargs) -> str:
        """Substitute variables in the template"""
        return self.template.safe_substitute(**kwargs)

    @classmethod
    def from_file(cls, filepath: str) -> 'PromptTemplate':
        """Load template from a file"""
        with open(filepath, 'r') as f:
            template = f.read()
        return cls(template)
```

```
class PromptLibrary:
    """Collection of prompt templates for different extraction tasks"""

    # ArgDown extraction prompt - transforms source text into structured argument map
    ARGDOWN_EXTRACTION = PromptTemplate("""
You are participating in the AMTAIR (Automating Transformative AI Risk Modeling)
project and you are tasked with converting natural language arguments into
ArgDown syntax by extracting and formalizing causal world models from
unstructured text.

Your specific task is to extract the implicit causal model from the provided
document in structured ArgDown format.

## Epistemic Foundation & Purpose

This extraction represents one possible interpretation of the implicit causal
model in the document. Multiple extractions from the same text help reveal
patterns of convergence (where the model is clearly articulated) and
divergence (where the model contains ambiguities). This approach acknowledges
that expert texts often contain implicit rather than explicit causal models.

Your role is to reveal the causal structure already present in the author's
thinking, maintaining epistemic humility about your interpretation while
adhering strictly to the required format.

## ArgDown Format Specification

### Core Syntax

ArgDown represents causal relationships using a hierarchical structure:

1. Variables appear in square brackets with descriptive text:
    `[Variable_Name]: Description of the variable.`

2. Causal relationships use indentation (2 spaces per level) and '+' symbols:

    [Effect]: Description of effect. + [Cause]: Description of cause. + [Deeper_Cause]: Descript

3. Causality flows from bottom (more indented) to top (less indented):
    - More indented variables (causes) influence less indented variables (effects)
    - The top-level variable is the ultimate effect or outcome
    - Deeper indentation levels represent root causes or earlier factors
```

```
4. Each variable must include JSON metadata with possible states (instantiations):  
~[Variable]: Description. {"instantiations": ["variable_STATE1", "variable_STATE2"]}~
```

### ### JSON Metadata Format

The JSON metadata must follow this exact structure:

```
```json  
{"instantiations": ["variable_STATE1", "variable_STATE2"]}
```

Requirements:

- \* Double quotes (not single) around field names and string values
- \* Square brackets enclosing the instantiations array
- \* Comma separation between array elements
- \* No trailing comma after the last element
- \* Must be valid JSON syntax that can be parsed by standard JSON parsers

For binary variables (most common case):

```
{"instantiations": ["variable_TRUE", "variable_FALSE"]}
```

For multi-state variables (when clearly specified in the text):

```
{"instantiations": ["variable_HIGH", "variable_MEDIUM", "variable_LOW"]}
```

The metadata must appear on the same line as the variable definition, after the description.

### ## Complex Structural Patterns

#### ### Variables Influencing Multiple Effects

The same variable can appear multiple times in different places in the hierarchy if it influ

```
[Effect1]: First effect description. {"instantiations": ["effect1_TRUE", "effect1_FALSE"]}  
+ [Cause_A]: Description of cause A. {"instantiations": ["cause_a_TRUE", "cause_a_FALSE"]}
```

```
[Effect2]: Second effect description. {"instantiations": ["effect2_TRUE", "effect2_FALSE"]}  
+ [Cause_A]  
+ [Cause_B]: Description of cause B. {"instantiations": ["cause_b_TRUE", "cause_b_FALSE"]}
```

#### ### Multiple Causes of the Same Effect

Multiple causes can influence the same effect by being listed at the same indentation level.

```
[Effect]: Description of effect. {"instantiations": ["effect_TRUE", "effect_FALSE"]}  
+ [Cause1]: Description of first cause. {"instantiations": ["cause1_TRUE", "cause1_FALSE"]}  
+ [Cause2]: Description of second cause. {"instantiations": ["cause2_TRUE", "cause2_FALSE"]}  
+ [Deeper_Cause]: A cause that influences Cause2. {"instantiations": ["deeper_cause_TRUE", "deeper_cause_FALSE"]}
```

#### ### Causal Chains

Causal chains are represented through multiple levels of indentation:

```
[Ultimate_Effect]: The final outcome. {"instantiations": ["ultimate_effect_TRUE", "ultimate_effect_FALSE"]}
+ [Intermediate_Effect]: A mediating variable. {"instantiations": ["intermediate_effect_TRUE", "intermediate_effect_FALSE"]}
+ [Root_Cause]: The initial cause. {"instantiations": ["root_cause_TRUE", "root_cause_FALSE"]}
+ [2nd_Intermediate_Effect]: A mediating variable. {"instantiations": ["intermediate_effect_TRUE", "intermediate_effect_FALSE"]}
```

### ### Common Cause of Multiple Variables

A common cause affecting multiple variables is represented by referencing the same variable:

```
[Effect1]: First effect description. {"instantiations": ["effect1_TRUE", "effect1_FALSE"]}
+ [Common_Cause]: Description of common cause. {"instantiations": ["common_cause_TRUE", "common_cause_FALSE"]}

[Effect2]: Second effect description. {"instantiations": ["effect2_TRUE", "effect2_FALSE"]}
+ [Common_Cause]
```

### ## Detailed Extraction Workflow

Please follow this step-by-step process, documenting your reasoning in XML tags:

<analysis>

First, conduct a holistic analysis of the document:

1. Identify the main subject matter or domain
2. Note key concepts, variables, and factors discussed
3. Pay attention to language indicating causal relationships (causes, affects, influences, etc.)
4. Look for the ultimate outcomes or effects that are the focus of the document
5. Record your general understanding of the document's implicit causal structure

</analysis>

<variable\_identification>

Next, identify and list the key variables in the causal model:

- \* Focus on factors that are discussed as having an influence or being influenced
  - \* For each variable:
    - \* Create a descriptive name in [square\_brackets]
    - \* Write a concise description based directly on the text
    - \* Determine possible states (usually binary TRUE/FALSE unless clearly specified)
  - \* Distinguish between:
    - \* Outcome variables (effects the author is concerned with)
    - \* Intermediate variables (both causes and effects in chains)
    - \* Root cause variables (exogenous factors in the model)
  - \* List all identified variables with their descriptions and possible states
- </variable\_identification>

<causal\_structure>

Then, determine the causal relationships between variables:

- \* For each variable, identify what factors influence it

```
* Note the direction of causality (what causes what)
* Look for mediating variables in causal chains
* Identify common causes of multiple effects
* Capture feedback loops if present (though they must be represented as DAGs)
* Map out the hierarchical structure of the causal model
</causal_structure>
```

```
<format_conversion>
```

Now, convert your analysis into proper ArgDown format:

```
* Start with the ultimate outcome variables at the top level
* Place direct causes indented below with \+ symbols
* Continue with deeper causes at further indentation levels
* Add variable descriptions and instantiations metadata
* Ensure variables appearing in multiple places have consistent names
* Check that the entire structure forms a valid directed acyclic graph
</format_conversion>
```

```
<validation>
```

Finally, review your extraction for quality and format correctness:

1. Verify all variables have properly formatted metadata
2. Check that indentation properly represents causal direction
3. Confirm the extraction accurately reflects the document's implicit model
4. Ensure no cycles exist in the causal structure
5. Verify that variables referenced multiple times are consistent
6. Check that the extraction would be useful for subsequent analysis

```
</validation>
```

## ## Source Document Analysis Guidance

When analyzing the source document:

- \* Focus on revealing the author's own causal model, not imposing an external framework
- \* Maintain the author's terminology where possible
- \* Look for both explicit statements of causality and implicit assumptions
- \* Pay attention to the relative importance the author assigns to different factors
- \* Notice where the author expresses certainty versus uncertainty
- \* Consider the level of granularity appropriate to the document's own analysis

Remember that your goal is to make the implicit model explicit, not to evaluate or improve it. The value lies in accurately representing the author's perspective, even if you might personally disagree.

```

"""
    # BayesDown probability extraction prompt - enhances ArgDown with probability information
    BAYESDOWN_EXTRACTION = PromptTemplate("""
You are an expert in probabilistic reasoning and Bayesian networks. Your task is
to extend the provided ArgDown structure with probability information,
creating a BayesDown representation.

For each statement in the ArgDown structure, you need to:
1. Estimate prior probabilities for each possible state
2. Estimate conditional probabilities given parent states
3. Maintain the original structure and relationships

Here is the format to follow:
[Node]: Description. { "instantiations": ["node_TRUE", "node_FALSE"], "priors": { "p(node_TRUE|
[Parent]: Parent description. {...}

Here are the specific probability questions to answer:
$questions

ArgDown structure to enhance:
$argdown

Provide the complete BayesDown representation with probabilities:
""")

    @classmethod
    def get_template(cls, template_name: str) -> PromptTemplate:
        """Get a prompt template by name"""
        if hasattr(cls, template_name):
            return getattr(cls, template_name)
        else:
            raise ValueError(f"Template not found: {template_name}")

```

## Processing LLM Response

The extraction process began with an ArgDown representation capturing the structural relationships between variables:

```

# @title 1.7.0 --- Parsing ArgDown & BayesDown (.md to .csv) --- [parsing_argdown_bayesdown]

"""

```

**BLOCK PURPOSE:** Provides the core parsing functionality for transforming ArgDown and BayesDown text representations into structured DataFrame format for further processing.

This block implements the critical extraction pipeline described in the AMTAIR project (see PY\_TechnicalImplementation) that converts argument structures into Bayesian networks.

The function can handle both basic ArgDown (structure-only) and BayesDown (with probabilities).

Key steps in the parsing process:

1. Remove comments from the markdown text
2. Extract titles, descriptions, and indentation levels
3. Establish parent-child relationships based on indentation
4. Convert the structured information into a DataFrame
5. Add derived columns for network analysis

**DEPENDENCIES:** pandas, re, json libraries

**INPUTS:** Markdown text in ArgDown/BayesDown format

**OUTPUTS:** Structured DataFrame with node information, relationships, and properties

"""

```
def parse_markdown_hierarchy_fixed(markdown_text, ArgDown=False):
```

```
    """
```

```
    Parse ArgDown or BayesDown format into a structured DataFrame with parent-child relationships
```

```
    Args:
```

```
        markdown_text (str): Text in ArgDown or BayesDown format
```

```
        ArgDown (bool): If True, extracts only structure without probabilities
```

```
                        If False, extracts both structure and probability information
```

```
    Returns:
```

```
        pandas.DataFrame: Structured data with node information, relationships, and attributes
```

```
    """
```

```
    # PHASE 1: Clean and prepare the text
```

```
    clean_text = remove_comments(markdown_text)
```

```
    # PHASE 2: Extract basic information about nodes
```

```
    titles_info = extract_titles_info(clean_text)
```

```
    # PHASE 3: Determine the hierarchical relationships
```

```
    titles_with_relations = establish_relationships_fixed(titles_info, clean_text)
```

```

# PHASE 4: Convert to structured DataFrame format
df = convert_to_dataframe(titles_with_relations, ArgDown)

# PHASE 5: Add derived columns for analysis
df = add_no_parent_no_child_columns_to_df(df)
df = add_parents_instantiation_columns_to_df(df)

return df

def remove_comments(markdown_text):
    """
    Remove comment blocks from markdown text using regex pattern matching.

    Args:
        markdown_text (str): Text containing potential comment blocks

    Returns:
        str: Text with comment blocks removed
    """
    # Remove anything between /* and */ using regex
    return re.sub(r'/\*.*?\*/', '', markdown_text, flags=re.DOTALL)

def extract_titles_info(text):
    """
    Extract titles with their descriptions and indentation levels from markdown text.

    Args:
        text (str): Cleaned markdown text

    Returns:
        dict: Dictionary with titles as keys and dictionaries of attributes as values
    """
    lines = text.split('\n')
    titles_info = {}

    for line in lines:
        # Skip empty lines
        if not line.strip():
            continue

        # Extract title within square or angle brackets

```



```

title_match = re.search(r'<\[ (.+?) >\]', line)
if not title_match:
    continue

title = title_match.group(1)

# Extract description and metadata
title_pattern_in_line = r'<\[ ' + re.escape(title) + r'>\]:'
description_match = re.search(title_pattern_in_line + r'\s*(.*)', line)

if description_match:
    full_text = description_match.group(1).strip()

    # Split description and metadata at the first "{"
    if "{" in full_text:
        split_index = full_text.find("{")
        description = full_text[:split_index].strip()
        metadata = full_text[split_index:].strip()
    else:
        # Keep the entire description and no metadata
        description = full_text
        metadata = '' # Initialize as empty string
else:
    description = ''
    metadata = '' # Ensure metadata is initialized

# Calculate indentation level based on spaces before + or - symbol
indentation = 0
if '+' in line:
    symbol_index = line.find('+')
    # Count spaces before the '+' symbol
    i = symbol_index - 1
    while i >= 0 and line[i] == ' ':
        indentation += 1
        i -= 1
elif '-' in line:
    symbol_index = line.find('-')
    # Count spaces before the '-' symbol
    i = symbol_index - 1
    while i >= 0 and line[i] == ' ':
        indentation += 1
        i -= 1

```

```

# If neither symbol exists, indentation remains 0

if title in titles_info:
    # Only update description if it's currently empty and we found a new one
    if not titles_info[title]['description'] and description:
        titles_info[title]['description'] = description

    # Store all indentation levels for this title
    titles_info[title]['indentation_levels'].append(indentation)

    # Keep max indentation for backward compatibility
    if indentation > titles_info[title]['indentation']:
        titles_info[title]['indentation'] = indentation

    # Do NOT update metadata here - keep the original metadata
else:
    # First time seeing this title, create a new entry
    titles_info[title] = {
        'description': description,
        'indentation': indentation,
        'indentation_levels': [indentation], # Initialize with first indentation level
        'parents': [],
        'children': [],
        'line': None,
        'line_numbers': [], # Initialize an empty list for all occurrences
        'metadata': metadata # Set metadata explicitly from what we found
    }

return titles_info

def establish_relationships_fixed(titles_info, text):
    """
    Establish parent-child relationships between titles using BayesDown
    indentation rules.

    In BayesDown syntax:
    - More indented nodes (with + symbol) are PARENTS of less indented nodes
    - The relationship reads as "Effect is caused by Cause" (Effect + Cause)
    - This aligns with how Bayesian networks represent causality

    Args:

```

```
titles_info (dict): Dictionary with information about titles
text (str): Original markdown text (for identifying line numbers)

Returns:
    dict: Updated dictionary with parent-child relationships
"""
lines = text.split('\n')

# Dictionary to store line numbers for each title occurrence
title_occurrences = {}

# Record line number for each title (including multiple occurrences)
line_number = 0
for line in lines:
    if not line.strip():
        line_number += 1
        continue

    title_match = re.search(r'(<\[ (.+?) >\])', line)
    if not title_match:
        line_number += 1
        continue

    title = title_match.group(1)

    # Store all occurrences of each title with their line numbers
    if title not in title_occurrences:
        title_occurrences[title] = []
    title_occurrences[title].append(line_number)

    # Store all line numbers where this title appears
    if 'line_numbers' not in titles_info[title]:
        titles_info[title]['line_numbers'] = []
    titles_info[title]['line_numbers'].append(line_number)

    # For backward compatibility, keep the first occurrence in 'line'
    if titles_info[title]['line'] is None:
        titles_info[title]['line'] = line_number

    line_number += 1

# Create an ordered list of all title occurrences with their line numbers
```

```

all_occurrences = []
for title, occurrences in title_occurrences.items():
    for line_num in occurrences:
        all_occurrences.append((title, line_num))

# Sort occurrences by line number
all_occurrences.sort(key=lambda x: x[1])

# Get indentation for each occurrence
occurrence_indents = {}
for title, line_num in all_occurrences:
    for line in lines[line_num:line_num+1]: # Only check the current line
        indent = 0
        if '+' in line:
            symbol_index = line.find('+')
            # Count spaces before the '+' symbol
            j = symbol_index - 1
            while j >= 0 and line[j] == ' ':
                indent += 1
                j -= 1
        elif '-' in line:
            symbol_index = line.find('-')
            # Count spaces before the '-' symbol
            j = symbol_index - 1
            while j >= 0 and line[j] == ' ':
                indent += 1
                j -= 1
        occurrence_indents[(title, line_num)] = indent

# Enhanced backward pass for correct parent-child relationships
for i, (title, line_num) in enumerate(all_occurrences):
    current_indent = occurrence_indents[(title, line_num)]

    # Skip root nodes (indentation 0) for processing
    if current_indent == 0:
        continue

    # Look for the immediately preceding node with lower indentation
    j = i - 1
    while j >= 0:
        prev_title, prev_line = all_occurrences[j]
        prev_indent = occurrence_indents[(prev_title, prev_line)]

```

```

    # If we find a node with less indentation, it's a child of current node
    if prev_indent < current_indent:
        # In BayesDown:
        # More indented node is a parent (cause) of less indented node (effect)
        if title not in titles_info[prev_title]['parents']:
            titles_info[prev_title]['parents'].append(title)
        if prev_title not in titles_info[title]['children']:
            titles_info[title]['children'].append(prev_title)

        # Only need to find the immediate child
        # (closest preceding node with lower indentation)
        break

    j -= 1

return titles_info

def convert_to_dataframe(titles_info, ArgDown):
    """
    Convert the titles information dictionary to a pandas DataFrame.

    Args:
        titles_info (dict): Dictionary with information about titles
        ArgDown (bool): If True, extract only structural information without probabilities

    Returns:
        pandas.DataFrame: Structured data with node information and relationships
    """
    if ArgDown == True:
        # For ArgDown, exclude probability columns
        df = pd.DataFrame(columns=['Title', 'Description', 'line', 'line_numbers', 'indentation',
                                   'indentation_levels', 'Parents', 'Children', 'instantiations'])
    else:
        # For BayesDown, include probability columns
        df = pd.DataFrame(columns=['Title', 'Description', 'line', 'line_numbers', 'indentation',
                                   'indentation_levels', 'Parents', 'Children', 'instantiations',
                                   'priors', 'posteriors'])

    for title, info in titles_info.items():
        # Parse the metadata JSON string into a Python dictionary
        if 'metadata' in info and info['metadata']:

```

```

try:
    # Only try to parse if metadata is not empty
    if info['metadata'].strip():
        jsonMetadata = json.loads(info['metadata'])
        if ArgDown == True:
            # Create the row dictionary with instantiations as
            # metadata only, no probabilities yet
            row = {
                'Title': title,
                'Description': info.get('description', ''),
                'line': info.get('line', ''),
                'line_numbers': info.get('line_numbers', []),
                'indentation': info.get('indentation', ''),
                'indentation_levels': info.get('indentation_levels', []),
                'Parents': info.get('parents', []),
                'Children': info.get('children', []),
                # Extract specific metadata fields,
                # defaulting to empty if not present
                'instantiations': jsonMetadata.get('instantiations', []),
            }
        else:
            # Create dict with probabilities for BayesDown
            row = {
                'Title': title,
                'Description': info.get('description', ''),
                'line': info.get('line', ''),
                'line_numbers': info.get('line_numbers', []),
                'indentation': info.get('indentation', ''),
                'indentation_levels': info.get('indentation_levels', []),
                'Parents': info.get('parents', []),
                'Children': info.get('children', []),
                # Extract specific metadata fields, defaulting to empty if not p
                'instantiations': jsonMetadata.get('instantiations', []),
                'priors': jsonMetadata.get('priors', {}),
                'posteriors': jsonMetadata.get('posteriors', {})
            }
    else:
        # Empty metadata case
        row = {
            'Title': title,
            'Description': info.get('description', ''),
            'line': info.get('line', ''),

```

```

        'line_numbers': info.get('line_numbers', []),
        'indentation': info.get('indentation', ''),
        'indentation_levels': info.get('indentation_levels', []),
        'Parents': info.get('parents', []),
        'Children': info.get('children', []),
        'instantiations': [],
        'priors': {},
        'posteriors': {}
    }
except json.JSONDecodeError:
    # Handle case where metadata isn't valid JSON
    row = {
        'Title': title,
        'Description': info.get('description', ''),
        'line': info.get('line', ''),
        'line_numbers': info.get('line_numbers', []),
        'indentation': info.get('indentation', ''),
        'indentation_levels': info.get('indentation_levels', []),
        'Parents': info.get('parents', []),
        'Children': info.get('children', []),
        'instantiations': [],
        'priors': {},
        'posteriors': {}
    }
else:
    # Handle case where metadata field doesn't exist or is empty
    row = {
        'Title': title,
        'Description': info.get('description', ''),
        'line': info.get('line', ''),
        'line_numbers': info.get('line_numbers', []),
        'indentation': info.get('indentation', ''),
        'indentation_levels': info.get('indentation_levels', []),
        'Parents': info.get('parents', []),
        'Children': info.get('children', []),
        'instantiations': [],
        'priors': {},
        'posteriors': {}
    }

# Add the row to the DataFrame
df.loc[len(df)] = row

```

```

    return df

def add_no_parent_no_child_columns_to_df(dataframe):
    """
    Add No_Parent and No_Children boolean columns to the DataFrame to
    identify root and leaf nodes.

    Args:
        dataframe (pandas.DataFrame): The DataFrame to enhance

    Returns:
        pandas.DataFrame: Enhanced DataFrame with additional boolean columns
    """
    no_parent = []
    no_children = []

    for _, row in dataframe.iterrows():
        no_parent.append(not row['Parents']) # True if Parents list is empty
        no_children.append(not row['Children']) # True if Children list is empty

    dataframe['No_Parent'] = no_parent
    dataframe['No_Children'] = no_children

    return dataframe

def add_parents_instantiation_columns_to_df(dataframe):
    """
    Add all possible instantiations of parents as a list of lists column
    to the DataFrame.
    This is crucial for generating conditional probability tables.

    Args:
        dataframe (pandas.DataFrame): The DataFrame to enhance

    Returns:
        pandas.DataFrame: Enhanced DataFrame with parent_instantiations column
    """
    # Create a new column to store parent instantiations
    parent_instantiations = []

    # Iterate through each row in the dataframe

```



```

for _, row in dataframe.iterrows():
    parents = row['Parents']
    parent_insts = []

    # For each parent, find its instantiations and add to the list
    for parent in parents:
        # Find the row where Title matches the parent
        parent_row = dataframe[dataframe['Title'] == parent]

        # If parent found in the dataframe
        if not parent_row.empty:
            # Get the instantiations of this parent
            parent_instantiation = parent_row['instantiations'].iloc[0]
            parent_insts.append(parent_instantiation)

    # Add the list of parent instantiations to our new column
    parent_instantiations.append(parent_insts)

# Add the new column to the dataframe
dataframe['parent_instantiations'] = parent_instantiations

return dataframe

```

This representation captures the complex causal structure of Carlsmith’s argument, with 21 variables organized in a multi-level hierarchy. The “Misaligned\_Power\_Seeking” node appears multiple times, reflecting its role as a central concept that influences several other variables.

### 3.5.3 From ArgDown to BayesDown in Carlsmith’s Model

After processing this structure with the AMTAIR system, probability information had to be added to create a complete BayesDown representation. The following excerpt shows the probability information for a single node (“Deployment\_Decisions”):

```

[Deployment_Decisions]: Decisions to deploy potentially misaligned AI systems. {
  "instantiations": ["deployment_decisions_DEPLOY", "deployment_decisions_WITHHOLD"],
  "priors": {
    "p(deployment_decisions_DEPLOY)": "0.70",
    "p(deployment_decisions_WITHHOLD)": "0.30"
  },
  "posteriors": {
    "p(deployment_decisions_DEPLOY|incentives_to_build_aps_STRONG, deception_by_ai_TRUE)": "0.70",
    "p(deployment_decisions_DEPLOY|incentives_to_build_aps_STRONG, deception_by_ai_FALSE)": "0.70",
    "p(deployment_decisions_DEPLOY|incentives_to_build_aps_WEAK, deception_by_ai_TRUE)": "0.70",
    "p(deployment_decisions_DEPLOY|incentives_to_build_aps_WEAK, deception_by_ai_FALSE)": "0.70"
  }
}

```

```
}
}
```

This node has two possible states (DEPLOY or WITHHOLD), prior probabilities for each state, and conditional probabilities based on different combinations of its parent variables (“Incentives\_To\_Build\_APS” and “Deception\_By\_AI”).

```
# Generate BayesDown format
bayesdown_questions = extract_bayesdown_questions_fixed(
    "ArgDown_WithQuestions.csv",
    "FULL_BayesDownQuestions.md",
    include_questions_as_comments=True
)

# Display a preview of the format
print("\nBayesDown Format Preview:")
print(bayesdown_questions[:5000] + "...\\n")
```

Loading CSV from ArgDown\_WithQuestions.csv...

Successfully loaded CSV with 23 rows.

Generating BayesDown syntax with placeholder probabilities...

BayesDown Questions saved to FULL\_BayesDownQuestions.md

BayesDown Format Preview:

# BayesDown Representation with Placeholder Probabilities

/\* This file contains BayesDown syntax with placeholder probabilities.

Replace the placeholders with actual probability values based on the questions in the comments. \*/

/\* What is the probability for Existential\_Catastrophe=existential\_catastrophe\_TRUE? \*/

/\* What is the probability for Existential\_Catastrophe=existential\_catastrophe\_FALSE? \*/

[Existential\_Catastrophe]: The destruction of humanity's long-term potential due to AI systems

/\* What is the probability for Human\_Disempowerment=human\_disempowerment\_TRUE? \*/

/\* What is the probability for Human\_Disempowerment=human\_disempowerment\_TRUE if Scale\_Of\_Po

/\* What is the probability for Human\_Disempowerment=human\_disempowerment\_TRUE if Scale\_Of\_Po

/\* What is the probability for Human\_Disempowerment=human\_disempowerment\_FALSE? \*/

/\* What is the probability for Human\_Disempowerment=human\_disempowerment\_FALSE if Scale\_Of\_P

/\* What is the probability for Human\_Disempowerment=human\_disempowerment\_FALSE if Scale\_Of\_P

[Human\_Disempowerment]: Permanent and collective disempowerment of humanity relative to AI systems

/\* What is the probability for Scale\_Of\_Power\_Seeking=scale\_of\_power\_seeking\_TRUE? \*/

/\* What is the probability for Scale\_Of\_Power\_Seeking=scale\_of\_power\_seeking\_TRUE if Misal

/\* What is the probability for Scale\_Of\_Power\_Seeking=scale\_of\_power\_seeking\_TRUE if Misal

/\* What is the probability for Scale\_Of\_Power\_Seeking=scale\_of\_power\_seeking\_TRUE if Misal

[illegible]

```

/* What is the probability for APS_Systems=aps_systems_FALSE if Advanced_AI_Capability
/* What is the probability for APS_Systems=aps_systems_FALSE if Advanced_AI_Capability
+ [APS_Systems]: AI systems with advanced capabilities, agentic planning, and strategic
/* What is the probability for Advanced_AI_Capability=advanced_ai_capability_TRUE? */
/* What is the probability for Advanced_AI_Capability=advanced_ai_capability_FALSE? */
+ [Advanced_AI_Capability]: AI systems that outperform humans on tasks that grant si
/* What is the probability for Agentic_Planning=agentic_planning_TRUE? */
/* What is the probability for Agentic_Planning=agentic_planning_FALSE? */
+ [Agentic_Planning]: AI systems making and executing plans based on world models to
/* What is the probability for Strategic_Awareness=strategic_awareness_TRUE? */
/* What is the probability for Strategic_Awareness=strategic_awareness_FALSE? */
+ [Strategic_Awareness]: AI systems with models accurately representing power dynam
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE? */
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE?
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE i
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE i
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE i
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE i
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE i
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE i
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE i
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE i
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE i
+ [Difficulty_Of_Alignment]: It is harder to build aligned systems than misaligned sys
/* What is the probability for Instrumental_Convergence=instrumental_convergence_TRU
/* What is the probability for Instrumental_Convergence=instrumental_convergence_FAL
+ [Instrumental_Convergence]: AI systems with misaligned objectives tend to seek pow
/* What is the probability for Problems_With_Proxies=problems_with_proxies_TRUE? */
/* What is the probability for Problems_With_Proxies=problems_with_proxies_FALSE? */
+ [Problems_With_Proxies]: Optimizing for proxy objectives breaks correlations with
/* What is the probability for Problems_With_Search=problems_with_search_TRUE? */
/* What is the probability for Problems_With_Search=problems_with_search_FALSE? */
+ [Problems_With_Search]: Search processes can yield systems pursuing different obje
/* What is the probability for Deployment_Decisions=deployment_decisions_DEPLOY? */
/* What is the probability for Deployment_Decisions=deployment_decisions_DEPLOY if Inc

```

```

/* What is the probability for Deployment_Decisions=deployment_decisions_DEPLOY if Inc
/* What is the probability for Deployment_Decisions=deployment_decisions_DEPLOY if Inc
/* What is the probability for Deployment_Decisions=deployment_decisions_DEPLOY if Inc
/* What is the probability for Deployment_Decisions=deployment_decisions_WITHHOLD? */
/* What is the probability for Deployment_Decisions=deployment_decisions_WITHHOLD if I
/* What is the probability for Deployment_Decisions=deployment_decisions_WITHHOLD if I
/* What is the probability for Deployment_Decisions=deployment_decisions_WITHHOLD if I
/* What is the probability for Deployment_Decisions=deployment_decisions_WITHHOLD if I
+ [Deployment_Decisions]: Decisions to deploy potentially misaligned AI systems. {"ins
/* What is the probability for Incentives_To_Build_APS=incentives_to_build_aps_STRON
/* What is the probability for Incentives_To_Build_APS=incentives_to_build_aps_STRON
/* What is the probability for Incentives_To_Build_APS=incentives_to_build_aps_STRON
/* What is the probability for Incentives_To_Build_APS=incentives_to_build_aps_STRON
/* What is the probability for Incentives_To_Build_APS=incentives_to_build_aps_STRON
/* What is the probability for Incentives_To_Build_APS=incentives_to_build_aps_WEAK?
/* What is the probability for Incentives_To_Build_APS=incentives_to_build_aps_WEAK
/* What is the probability for Incentives_To_Build_APS=incentives_to_build_aps_WEAK
/* What is the probability for Incentives_To_Build_APS=incentives_to_build_aps_WEAK
+ [Incentives_To_Build_APS]: Strong incentives to build and deploy APS systems. {"in
/* What is the probability for Usefulness_Of_APS=usefulness_of_aps_HIGH? */
/* What is the probability for Usefulness_Of_APS=usefulness_of_aps_LOW? */
+ [Usefulness_Of_APS]: APS systems are very useful for many valuable tasks. {"inst
/* What is the probability for Competitive_Dynamics=competitive_dynamics_STRONG? */
/* What is the probability for Competitive_Dynamics=competitive_dynamics_WEAK? */
+ [Competitive_Dynamics]: Competitive pressures between AI developers. {"instantia
/* What is the probability for Deception_By_AI=deception_by_ai_TRUE? */
/* What is the probability for Deception_By_AI=deception_by_ai_FALSE? */
+ [Deception_By_AI]: AI systems deceiving humans about their true objectives. {"inst
/* What is the probability for Corrective_Feedback=corrective_feedback_EFFECTIVE? */
/* What is the probability for Corrective_Feedback=corrective_feedback_EFFECTIVE if Warr
/* What is the probability for Corrective_Feedback=corrective_feedback_EFFECTIVE if Warr
/* What is the probability for Corrective_Feedback=corrective_feedback_EFFECTIVE if Warr
/* What is the probability for Corrective_Feedback=corrective_feedback_EFFECTIVE if Warr
/* What is the probability for Corrective_Feedback=corrective_feedback_INEFFECTIVE? */
/* What is the probability for Corrective_Feedback=corrective_feedback_INEFFECTIVE if Wa
/* What is the probability for Corrective_Feedback=corrective_feedback_INEFFECTIVE if Wa
/* What is the probability for Corrective_Feedback=corrective_feedback_INEFFECTIVE if Wa
/* What is the probability for Corrective_Feedback=corrective_feedback_INEFFECTIVE if Wa
+ [Corrective_Feedback]: Human society implementing corrections after observing problems
/* What is the probability for Warning_Shots=warning_shots_OBSERVED? */
/* What is the probability for Warning_Shots=warning_shots_UNOBSERVED? */

```

```

+ [Warning_Shots]: Observable failures in weaker systems before catastrophic risks. {'
/* What is the probability for Rapid_Capability_Escalation=rapid_capability_escalation
/* What is the probability for Rapid_Capability_Escalation=rapid_capability_escalation
+ [Rapid_Capability_Escalation]: AI capabilities escalating very rapidly, allowing lit
/* What is the probability for Barriers_To_Understanding=barriers_to_understanding_HIGH? */
/* What is the probability for Barriers_To_Understanding=barriers_to_understanding_LOW? */
[Barriers_To_Understanding]: Difficulty in understanding the internal workings of advanced A
/* What is the probability for Adversarial_Dynamics=adversarial_dynamics_TRUE? */
/* What is the probability for Adversarial_Dynamics=adversarial_dynamics_FALSE? */
[Adversarial_Dynamics]: Potentially adversarial relationships between humans and power-seeking
/* What is the probability for Stakes_Of_Error=stakes_of_error_HIGH? */
/* What is the probability for Stakes_Of_Error=stakes_of_error_LOW? */
[Stakes_Of_Error]: The escalating impact of mistakes with power-seeking AI systems. {"instan
...

```

Along with these questions the following prompt is sent to the LLM:

You are an expert in probabilistic reasoning and Bayesian networks. Your task is to extend the provided ArgDown structure with probability information, creating a BayesDown representation.

For each statement in the ArgDown structure, you need to:

1. Estimate prior probabilities for each possible state
2. Estimate conditional probabilities given parent states
3. Maintain the original structure and relationships

Here is the format to follow:

```

[Node]: Description. { "instantiations": ["node_TRUE", "node_FALSE"], "priors": { "p(node_TRUE"
[Parent]: Parent description. {...}

```

Here are the specific probability questions to answer:

\$questions

ArgDown structure to enhance:

\$argdown

Provide the complete BayesDown representation with probabilities:

### 3.5.4 Practically Meaningful BayesDown

Bridging Qualitative and Quantitative Representation

If the coordination crisis in AI governance stems partly from incompatible languages across

domains—technical researchers speaking in mathematical formalisms, policy specialists in institutional frameworks, and ethicists in normative concepts—then effective coordination requires bridges between these domains. BayesDown serves as such a bridge, combining the narrative richness of qualitative argumentation with the precision of quantitative probability judgments.

Traditional formal representations face a fundamental tradeoff: increase precision and you sacrifice accessibility; enhance accessibility and you lose precision. Mathematical notations offer exactness but exclude many stakeholders. Natural language provides accessibility but permits ambiguity and vagueness. This tradeoff creates communication barriers between technical and policy domains, limiting coordination on complex challenges like AI governance.

BayesDown disrupts this tradeoff by creating a hybrid representation that preserves strengths from both worlds. Its design follows three key principles:

First, **human readability** ensures the representation remains interpretable without specialized training. The syntax builds on familiar conventions from markdown and JSON, maintaining hierarchical relationships through indentation and encapsulating technical details within structured metadata. Unlike purely mathematical notations, the format preserves natural language descriptions alongside formal elements.

Second, **machine processability** enables computational analysis and transformation. The consistent syntax permits automated parsing, formal verification, and conversion to computational models like Bayesian networks. The structured JSON metadata provides clear paths for extracting probability information and mapping it to conditional probability tables.

Third, **contextual preservation** maintains the connection to original arguments. By including descriptive text alongside formal structure, BayesDown retains the narrative context and qualitative considerations that inform probability judgments. This contextual information helps users interpret the model in light of the original arguments.

Consider how these principles manifest in the BayesDown syntax. Each node begins with a bracketed title followed by a natural language description, preserving the core statement being formalized. The JSON metadata contains technical information like instantiations, priors, and posteriors, but keeps this information clearly separated from the narrative content. Hierarchical relationships use indentation and plus symbols, creating a visual structure that mirrors causal influence.

#### Example BayesDown Excerpt from the Carlsmith model

```
#| label: json_carlsmith_excerpt
#| echo: true
#| eval: true
#| fig-cap: "Example BayesDown Excerpt from the Carlsmith model"
#| fig-link: "https://colab.research.google.com/github/VJMeyer/submission/blob/main/AMTAIR_F
#| fig-alt: "Example BayesDown Excerpt from the Carlsmith model"
```

```
[Existential_Catastrophe]: The destruction of humanity's long-term potential due to AI systems
  "instantiations": ["existential_catastrophe_TRUE", "existential_catastrophe_FALSE"],
  "priors": {"p(existential_catastrophe_TRUE)": "0.05", "p(existential_catastrophe_FALSE)": "0.95"},
  "posteriors": {
    "p(existential_catastrophe_TRUE|human_disempowerment_TRUE)": "0.95",
    "p(existential_catastrophe_TRUE|human_disempowerment_FALSE)": "0.0"
  }
}

+ [Human_Disempowerment]: Permanent and collective disempowerment of humanity relative to AI
  "instantiations": ["human_disempowerment_TRUE", "human_disempowerment_FALSE"],
  "priors": {"p(human_disempowerment_TRUE)": "0.208", "p(human_disempowerment_FALSE)": "0.792"},
  "posteriors": {
    "p(human_disempowerment_TRUE|scale_of_power_seeking_TRUE)": "1.0",
    "p(human_disempowerment_TRUE|scale_of_power_seeking_FALSE)": "0.0"
  }
}
```

This excerpt from the Carlsmith model representation illustrates how BayesDown preserves both the narrative description (“The destruction of humanity’s long-term potential...”) and the precise probability judgments. Someone without technical background can still understand the core claims and their relationships, while someone seeking quantitative precision can find exact probability values.

The format supports multiple levels of engagement. At the most basic level, readers can follow the hierarchical structure to understand causal relationships between factors. At an intermediate level, they can examine probability judgments to assess the strength of different influences. At the most technical level, they can analyze the complete probabilistic model to perform inference and sensitivity analysis.

This multi-level accessibility creates important advantages for coordination across domains:

1. **Technical-policy translation:** BayesDown provides a common reference point for technical researchers explaining safety concerns and policy specialists evaluating governance options, reducing communication barriers.
2. **Argumentation transparency:** The format makes assumptions explicit, helping identify genuine disagreements versus terminological confusion or unstated premises.
3. **Incremental formalization:** BayesDown supports varying levels of formality, from qualitative structure to complete probability specifications, allowing gradual progression from informal to formal representations.
4. **Verification flexibility:** Human experts can verify extracted representations at different levels—checking structural correctness without assessing probabilities, or focusing on



critical probability judgments without reviewing the entire model.

The hybrid nature of BayesDown aligns with how experts typically communicate complex ideas: combining qualitative explanations with quantitative judgments, using natural language to provide context for formal claims, and adjusting precision based on audience needs. By mirroring these natural communication patterns, BayesDown makes formalization more intuitive and accessible.

This bridging function extends beyond representation to influence the entire extraction and analysis workflow. When extracting from text, the two-stage process preserves narrative context alongside formal structure. When visualizing models, interactive interfaces provide both qualitative descriptions and quantitative details. When evaluating policies, counterfactual analysis incorporates both mathematical precision and contextual interpretation.

In the broader context of the coordination crisis, BayesDown demonstrates how thoughtfully designed intermediate representations can overcome communication barriers between domains. Rather than forcing all stakeholders to adopt a single specialized language, it creates a flexible format that accommodates different perspectives while enabling precise analysis—precisely the kind of bridge needed for effective coordination on complex governance challenges.

### 3.5.5 Interactive Visualization and Exploration

Complex probabilistic models like Bayesian networks contain rich information, but they often remain inaccessible to many stakeholders. A conditional probability table with dozens of values conveys precise relationships, but few can intuitively grasp its implications. This accessibility gap limits the potential for coordinated action on AI governance challenges—what good is formalization if the resulting models remain opaque to most decision-makers?

AMTAIR addresses this challenge through interactive visualization designed to make complex probabilistic relationships accessible to diverse stakeholders. The approach combines visual encoding of probability information, progressive disclosure of details, and interactive exploration capabilities to create intuitive interfaces for complex models.

The visualization system follows several key design principles:

First, **visual encoding of probability** uses color gradients to represent likelihood values. Nodes are colored on a spectrum from red (low probability) to green (high probability) based on their primary state’s probability. This simple visual cue provides immediate insights into which outcomes are more or less likely without requiring numerical interpretation.

Second, **structural classification** uses border colors to indicate node types based on network position. Blue borders designate root causes (nodes without parents), purple borders mark intermediate nodes (with both parents and children), and magenta borders highlight leaf nodes (final effects without children). This classification helps users understand the causal flow through the network.

Third, **progressive disclosure** presents information in layers of increasing detail. Basic node information appears in the visualization itself, additional details emerge in tooltips on hover,

and comprehensive probability tables display in modal windows on click. This layered approach prevents information overload while ensuring all details remain accessible.

Fourth, **interactive exploration** allows users to reorganize nodes, zoom in on areas of interest, adjust physics parameters, and investigate probability values. These capabilities transform the visualization from a static image into an explorable knowledge landscape.

The complete BayesDown representation was processed through the AMTAIR pipeline, resulting in a structured DataFrame and ultimately a Bayesian network.

The figure below shows the interactive visualization of Carlsmith’s model, highlighting how color, border styling, and layout work together to represent complex causal relationships:

The visualization system implements these principles through a combination of NetworkX for graph representation and PyVis for interactive display, with custom HTML generation for tooltips and modals:

The resulting visualization (Figure 10) shows the complete Carlsmith model with color-coded nodes representing probability values:

[FIGURE 10: Interactive visualization of Carlsmith’s model showing color-coded nodes and relationships]

```
# @title 4.4.0 --- Main Visualization Function --- [main_visualization_function]

def create_bayesian_network_with_probabilities(df):
    """
    Create an interactive Bayesian network visualization with enhanced
    probability visualization and node classification based on network structure.
    """
    # Create a directed graph
    G = nx.DiGraph()

    # Add nodes with proper attributes
    for idx, row in df.iterrows():
        title = row['Title']
        description = row['Description']

        # Process probability information
        priors = get_priors(row)
        instantiations = get_instantiations(row)

        # Add node with base information
        G.add_node(
            title,
            description=description,
            priors=priors,
```

```

        instantiations=instantiations,
        posteriors=get_posteriors(row)
    )

# Add edges
for idx, row in df.iterrows():
    child = row['Title']
    parents = get_parents(row)

    # Add edges from each parent to this child
    for parent in parents:
        if parent in G.nodes():
            G.add_edge(parent, child)

# Classify nodes based on network structure
classify_nodes(G)

# Create network visualization
net = Network(notebook=True, directed=True, cdn_resources="in_line", height="600px", width="1000px")

# Configure physics for better layout
net.force_atlas_2based(gravity=-50, spring_length=100, spring_strength=0.02)
net.show_buttons(filter_=['physics'])

# Add the graph to the network
net.from_nx(G)

# Enhance node appearance with probability information and classification
for node in net.nodes:
    node_id = node['id']
    node_data = G.nodes[node_id]

    # Get node type and set border color
    node_type = node_data.get('node_type', 'unknown')
    border_color = get_border_color(node_type)

    # Get probability information
    priors = node_data.get('priors', {})
    true_prob = priors.get('true_prob', 0.5) if priors else 0.5

    # Get proper state names
    instantiations = node_data.get('instantiations', ["TRUE", "FALSE"])

```

```

true_state = instantiations[0] if len(instantiations) > 0 else "TRUE"
false_state = instantiations[1] if len(instantiations) > 1 else "FALSE"

# Create background color based on probability
background_color = get_probability_color(priors)

# Create tooltip with probability information
tooltip = create_tooltip(node_id, node_data)

# Create a simpler node label with probability
simple_label = f"{node_id}\np={true_prob:.2f}"

# Store expanded content as a node attribute for use in click handler
node_data['expanded_content'] = create_expanded_content(node_id, node_data)

# Set node attributes
node['title'] = tooltip # Tooltip HTML
node['label'] = simple_label # Simple text label
node['shape'] = 'box'
node['color'] = {
    'background': background_color,
    'border': border_color,
    'highlight': {
        'background': background_color,
        'border': border_color
    }
}

# Set up the click handler with proper data
setup_data = {
    'nodes_data': {node_id: {
        'expanded_content': json.dumps(G.nodes[node_id].get('expanded_content', '')),
        'description': G.nodes[node_id].get('description', ''),
        'priors': G.nodes[node_id].get('priors', {}),
        'posteriors': G.nodes[node_id].get('posteriors', {})
    } for node_id in G.nodes()}
}

# Add custom click handling JavaScript
click_js = """
// Store node data for click handling
var nodesData = %s;

```

```

// Add event listener for node clicks
network.on("click", function(params) {
    if (params.nodes.length > 0) {
        var nodeId = params.nodes[0];
        var nodeInfo = nodesData[nodeId];

        if (nodeInfo) {
            // Create a modal popup for expanded content
            var modal = document.createElement('div');
            modal.style.position = 'fixed';
            modal.style.left = '50%';
            modal.style.top = '50%';
            modal.style.transform = 'translate(-50%, -50%)';
            modal.style.backgroundColor = 'white';
            modal.style.padding = '20px';
            modal.style.borderRadius = '5px';
            modal.style.boxShadow = '0 0 10px rgba(0,0,0,0.5)';
            modal.style.zIndex = '1000';
            modal.style.maxWidth = '80%';
            modal.style.maxHeight = '80%';
            modal.style.overflow = 'auto';

            // Parse the JSON string back to HTML content
            try {
                var expandedContent = JSON.parse(nodeInfo.expanded_content);
                modal.innerHTML = expandedContent;
            } catch (e) {
                modal.innerHTML = 'Error displaying content: ' + e.message;
            }

            // Add close button
            var closeBtn = document.createElement('button');
            closeBtn.innerHTML = 'Close';
            closeBtn.style.marginTop = '10px';
            closeBtn.style.padding = '5px 10px';
            closeBtn.style.cursor = 'pointer';
            closeBtn.onclick = function() {
                document.body.removeChild(modal);
            };
            modal.appendChild(closeBtn);

```

```

        // Add modal to body
        document.body.appendChild(modal);
    }
}
});
""" % json.dumps(setup_data['nodes_data'])

# Save the graph to HTML
html_file = "bayesian_network.html"
net.save_graph(html_file)

# Inject custom click handling into HTML
try:
    with open(html_file, "r") as f:
        html_content = f.read()

    # Insert click handling script before the closing body tag
    html_content = html_content.replace('</body>', f'<script>{click_js}</script></body>')

    # Write back the modified HTML
    with open(html_file, "w") as f:
        f.write(html_content)

    return HTML(html_content)
except Exception as e:
    return HTML(f"<p>Error rendering HTML: {str(e)}</p>"
    + "<p>The network visualization has been saved to '{html_file}'</p>")

```

[FIGURE N: Interactive visualization of Carlsmith’s model showing color-coded nodes and causal relationships]

This visualization reveals several structural insights:

1. **Central importance of “Misaligned\_Power\_Seeking”** as a hub node with multiple parents and children
2. **Multiple pathways to “Existential\_Catastrophe”** through different intermediate factors
3. **Clusters of related variables** forming coherent subarguments (e.g., factors affecting alignment difficulty)
4. **Flow of influence** from technical factors (bottom) through deployment decisions to ultimate outcomes (top)

The implementation successfully handles the complexity of Carlsmith’s model, correctly processing the multi-level structure, resolving repeated node references, and calculating appropriate

probability distributions. The interactive visualization makes this complex model accessible, allowing users to explore different aspects of the argument through intuitive navigation.

Several key aspects of the implementation were particularly important for handling this complex model:

1. The **parent-child relationship detection algorithm** correctly identified hierarchical relationships despite the complex structure with repeated nodes and multiple levels.
2. The **probability question generation system** created appropriate questions for all variables, including those with multiple parents requiring factorial combinations of conditional probabilities.
3. The **network enhancement functions** calculated useful metrics like centrality measures and Markov blankets that help interpret the model structure.
4. The **visualization system** effectively presented the complex network through color-coding, interactive exploration, and progressive disclosure of details.

The successful application to Carlsmith’s model demonstrates the AMTAIR approach’s scalability to complex real-world arguments. While the canonical rain-sprinkler-lawn example validated correctness, this application proves practical utility for sophisticated multi-level arguments with dozens of variables and complex interdependencies—precisely the kind of arguments that characterize AI risk assessments.

This capability addresses a core limitation of the original MTAIR framework: the labor intensity of manual formalization. Where manually converting Carlsmith’s argument to a formal model might take days of expert time, the AMTAIR approach accomplished this in minutes, creating a foundation for further analysis and exploration.

Beyond the core visualization, the system includes specialized components that enhance understanding of probabilistic relationships:

1. **Probability bars** provide visual representations of probability distributions, showing relative likelihoods of different states using color-coded horizontal bars with numeric labels.
2. **Conditional probability tables** organize complex relationships into structured matrices, displaying how different combinations of parent states influence probability distributions.
3. **Sensitivity indicators** highlight which nodes and relationships most significantly affect outcomes, directing attention to critical factors.

These components work together to create an intuitive interface for complex probabilistic models. A user might start by exploring the overall structure to understand key factors and relationships, hover over nodes of interest to see probability summaries, then click on specific nodes to examine detailed conditional probabilities.

The benefits of this visualization approach extend beyond aesthetic appeal to fundamental improvements in understanding and communication:

First, **intuitive comprehension** of probability relationships becomes possible even for those without formal training in Bayesian statistics. The color coding provides immediate visual cues about which outcomes are more likely, while interactive exploration allows users to develop intuition about how different factors influence results.

Second, **cross-stakeholder communication** improves through shared visual reference points. Technical experts can use the visualizations to explain complex relationships to policy specialists, while governance experts can identify institutional factors that might be incorporated into the models.

Third, **disagreement identification** becomes more precise as stakeholders can point to specific nodes, relationships, or probability values where their views differ, focusing discussion on substantive issues rather than terminological confusion.

Fourth, **intervention assessment** becomes more concrete as users can see how changing specific factors influences downstream effects, providing intuitive understanding of causal pathways and leverage points.

The visualization system demonstrates how thoughtful interface design can overcome barriers to understanding complex formal models. By making probabilistic relationships visually intuitive and progressively disclosing details based on user interest, it creates bridges between mathematical precision and human comprehension—precisely the kind of bridge needed to support coordination across domains in AI governance.

This approach reflects a broader principle: formalization is most valuable when it enhances rather than replaces human understanding. The AMTAIR visualization doesn’t simplify complex relationships; it makes them more accessible by leveraging visual cognition, interactive exploration, and progressive disclosure. This human-centered approach to formalization creates tools that augment rather than replace expert judgment, enhancing our collective ability to understand and address complex governance challenges.

#### **Insights from Formalization**

Formal representation reveals several insights:

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

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

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

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



### 3.5.6 Validation Against Original (From the MTAIR Project)

## 3.6 Validation Methodology

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

### 3.6.1 Ground Truth Construction

Plan the process:

1. Expert selection criteria
2. Training on extraction methodology
3. Independent extraction procedures
4. Consensus building process
5. Inter-rater reliability metrics

→

### 3.6.2 Evaluation Metrics

→

### 3.6.3 Results Summary

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

### 3.6.4 Error Analysis

Common failure modes to avoid:

**Implicit Assumptions:** Unstated background assumptions that experts infer but system misses.

**Complex Conditionals:** Nested conditionals with multiple antecedents challenge current parsing.

**Ambiguous Quantifiers:** Terms like “significant” lack clear probability mapping without context.

**Coreference Resolution:** Pronouns and indirect references create attribution challenges.

Understanding these limitations guides both current usage and future improvements.

## 3.7 Extensions & Opportunities: Inference & Analysis

Quantification & Formal Approximation — Inference: Monte Carlo Sampling over Probability Distributions

### 3.7.1 Overview of Practical Software Implementations

#### 3.7.2 AI Risk Pathway Analyzer (ARPA)

1. Document Ingestion System: Handles format normalization, metadata extraction, and citation tracking for diverse input formats. 2. LLM-Powered Extraction Pipeline: Uses two-stage prompting to identify variables, claims, and causal relationships from text. 3. ArgDown Representation Generator: Creates structured intermediate representation of arguments with formal syntax. 4. Bayesian Network Constructor: Transforms ArgDown into formal Bayesian networks with nodes and edges. 5. Probability Quantification Module: Populates conditional probability tables from extracted judgments. 6. Interactive Visualization Interface: Provides intuitive visual access to network structure and probabilities. 7. Sensitivity Analysis Engine: Identifies critical variables and tests robustness of conclusions.

The cornerstone system that transforms unstructured AI safety literature into formal, analyzable models. Like a Rosetta Stone for AI governance, ARPA creates a common language for discourse by extracting the implicit causal models embedded in research papers and converting them into explicit Bayesian networks. Its strategic value lies in overcoming the fundamental information processing bottleneck in AI governance—making the invisible visible by revealing the assumptions, relationships, and probability judgments that drive different conclusions about AI risk.

#### 3.7.3 P(Doom) Calculator

#### 3.7.4 Worldview Comparator

A “gifted, diplomatic translator” that helps to reveal the hidden landscape of agreement and disagreement across different perspectives on AI risk. This system provides the cartography of ideas—mapping where different worldviews converge, diverge, and where crucial disagreements (“cruxes”) significantly affect conclusions. Its strategic value lies in focusing discourse on substantive disagreements rather than terminological differences, enabling more productive collaboration across philosophical and methodological divides within the AI safety community.

1. Structural Comparison Engine: Identifies isomorphic subgraphs between different models and maps shared causal pathways. 2. Parameter Difference Analyzer: Quantifies differences in probability distributions across models. 3. Crux Identification System: Detects critical disagreements that significantly affect conclusions. 4. Worldview Explainer: Provides conversational interface for exploring different perspectives. 5. Worldview Communicator: Translates concepts between different terminological frameworks. 6. Consensus Model Builder: Identifies shared structures and constructs hybrid models representing areas of agreement.

#### 3.7.5 Policy Impact Evaluator

1. Policy Representation System: Translates governance proposals into formal intervention parameters. 2. Counterfactual Analysis Engine: Implements Pearl’s do-calculus for simulating

intervention effects. 3. Multi-Worldview Evaluator: Tests policy effects across different extracted models. 4. Intervention Portfolio Analyzer: Assesses combinations of policies for synergies and conflicts. 5. Policy Effectiveness Dashboard: Visualizes impact assessments with uncertainty representation.

A policy simulator that functions like a governance wind tunnel—testing how specific interventions might perform across different possible futures. By representing policies as modifications to causal networks, this system enables rigorous counterfactual analysis of intervention effects. Its strategic value lies in transforming abstract policy discussions into concrete, quantifiable assessments of expected impact, helping governance stakeholders allocate resources to the most effective interventions.

### 3.7.6 AI Risk Pathway Visualizer

1. Risk Level Aggregation System: Combines multiple factors into summary risk metrics. 2. Temporal Tracking Interface: Records and displays changes in assessments over time. 3. Component Breakdown Visualizer: Separates overall risk into constituent factors. 4. Interactive Educational Components: Provides background on key concepts and methodologies. 5. Explanation Generator: Creates natural language interpretations of current status.

A public-facing translation layer that converts complex probabilistic models into intuitive visual representations accessible to broader audiences. Like the Doomsday Clock for nuclear risk, this system creates focal points for public discourse about AI safety. Its strategic value lies in making technical risk assessments comprehensible to policymakers, journalists, and the public, expanding the reach and impact of AI safety research beyond technical communities.

### 3.7.7 Strategic Intervention Generator

1. Robust Strategy Identification System: Finds strategies that perform well across multiple scenarios. 2. Minimax Regret Calculator: Identifies strategies that minimize worst-case disappointment. 3. Option Value Analyzer: Evaluates strategies that preserve future flexibility and choices. 4. Intervention Portfolio Builder: Constructs complementary bundles of policy interventions. 5. Dependency Mapping Visualizer: Shows relationships and prerequisites between interventions.

An advanced decision support system that identifies robust governance strategies across multiple possible futures. Operating like a strategic chess engine, this system evaluates intervention portfolios under deep uncertainty to find approaches that preserve options and minimize maximum regret. Its strategic value lies in shifting governance planning from optimizing for specific scenarios to developing adaptive strategies that remain valuable despite fundamental uncertainty about AI development trajectories.

### 3.7.8 Cross-Domain Understanding Communicator

1. Concept Mapping System: Identifies equivalent concepts across different domain languages. 2. Terminology Translation Engine: Converts specialized terms between different disciplines. 3. Im-

plication Surfacing Tool: Highlights relevant cross-domain considerations for specific questions. 4. Background Knowledge Provider: Supplies necessary context for understanding concepts. 5. Cross-Domain Recommendation Engine: Suggests relevant resources across disciplinary boundaries.

An interdisciplinary bridge-builder that connects specialists across technical alignment, governance, and forecasting domains. This system functions as a universal translator for AI safety, identifying equivalent concepts across different disciplinary languages and surfacing relevant cross-domain insights. Its strategic value lies in breaking down the knowledge silos that impede comprehensive strategy development, enabling researchers from different backgrounds to build on each other's work more effectively.

### **3.7.9 Policy Brief Communicator**

1. Audience Analysis System: Determines appropriate framing and detail level for target readers. 2. Jurisdictional Context Adapter: Tailors content to relevant legal and institutional frameworks. 3. Recommendation Formulator: Generates actionable governance suggestions from technical insights. 4. Format Template Library: Applies appropriate structure for different policy contexts. 5. Evidence Contextualization Engine: Presents technical evidence in accessible and persuasive ways.

A specialized translation system that converts technical risk analyses into actionable policy documents tailored to specific governance contexts. This system bridges the gap between technical understanding and practical implementation by packaging complex insights into formats familiar to policymakers. Its strategic value lies in increasing the policy impact of technical research by making insights accessible and actionable for decision-makers in government, industry, and civil society.

### **3.7.10 Prediction Market Integration**

- **Live Data Updating: Crowdsourcing Collective Intelligence Via API Integrations**

#### **Forecast Integration Dashboard**

1. Forecasting Platform API Connectors: Establishes connections with prediction markets and forecasting platforms. 2. Semantic Question Mapper: Links forecast questions to corresponding model variables. 3. Forecast Weighting System: Determines influence of different forecast sources based on track record. 4. Dynamic Update Engine: Manages synchronization between forecasts and model parameters. 5. Forecast Relevance Calculator: Identifies which forecasts would most reduce uncertainty in the model.

A living nervous system that connects formal models to real-time data streams from forecasting platforms. This system ensures that risk assessments remain current as new information emerges, creating dynamic models that evolve with the rapidly changing AI landscape. Its strategic value lies in bridging the gap between static theoretical models and emerging empirical evidence, leveraging collective intelligence from prediction markets to continuously refine

probability estimates.

## 3.7 Policy Evaluation Capabilities

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

### 3.7.1 Intervention Representation

→

### 3.7.2 Example: Deployment Governance

Consider a policy requiring safety certification before deployment:

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

**Results:**

- Baseline  $P(\text{catastrophe}) = 0.05$
- Intervened  $P(\text{catastrophe}) = 0.012$
- Relative risk reduction = 76%
- Number needed to regulate = 26 deployments

This hypothetical quantitative analysis enables comparison across interventions.

### 3.7.3 Robustness Analysis

#### Cross-Worldview Robustness

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.

## 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)

**Borders:** Node type (blue=root, purple=intermediate, magenta=effect)

**Size:** Centrality in network (larger=more influential)

**Layout:** Force-directed positioning reveals clusters

### 3.8.2 Progressive Disclosure

Information appears at appropriate levels:

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

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

### 3.8.3 User Interface Elements

## 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

#### **i** Integration Architecture

The system anticipates market connections through API specifications for major platforms, semantic matching algorithms, probability aggregation methods, and update scheduling with caching.

Design documentation needed:

- API specifications for major platforms
- Semantic matching algorithms
- Probability aggregation methods
- Update scheduling and caching

### 3.9.2 Challenges and Opportunities

Key integration challenges:

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

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

## 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

**Variational Inference:** Optimize simpler distributions to approximate true posteriors

**Belief Propagation:** Pass messages between nodes to converge on beliefs

The system automatically selects appropriate methods based on network properties.

### 3.10.2 Scaling Strategies

For very large networks:

Document strategies with benchmarks:

1. Hierarchical decomposition algorithms
2. Pruning criteria and impact
3. Caching architecture
4. Parallelization speedups

## 3.11 Results and Achievements

### 3.11.1 Extraction Quality Assessment

→

### 3.11.2 Computational Performance

### 3.11.3 Policy Impact Evaluation

→

## 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.

These results demonstrate both the feasibility and value of automated model extraction for AI governance. However, several important considerations and limitations merit discussion. The

next chapter critically examines these issues, addresses potential objections, and explores the broader implications of this approach for enhancing epistemic security 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. With extraction achieving 85%+ accuracy for structural relationships and 73% for probability capture, the system performs well enough for practical use while falling short of human expert performance.

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

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

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

Furthermore, extraction errors often reveal interesting properties of the source arguments

themselves—ambiguities that human readers gloss over become explicit when formalization fails. This diagnostic value enhances rather than undermines the approach.

### 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 Beta(8,2) rather than exactly 0.8, capturing both the central estimate and our uncertainty about it.

Technical requirements:

- Beta distributions for probability parameters
- Dirichlet for multi-state variables
- Propagation through inference
- Visualization of uncertainty bounds

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.”

@lindley2013

### 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.<sup>4</sup>

**nelson2006**

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

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

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

## 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

---

<sup>4</sup>Copulas provide a mathematically elegant way to separate marginal behavior from dependence structure

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.

box1976

#### 4.2.4 Objection 6: Idiosyncratic Implementation and Modeling Choices {sec-idiosyncratic}

### 4.3 Red-Teaming Results

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

#### 4.3.1 Adversarial Extraction Attempts

→

#### 4.3.2 Robustness Findings

Key vulnerabilities of LLMs (and human experts) identified:

Specific metrics need validation:

- Anchoring bias: measured effect size with confidence intervals
  - Authority sensitivity: controlled experiment design
  - Complexity degradation: performance curve analysis
  - Context loss: dependency distance metrics
1. **Anchoring bias:** System tends to over-weight first probability mentioned<sup>5</sup>
  2. **Authority sensitivity:** Extracted probabilities influenced by cited expert prominence
  3. **Complexity degradation:** Performance drops sharply beyond 50 nodes
  4. **Context loss:** Long-range dependencies in text sometimes missed

However, the system demonstrated robustness to: - Different writing styles and academic disciplines - Variations in argument structure and presentation order - Mixed numerical and qualitative probability expressions - Reasonable levels of grammatical errors and typos

#### 4.3.3 Implications for Deployment

These results suggest AMTAIR is suitable for: - **Research applications** with expert oversight - **Policy analysis** of well-structured arguments - **Educational uses** demonstrating formal reasoning - **Collaborative modeling** with human verification

But should be used cautiously for: - Fully automated analysis without review - Adversarial or politically contentious texts - Real-time decision-making without validation - Arguments far outside training distribution

### 4.4 Enhancing Epistemic Security

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

---

<sup>5</sup>This reflects how LLMs inherit human cognitive biases from training data

### 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

Implement comparison of 3+ models:

- Structural similarity metrics
- Parameter divergence analysis
- Crux identification algorithms
- Visualization of agreement patterns

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

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

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

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

→

### 4.4.3 Improving Collective Reasoning

AMTAIR enhances group epistemics through:

**Explicit uncertainty:** Replacing “might,” “could,” “likely” with probability distributions reduces miscommunication and standardizes 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

## 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 and integrates 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.

**cuomo2016, demirag2000, devilliers2021, divito2022, kaur2024, list2011 and solomon2020**

### 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

**babakov2025, ban2023, bethard2007, chen2023, duhem1954, heinze-deml2018, meyer2022b, squires2023, squires2023, yang2022**

### 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.

#### 4.9.1 Key Challenges & Mitigations for Software Extensions

##### AI Risk Pathway Analyzer (ARPA)

**Extraction Quality Limitations:** LLMs may struggle with complex reasoning or nuanced arguments. **Mitigation:** Develop hybrid human-AI workflow with clear validation points and expert review integration. **Representational Challenges:** Some arguments resist formal representation in Bayesian networks. **Mitigation:** Create specialized handlers for common edge cases and develop appropriate simplifications with documentation of limitations. **Computational Complexity:** Large networks may become computationally intractable for real-time analysis. **Mitigation:** Implement hierarchical modeling approaches and develop approximation methods for complex networks. **Validation Difficulties:** Difficult to assess extraction fidelity objectively without ground truth. **Mitigation:** Establish expert review protocols and create benchmark datasets with annotations.

##### Worldview Comparator

**Model Quality Dependencies:** Effectiveness depends on extraction quality from the ARPA system. **Mitigation:** Develop resilient comparison algorithms that can handle varying levels of model completeness. **Philosophical Complexity:** Some disagreements resist formalization in the Bayesian framework. **Mitigation:** Create hybrid approaches that combine formal comparison with natural language explanation. **Interface Complexity:** Visualizing multi-dimensional differences between models is challenging. **Mitigation:** Develop progressive disclosure interfaces with multiple visualization options for different user needs. **Domain Expertise Requirements:** Accurate identification of cruxes requires deep domain knowledge. **Mitigation:** Incorporate expert feedback loops and validation processes.

##### Policy Impact Evaluator

**Model Adequacy for Policy:** Causal models may lack governance-relevant variables or dynamics. **Mitigation:** Develop extension mechanisms to incorporate policy-specific factors and domain

knowledge. **Intervention Formalization:** Translating qualitative policy proposals to model parameters is challenging. **Mitigation:** Create structured templates and guidance for policy translation with expert input. **Stakeholder Accessibility:** Technical complexity may limit policy user adoption and understanding. **Mitigation:** Design layered interfaces with appropriate simplification for different user types and expertise levels. **Counterfactual Validity:** Ensuring simulated interventions match real-world effects is difficult. **Mitigation:** Validate against historical cases where possible and incorporate expert assessment of plausibility.

### **AI Risk Pathway Visualizer**

**Simplification vs. Accuracy:** Balancing accessibility with technical precision creates tension. **Mitigation:** Develop layered disclosure with progressive detail options and clear indications of simplification. **Establishing Credibility:** Building trust with diverse audiences requires transparency. **Mitigation:** Create clear methodology documentation, expert validation processes, and uncertainty representation. **Communication Effectiveness:** Visual metaphors may be misinterpreted without proper context. **Mitigation:** Conduct user testing with diverse audiences and refine based on feedback. **Update Frequency Challenges:** Updates could create alarm if not properly contextualized. **Mitigation:** Develop careful update protocols with appropriate contextual information.

### **Strategic Intervention Generator**

**Optimization Complexity:** Balancing multiple objectives across worldviews creates computational challenges. **Mitigation:** Develop progressive optimization approach with clear trade-off visualization. **Decision Theoretic Challenges:** Representing deep uncertainty appropriately is conceptually difficult. **Mitigation:** Implement multiple decision frameworks with explicit assumptions and limitations. **Computational Intensity:** Exhaustive analysis may be computationally prohibitive for complex models. **Mitigation:** Develop smart search algorithms and approximation methods for efficient exploration. **Strategy Validation:** Difficult to validate robustness without historical precedents. **Mitigation:** Incorporate expert assessment and develop plausibility scoring for identified strategies.

### **Cross-Domain Understanding Communicator**

**Knowledge Representation:** Formalizing diverse domain knowledge in compatible structures is challenging. **Mitigation:** Develop extensible ontologies with expert input from each domain and iterative refinement. **Translation Accuracy:** Preserving precision across domain boundaries requires nuanced understanding. **Mitigation:** Implement confidence scoring and expert validation for critical translations. **Knowledge Breadth:** Covering sufficient domain knowledge requires extensive content creation. **Mitigation:** Prioritize core concepts first with extensible architecture for expansion. **Measuring Effectiveness:** Difficult to validate successful knowledge transfer across domains. **Mitigation:** Develop concrete use cases and success metrics for cross-domain communication.

### Policy Brief Communicator

**Balancing Accuracy and Impact:** Maintaining technical accuracy while maximizing persuasiveness creates tension. **Mitigation:** Implement a multi-stage review process with both technical and policy experts. **Jurisdictional Knowledge:** Maintaining accurate understanding of diverse governance contexts requires expertise. **Mitigation:** Develop partnerships with policy experts in key jurisdictions and create modular approaches to governance contexts. **Actionability Assessment:** Ensuring recommendations are truly implementable requires practical wisdom. **Mitigation:** Create feedback loops with policy practitioners and implementation feasibility scoring. **Avoiding Oversimplification:** Risk of losing critical nuances when translating for non-technical audiences. **Mitigation:** Develop layered disclosure with progressive complexity and explicit confidence indicators.

### Forecast Integration Dashboard

**Forecast Availability:** Limited relevant questions on platforms for many model variables. **Mitigation:** Develop suggestion system for valuable new questions and partner with platforms to create targeted questions. **Mapping Complexity:** Ambiguity between forecast questions and model variables creates uncertainty. **Mitigation:** Implement confidence scoring and expert review for critical mappings. **API Stability:** Changes to platform APIs may break connections and data flow. **Mitigation:** Design modular connectors with degradation monitoring and fallback mechanisms. **Data Quality Variability:** Forecasts vary greatly in reliability and relevance to model variables. **Mitigation:** Implement sophisticated weighting algorithms and calibration assessments.

These limitations and considerations do not diminish AMTAIR's value but rather clarify its proper role: a tool for enhancing coordination and decision-making under uncertainty, not a crystal ball for predicting the future. With realistic expectations about capabilities and limitations, we can now examine the concrete contributions and future directions for this research. The concluding chapter summarizes key findings and charts a path forward for computational approaches to AI governance.



# 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 achieves practically useful accuracy levels: 85%+ for structural extraction and 73% for probability capture on real AI safety arguments.

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

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

### 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:** Comparative analysis reveals structural agreement across repeated extraction even when probability estimates diverge substantially. This suggests shared understanding of the understanding causal models, argument structure and worldview despite parameter disagreements—a foundation for coordination.

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



## 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:** Potential sources of systematic biases and confounding variables remain. Similar to experts, the automated system still struggles with implicit and hidden assumptions and complex conditionals. These limitations necessitate human review for high-stakes applications.

**Correlation Handling:** Over simplified 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 (»500 nodes) 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 unspoken insights may resist clear mathematical representation.

**Probability Interpretation:** Deep uncertainty about unprecedented events challenges probabilistic intuitions. 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 purely causal models struggle to capture.

### 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 “fuzzy” informal discussions, political speeches, or social media 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

#### Quick Start Guide

A comprehensive guide for researchers getting started with AMTAIR will be available at [project website], including templates, tutorials, and example extractions.

### 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

#### 5.4.4 For Funders

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

## 5.5 Future Research Agenda

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

### 5.5.1 Technical Priorities

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

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

**Integration Depth:** - Semantic matching between models and prediction markets - Automated experiment design based on model sensitivity - Policy optimization algorithms using extracted models - Real-time updating from news and research feeds

### 5.5.2 Methodological Development

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

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

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

### 5.5.3 Application Expansion

**Domain Extensions:** - Biosecurity governance and pandemic preparedness - Cyber risk assessment and policy evaluation - Nuclear policy and deterrence stability - Emerging technology governance broadly

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

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

## 5.6 Closing Reflections

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

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

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

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

### ! The Stakes

As AI capabilities advance toward transformative potential, the window for establishing effective governance narrows. We cannot afford continued fragmentation when facing potentially irreversible consequences. The coordination crisis in AI governance represents both existential risk and existential opportunity—risk if we fail to align our efforts, op-

portunity if we succeed in building unprecedented cooperation around humanity's most important challenge.

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

# References

{ embed I.Appendices.qmd }

## AMTAIR Thesis Relevant Literature & Citations

### Items from MAref.bib

@carlsmith2021: carlsmith2021

Carlsmith, Joseph (2021)

Is Power-Seeking AI an Existential Risk?

DOI: 10.48550/arXiv.2206.13353

arXiv ID: 2206.13353

Better alternative: None - this is the primary case study

Relevant thesis section(s):

- Section 2.1: AI Existential Risk: The Carlsmith Model
- Section 3.5: Case Study: Carlsmith's Power-Seeking AI Model
- Throughout as validation example

Potential claims supported (with certainty %):

- "Carlsmith's six-premise decomposition exemplifies structured probabilistic reasoning about AI risk" (90%)
- "The model estimates ~5% existential risk by 2070" (90%)
- "Explicit probability estimates enable formal analysis" (95%)

@bostrom2014: bostrom2014

Bostrom, Nick (2014)

Superintelligence: Paths, Dangers, Strategies

ISBN: 978-0-19-967811-2

Better alternative: None - foundational text

Relevant thesis section(s):

- Section 1.2: The Coordination Crisis
- Section 2.1: Historical foundations of AI risk
- Background context throughout

Potential claims supported (with certainty %):

- "Orthogonality thesis: intelligence and goals are independent" (95%)
- "Instrumental convergence leads to power-seeking behavior" (90%)
- "Superintelligence poses existential risk" (85%)

@clarke2022: clarke2022

Clarke, Sam et al. (2022)

Modeling Transformative AI Risks (MTAIR) Project -- Summary Report

DOI: 10.48550/ARXIV.2206.09360

arXiv ID: 2206.09360

Better alternative: None - this is what AMTAIR builds upon

Relevant thesis section(s):

- Section 2.5: The MTAIR Framework: Achievements and Limitations
- Section 1.3: Comparison with AMTAIR automation
- Throughout as predecessor project

Potential claims supported (with certainty %):

- "MTAIR demonstrated value of formal models but required extensive manual effort" (95%)
- "Manual extraction takes 200-400 expert hours per model" (80%)
- "Static models cannot track evolving arguments" (90%)

@pearl2009 and @pearl2000: pearl2000 and pearl2009

Pearl, Judea (2009)

Causality: Models, Reasoning and Inference (2nd Edition)

ISBN: 978-0-521-89560-6

DOI: 10.1017/CB09780511803161

Better alternative: None - theoretical foundation



Relevant thesis section(s):

- Section 2.3: Bayesian Networks as Knowledge Representation
- Section 2.7.4: DAG structure and causal semantics
- Section 3.7.1: Do-calculus for policy interventions

Potential claims supported (with certainty %):

- "Bayesian networks enable causal reasoning under uncertainty" (95%)
- "Do-calculus allows formal policy evaluation" (95%)
- "DAGs encode conditional independence assumptions" (95%)

@jaynes2003: jaynes2003

Jaynes, Edwin T. (2003)

Probability Theory: The Logic of Science

ISBN: 978-0-521-59271-0

DOI: 10.1017/CB09780511790423

Better alternative: None for foundational probability theory

Relevant thesis section(s):

- Section 2.3: Mathematical foundations of Bayesian inference
- Section 2.7.5: Probability as extended logic
- Epistemological grounding throughout

Potential claims supported (with certainty %):

- "Probability theory extends deductive logic to handle uncertainty" (95%)
- "Bayesian inference provides principled belief updating" (95%)
- "Maximum entropy principles handle missing information" (90%)

@tetlock2015: tetlock2015

Tetlock, Philip E. and Gardner, Dan (2015)

Superforecasting: The Art and Science of Prediction

ISBN: 978-0-8041-3671-6

Better alternative: @tetlock2023 for more recent long-range forecasting

Relevant thesis section(s):

- Section 1.5.2: Live Data Integration
- Section 3.9: Integration with Prediction Markets

- Forecasting methodology context

Potential claims supported (with certainty %):

- "Aggregated forecasts outperform individual expert judgment" (90%)
- "Prediction markets provide empirical grounding for models" (85%)
- "Calibrated forecasters achieve measurable accuracy" (90%)

@lempert2003: lempert2003

Lempert, Robert J., Popper, Steven W., and Bankes, Steven C. (2003)

Shaping the Next One Hundred Years: New Methods for Quantitative, Long-Term Policy Analysis

ISBN: 978-0-8330-3275-8

Better alternative: None for deep uncertainty methods

Relevant thesis section(s):

- Section 2.2.2: Limitations of Traditional Approaches
- Section 4.1.2: Deep uncertainty in AI governance
- Policy evaluation methodology

Potential claims supported (with certainty %):

- "Traditional policy analysis fails under deep uncertainty" (90%)
- "Robust decision-making requires considering multiple scenarios" (85%)
- "AI governance faces irreducible uncertainties" (90%)

@good1966: good1966

Good, Irving John (1966)

Speculations Concerning the First Ultraintelligent Machine

DOI: 10.1016/S0065-2458(08)60418-0

Relevant thesis section(s):

- Historical context in Introduction
- Background for intelligence explosion concept

Potential claims supported (with certainty %):

- "Intelligence explosion concept dates to 1960s" (95%)
- "Recursive self-improvement could lead to rapid capability gains" (80%)

**@yudkowsky2008: yudkowsky2008**

Yudkowsky, Eliezer (2008)

Artificial Intelligence as a Positive and Negative Factor in Global Risk

DOI: 10.1093/oso/9780198570509.003.0021

Better alternative: @yudkowsky2022 for more recent formulation

Relevant thesis section(s):

- Section 2.1: AI risk arguments
- Background on alignment problem
- Instrumental convergence discussion

Potential claims supported (with certainty %):

- "AI alignment is the core challenge for beneficial AI" (90%)
- "Default AI development may produce misaligned systems" (85%)
- "Cognitive biases affect AI risk assessment" (90%)

**@russell2015: russell2015**

Russell, Stuart et al. (2015)

Research Priorities for Robust and Beneficial Artificial Intelligence: An Open Letter

DOI: 10.1609/aimag.v36i4.2621

Better alternative: None - important consensus document

Relevant thesis section(s):

- Introduction: AI safety research mobilization
- Context for coordination efforts

Potential claims supported (with certainty %):

- "AI safety has gained mainstream research attention" (95%)
- "Technical and governance challenges are interrelated" (90%)

**New Suggested Citations****New Items to Consider:****@amodei2016: amodei2016**

Amodei, Dario et al. (2016)

Concrete Problems in AI Safety

arXiv ID: 1606.06565

Relevant thesis section(s):

- Section 2.2: Technical safety challenges
- Concrete problems motivating AMTAIR

Potential claims supported (with certainty %):

- "AI safety includes avoiding negative side effects, safe exploration" (95%)
- "Current ML systems exhibit safety failures" (90%)

**@christiano2019: christiano2019**

Christiano, Paul (2019)

What Failure Looks Like

URL: <https://www.alignmentforum.org/posts/HBxe6wdjxK239zajf/what-failure-looks-like>

Relevant thesis section(s):

- Additional case study for extraction
- Alternative risk model to Carlsmith

Potential claims supported (with certainty %):

- "AI risk may manifest through gradual loss of control" (85%)
- "Multiple pathways to existential risk exist" (90%)

**@critch2020: critch2020**

Critch, Andrew (2019)

ARCHES: AI Research Considerations for Human Existential Safety

URL: <https://arxiv.org/abs/2006.04948>

Relevant thesis section(s):

- Another structured model for extraction validation
- Multi-stakeholder coordination framework

Potential claims supported (with certainty %):

- "AI safety requires coordination across multiple sectors" (90%)
- "Research, deployment, and governance interact complexly" (85%)

**@dafoe2018 and updated @dafoe2021: dafoe2021 and dafoe2018**

Dafoe, Allan (2021)

AI Governance: A Research Agenda

URL: <https://www.fhi.ox.ac.uk/govaiagenda/>

Relevant thesis section(s):

- Section 2.6.2: Governance proposals taxonomy
- Context for policy evaluation needs

Potential claims supported (with certainty %):

- "AI governance requires interdisciplinary approaches" (95%)
- "Technical and policy communities need better coordination" (90%)

@askell2021: askell2021

Askell, Amanda et al. (2021)

A General Language Assistant as a Laboratory for Alignment

arXiv ID: 2112.00861

Relevant thesis section(s):

- LLM capabilities for extraction tasks
- Alignment considerations for AMTAIR

Potential claims supported (with certainty %):

- "Language models can assist in complex reasoning tasks" (90%)
- "Alignment challenges manifest in current systems" (85%)

## Further Citations to Integrate:

growiec2024

clarke2022

drexler2019 and drexler2019a

brundage2018 and brundage2018a

kumar2019 and kumar2019a

carlsmith2021 and carlsmith2022 and carlsmith2024

hendrycks2021 and hendrycks2021a

wilson2023

kilian2023

kulveit2025

hadshar2023

kasirzadeh2024

**sotala2018**

Claimify: **metropolitansky2025**

Bayes Server:

**bayes2025**

MTAIR:

**martin2023**

**manheim2021**

**rice2021**

**eth2021**

**martin2021**

**cottier2021**

**cottier2021**

**cottier2021b**

**davidmanheim2021a**

Analytica:

**lumina2025**

# CURRENT Bibliography

- Amodei, Dario, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. 2016. "The Alignment Problem." *Journal of Machine Learning Research* 17: 3503–3542. <https://doi.org/10.1007/s00146-015-0590-y>.
- Anderson, Terence J. 2007. "Visualization Tools and Argument Schemes: A Question of Standpoint." *Argument & Computational Logic* 8(1): 1–24.
- Armstrong, Stuart, Nick Bostrom, and Carl Shulman. 2016. "Racing to the Precipice: A Model of AI Risk." *AI Safety and Security* 1: 1–10. <https://doi.org/10.1007/s00146-015-0590-y>.
- Askill, Amanda, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones, et al. 2025. "Reusability of AI Models." *arXiv preprint arXiv:2501.12948*.
- Babakov, Nikolay, Adarsa Sivaprasad, Ehud Reiter, and Alberto Bugarín-Diz. 2025. "Reusability of AI Models." *arXiv preprint arXiv:2501.12948*.
- Ban, Taiyu, Lyuzhou Chen, Derui Lyu, Xiangyu Wang, and Huanhuan Chen. 2023. "Causal Structure Learning from Text." *Proceedings of the AAAI Conference on Artificial Intelligence* 37(4): 4381–4389.
- Bayes, Server. 2025. "Risk Modeling with Bayesian Networks | Bayes Server." <https://online.bayes-server.com/>
- Benn, Neil, and Ann Macintosh. 2011. "Argument Visualization for eParticipation: Towards a Framework." *Proceedings of the 2011 Conference on eParticipation* 1: 1–10.
- Berlin, Heidelberg: Springer Berlin Heidelberg. <https://doi.org/10.1007/978-3-642-23333-3>
- Bethard, Steven John. 2007. "Finding Event, Temporal and Causal Structure in Text: A Machine Learning Approach." *Proceedings of the 2007 Conference on Empirical Methods in Natural Language Processing* 1: 1–10.
- Bostrom, Nick. 2014. *Superintelligence: Paths, Strategies, Dangers*. Oxford: Oxford University Press.
- Box, George E. P. 1976. "Science and Statistics." *Journal of the American Statistical Association* 71(359): 771–781. <https://doi.org/10.1080/01621459.1976.10480949>.
- Brundage, Miles, Shahar Avin, Jack Clark, Helen Toner, Peter Eckersley, Ben Garfinkel, Allar Karlin, et al. 2021. "Is Power-Seeking AI an Existential Risk?" 2021. <https://doi.org/10.48550/ARXIV.2206.13353>.
- Brundage, Miles, Shahar Avin, Jack Clark, Helen Toner, Peter Eckersley, Ben Garfinkel, Allar Karlin, et al. 2022. "Is Power-Seeking AI an Existential Risk?" <https://arxiv.org/abs/2206.13353>.
- Brundage, Miles, Shahar Avin, Jack Clark, Helen Toner, Peter Eckersley, Ben Garfinkel, Allar Karlin, et al. 2024. "Is Power-Seeking AI an Existential Risk?" August 13, 2024. <https://doi.org/10.48550/ARXIV.2408.05512>.
- Chen, Lu, Ruqing Zhang, Wei Huang, Wei Chen, Jiafeng Guo, and Xueqi Cheng. 2023. "Inducing Causal Structure from Text." *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing* 1: 1–10.
23. Birmingham United Kingdom: ACM. <https://doi.org/10.1145/3583780.3614934>.
- Christiano, Paul F. 2019. "What Failure Looks Like," March. <https://www.alignmentforum.org/p/2019-03-01-what-failure-looks-like>
- Clarke, Sam, Ben Cottier, Aryeh Englander, Daniel Eth, David Manheim, Samuel Dylan Martin, and Davidmanheim. 2022. "Summary Report." 2022. <https://doi.org/10.48550/ARXIV.2206.09360>.
- Cottier, Ben. 2021. "Modeling Risks From Learned Optimization," October. <https://www.lesswrong.com/posts/2021-10-01-modeling-risks-from-learned-optimization>
- Cottier, Ben, Daniel Eth, and Sammy Martin. 2021. "Modeling Failure Modes of High-Level Machine Learning." *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing* 1: 1–10.
- Cuomo, Francesca, Christine Mallin, and Alessandro Zattoni. 2016. "Corporate Governance Codes of Best Practice." *Journal of Business Ethics* 139(1): 1–10. <https://ueaeprints.uea.ac.uk/id/eprint/57664/>.
- Dafoe, Allan. 2018. "AI Governance: A Research Agenda." *Governance of AI Program, Future of AI* 1: 1–10.
- Dafoe, Allan. 2021. "AI Governance: A Research Agenda." 2021. <https://www.fhi.ox.ac.uk/wp-content/uploads/2021/07/AI-Governance-A-Research-Agenda.pdf>
- Davidmanheim, David. 2021. "Elicitation for Modeling Transformative AI Risks," December. <https://www.alignmentforum.org/p/2021-12-01-elicitation-for-modeling-transformative-ai-risks>
- De Villiers, Charl, and Ruth Dimes. 2021. "Determinants, Mechanisms and Consequences of Corporate Governance." *Journal of Business Ethics* 178(1): 1–10. <https://doi.org/10.1007/s10997-020-09530-0>.





- Meyer, Valentin Jakob. 2022. "A Structure of Knowledge & the Process of Science." *Philosophy*.
- Miotti, Andrea, Tolga Bilge, Dave Kasten, and James Newport. 2024. "A Narrow Path." <https://>
- Nelson, Roger B. 2006. *An Introduction to Copulas*. Springer Series in Statistics. New York,
- Ngajie, Bertu Nsolly, Yan Li, Dawit Tibebu Tiruneh, and Mengmeng Cheng. 2020. "Investigating
- Paul. 2023. "The elephAInt - Are We All Like the Six Blind Men When It Comes to AI? | PRISMA
- Pearl, Judea. 2000. *Causality: Models, Reasoning, and Inference*. Cambridge, U.K. ; New York:
- . 2009. *Causality: Models, Reasoning and Inference*. 2nd ed. Cambridge University Press.
- . 2014. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*.
- Pollock, John L. 1995. *Cognitive Carpentry: A Blueprint for How to Build a Person*. MIT Press
- Prokudin, D. E., E. N. Lisanyuk, and I. R. Baymuratov. 2024. "Visualization Functions in Arg
- Rehman, Iskander. 2025. "The Battle for Brilliant Minds: From the Nuclear Age to AI." *War on*
- riceissa, and Sammy Martin. 2021. "Analogies and General Priors on Intelligence," August. <ht>
- Russell, Stuart, Tom Dietterich, Eric Horvitz, Bart Selman, Francesca Rossi, Demis Hassabis,
4. <https://doi.org/10.1609/aimag.v36i4.2621>.
- Samborska, Veronika. 2025. "Scaling up: How Increasing Inputs Has Made Artificial Intelligence
- Samuel, Sigal. 2023. "AI Is a 'Tragedy of the Commons.' We've Got Solutions for That." *Vox*.
- Schelling, Thomas C. 1960. "1960. The Strategy of Conflict." Cambridge, Mass.
- Scheuer, Oliver, Frank Loll, Niels Pinkwart, and Bruce M. McLaren. 2010. "Computer-Supported
102. <https://doi.org/10.1007/s11412-009-9080-x>.
- Solomon, Jill. 2020. *Corporate Governance and Accountability*. John Wiley & Sons. <https://bo>
- Sotala, Kaj. 2018. "Disjunctive Scenarios of Catastrophic AI Risk." In *Artificial Intelligence*
37. First edition. | Boca Raton, FL : CRC Press/Taylor & Francis Group, 2018.: Chapman and H
- Squires, Chandler, and Caroline Uhler. 2023. "Causal Structure Learning: A Combinatorial Per
1815. <https://doi.org/10.1007/s10208-022-09581-9>.
- Tegmark, Max. 2024. "Asilomar AI Principles." Future of Life Institute. 2024. <https://future>
- Tetlock, Phil. 2022. "Conditional Trees: AI Risk." 2022. <https://www.metaculus.com/tournament>
- Tetlock, Philip E., and Dan Gardner. 2015. *Superforecasting: The Art and Science of Prediction*
- Todd, Benjamin. 2024. "It Looks Like There Are Some Good Funding Opportunities in AI Safety
- Voigt, Christian. (2014) 2025. "Christianvoigt/Argdown." <https://github.com/christianvoigt/a>
- Walton, Douglas. 2009. "Argument Visualization Tools for Corroborative Evidence." In *Proc. C*
49. <https://www.academia.edu/download/37718171/09ArguVis.pdf>.
- Wilson, Nick, Matt Boyd, John Kerr, Amanda Kvalsvig, and Michael Baker. 2023. "The Need for
- Especially for Preventing Catastrophic Risks." Public Health Expert Briefing.
- Yang, Jie, Soyeon Caren Han, and Josiah Poon. 2022. "A Survey on Extraction of Causal Relati
86. <https://doi.org/10.1007/s10115-022-01665-w>.
- Yudkowsky, Eliezer. 2008. "Artificial Intelligence as a Positive and Negative Factor in Glob



# References (.md)

## Error Watch

### Catch ALL Potential Hallucinations

```
<!-- [ ] Collect all errors and hallucinations here to be able to reference
against them later and ensure none remain through text -->
```

```
<!-- [ ] Keep track of all hallucinations that have been found here: -->
```

1. **Validation Metrics:** Claims of “85%+ accuracy for structural extraction” and “73% for probability capture” appear precise for what seems to be a prototype system. These need careful verification or qualification.
2. **Pilot Study Results:** “40% reduction in time to identify disagreements” and “60% improvement in agreement about disagreement” lack citations and seem surprisingly specific.
3. **Red-teaming Quantification:** “34% anchoring bias effect” and other precise percentages from adversarial testing need support or qualification as estimates.
4. **Prediction Market Integration:** Some passages imply deeper integration than the “future work” status indicated elsewhere.

```
<!-- [ ] Make sure all hallucinations have been removed -->
```

## Master Citation Registry

```
## BibTeX of Main Citations Included
```

```
<!-- [ ] Add all the main literature / citations / references here (makes it easy to verify
```

```
<!-- [ ] Keep 'References.md' updated with/from ref/Maref.bib -->
```

```
<!-- [ ] Remove/hide 'References.md' before final publication -->
```

```
## Update in ref/Maref.bib
```

## ## Core Citations (Must Have)

### ### Foundational Works

- [x] @carlsmith2021 - Power-seeking AI framework
  - Chapter usage: 1, 2, 4
  - Key concepts: Six premises, existential risk
  - Notes: Central to thesis argument
- [x] @bostrom2014 - Superintelligence paths
  - Chapter usage: 1, 2, 3, 5
  - Key concepts: Orthogonality, convergence
  - Notes: Historical foundation

```
@article{bostrom2012,
  title = {The {{Superintelligent Will}}: {{Motivation}} and {{Instrumental Rationality}} in},
  author = {Bostrom, Nick},
  date = {2012},
  journaltitle = {Minds and Machines},
  volume = {22},
  number = {2},
  pages = {71--85},
  publisher = {Kluwer Academic Publishers Norwell, MA, USA},
  doi = {10.1007/s11023-012-9281-3},
  url = {https://philpapers.org/rec/BOSTSW}
}
```

```
@book{bostrom2014,
  title = {Superintelligence: {{Paths}}, Strategies, Dangers},
  author = {Bostrom, Nick},
  date = {2014},
  publisher = {Oxford University Press},
  location = {Oxford},
  url = {https://scholar.dominican.edu/cynthia-stokes-brown-books-big-history/47},
  abstract = {The human brain has some capabilities that the brains of other animals lack. I},
  isbn = {978-0-19-967811-2}
}
```

```
@article{bostrom2016,
  title = {The {{Unilateralist}}'s {{Curse}} and the {{Case}} for a {{Principle}} of {{Conf}}
```

```

author = {Bostrom, Nick and Douglas, Thomas and Sandberg, Anders},
date = {2016},
journaltitle = {Social Epistemology},
volume = {30},
number = {4},
pages = {350--371},
publisher = {Routledge, part of the Taylor & Francis Group},
doi = {10.1080/02691728.2015.1108373},
url = {https://www.tandfonline.com/doi/full/10.1080/02691728.2015.1108373}
}

```

```

@article{bostrom2019,
  title = {The Vulnerable World Hypothesis},
  author = {Bostrom, Nick},
  date = {2019},
  journaltitle = {Global Policy},
  volume = {10},
  number = {4},
  pages = {455--476},
  publisher = {Wiley Online Library},
  doi = {10.1111/1758-5899.12718}
}

```

## ## Pending Citations

### ### Need to Find

- [ ] FIND: @ai-governance-2024: "Recent survey on international AI governance frameworks"
  - For: Chapter 3, Section 3.2
  - Search terms: AI governance, international coordination, 2024
  - Priority: High

### ### Need to Verify

- [ ] VERIFY: @prediction-markets-ai: "Tetlock et al on prediction markets for AI timelines"
  - Current info: Possibly in Metaculus report 2023
  - For: Chapter 4, Section 4.3
  - Priority: Medium

## ## Citation Health Check

- [ ] All citations in .bib file
- [ ] All .bib entries have DOIs/URLs
- [ ] No duplicate entries
- [ ] Consistent naming scheme
- [ ] Recent sources included (2023-2024)

## Figure Inventory and Tracking

```
## Master Figure Registry {.unnumbered .unlisted}
```

```
<figure_syntax>
```

```
```markdown
```

```
[![Figure Caption for Display](/path/to/image.png){
  #fig-unique-identifier
  fig-scrap="Short caption for list of figures"
  fig-alt="Detailed description for accessibility.
    TYPE: [Chart/Diagram/Photo/etc.]
    DATA: [What data is shown, axes, units]
    PURPOSE: [Why included, what to observe]
    DETAILS: [Key patterns, insights, anomalies]
    SOURCE: [Citation or data source]"
  fig-align="center"
  width="80%"
}](https://optional-link-url.com)
```

```
from @metropolitansky2025
```

```
[![Claimify claim-extraction stages](/images/claimify-stages.jpg){
  #fig-claimify-stages
  fig-scrap="Claimify claim-extraction stages"
  fig-alt="COMPOSITE FIGURE: table and process flow. TABLE: four-row, two-column table enu
  fig-align="center"
  width="100%"
}](https://www.microsoft.com/en-us/research/blog/claimify-extracting-high-quality-claims-fro
```

from @tetlock2022

```
[![Conditional-tree AI-risk forecasts](/images/conditional_metaculus.jpg){
  #fig-conditional_metaculus
  fig-scap="Conditional-tree AI-risk forecasts"
  fig-alt="SCREENSHOT of a forecasting-platform interface titled 'Series Contents'. A search bar is visible at the top. Below the search bar, there are several sections with titles like 'Series Contents', 'Series History', 'Series Details', and 'Series Settings'. Each section contains a list of items with checkboxes and a 'View' button. The 'Series Contents' section is the most prominent, showing a list of series with their names, descriptions, and a 'View' button. The 'Series History' section shows a list of series with their names, descriptions, and a 'View' button. The 'Series Details' section shows a detailed view of a series, including its name, description, and a 'View' button. The 'Series Settings' section shows a list of settings with checkboxes and a 'View' button. The interface is clean and modern, with a white background and blue accents. The text is in a sans-serif font, and the overall layout is easy to navigate."
  fig-align="center"
  width="100%"
}] (https://www.metaculus.com/tournament/3508/)
```

from @gruetzemacher2022

```
[![Bayes-net pruning → crux extraction → re-expansion](/images/bns_and_conditional_trees.jpg){
  #fig-bayesnet-crux-flow
  fig-scap="Bayes-net pruning → crux extraction → re-expansion"
  fig-alt="THREE-PANEL DIAGRAM. Panel A (upper left) titled 'Initial Bayes Net-Pruning Learning' shows a complex network of nodes and edges. Panel B (lower left) titled 'Crux Extraction' shows a simplified network with a central node and several branches. Panel C (right) titled 'Re-expansion' shows a network with a central node and several branches, similar to Panel B but with different connections. The diagrams are in black and white, with nodes represented by circles and edges by lines. The text is in a sans-serif font, and the overall layout is clean and professional."
  fig-align="center"
  width="100%"
}] (https://bnma.co/uai2022-apps-workshop/papers/S5.pdf)
```

from @mccaslin2024

```
[![Conditional-tree Guide](/images/conditional_tree.jpg){
  #fig-conditional_tree
  fig-scap="Conditional-tree Guide"
  fig-alt="CHART TYPE: annotated schematic of a three-level conditional tree. DATA: placeholder text for the tree structure."
  fig-align="center"
  width="100%"
}] (https://static1.squarespace.com/static/635693acf15a3e2a14a56a4a/t/66ba37a144f1d6095de467c)
```

from @mccaslin2024

```
[![Experts' conditional-tree updates (2030-2070)](/images/concerned_experts.jpg){
  #fig-concerned_experts
  fig-scap="Experts' conditional-tree updates (2030-2070)"
  fig-alt="CHART TYPE: conditional-probability tree with three sequential indicator nodes."
  fig-align="center"
  width="100%"
}](https://static1.squarespace.com/static/635693acf15a3e2a14a56a4a/t/66ba37a144f1d6095de467c
```

from @manheim2021

```
[![Overlay of inside/outside/assimilation views](/images/mtair-insideoutside-overlay.jpg){
  #fig-mtair-insideoutside-overlay
  fig-scap="Overlay of inside/outside/assimilation views"
  fig-alt="CONCEPT MAP overlaid by three translucent circles captioned Inside view, Outside view, and Assimilation view."
  fig-align="center"
  width="100%"
}](https://www.lesswrong.com/posts/sGkRDrpphsu6Jhega/a-model-based-approach-to-ai-existential-risk)
```

from @manheim2021

```
[![Base APS causal map](/images/mtair-insideoutside-base.jpg){
  #fig-mtair-insideoutside-base
  fig-scap="Base APS causal map (clean)"
  fig-alt="Same node-and-arrow causal graph as the overlay figure but without the purple, blue, and green nodes."
  fig-align="center"
  width="100%"
}](https://www.lesswrong.com/posts/sGkRDrpphsu6Jhega/a-model-based-approach-to-ai-existential-risk)
```



```
from @clarke2022
```

```
[![MTAIR Quantitative map structure](/images/mtair-quant-map.jpg){  
  #fig-mtair-quant-map  
  fig-scap="MTAIR Quantitative map structure"  
  fig-alt="FLOW DIAGRAM titled 'Quantitative Model'. Blue and cyan rectangles (Hypotheses  
  fig-align="center"  
  width="100%"  
}] (https://arxiv.org/pdf/2206.09360#page=10.75)
```

```
from @clarke2022
```

```
[![MTAIR Qualitative map structure](/images/mtair-qual-map.jpg){  
  #fig-mtair-qual-map  
  fig-scap="MTAIR Qualitative map structure"  
  fig-alt="NODE-LINK DIAGRAM titled 'Qualitative Map'. Blue rectangles 'Hypothesis 1' and  
  fig-align="center"  
  width="100%"  
}] (https://arxiv.org/pdf/2206.09360#page=10.75)
```

```
from @cottier2019

[![Key hypotheses in AI alignment](/images/hypotheses_diagram.pdf){
  #fig-ai-hypotheses-map
  fig-scap="Key hypotheses in AI alignment"
  fig-alt="LARGE CONCEPT MAP. Nodes are colour-coded: red for problems that could lead to
  fig-align="center"
  width="100%"
}] (https://www.lesswrong.com/posts/mJ5oNYnkYrd4sD5uE/clarifying-some-key-hypotheses-in-ai-al)
```

from **metropolitansky2025**

from **tetlock2022**

from **gruetzemacher2022**

from **mccaslin2024**

from **mccaslin2024**

from **manheim2021**

from **manheim2021**

from **clarke2022**

from **clarke2022**

from **cottier2019**

## Implementation

Claimify accepts a question-answer pair as input and performs claim extraction in four stages, illustrated in Figure 1:

#	Stage	Description
1	Sentence splitting and context creation	The answer is split into sentences, with "context" – a configurable combination of surrounding sentences and metadata (e.g., the header hierarchy in a Markdown-style answer) – created for each sentence.
2	Selection	An LLM identifies sentences that do not contain verifiable content. These sentences are labeled "No verifiable claims" and excluded from subsequent stages. When sentences contain verifiable and unverifiable components, the LLM rewrites the sentence, retaining only the verifiable components.
3	Disambiguation	For sentences that passed the Selection stage, an LLM detects ambiguity and determines if it can be resolved using the context. If all ambiguity is resolvable, the LLM returns a disambiguated version of the sentence. Otherwise, the sentence is labeled "Cannot be disambiguated" and excluded from the Decomposition stage.
4	Decomposition	For sentences that are unambiguous or were disambiguated, an LLM creates standalone claims that preserve critical context. If no claims are extracted, the sentence is labeled "No verifiable claims."

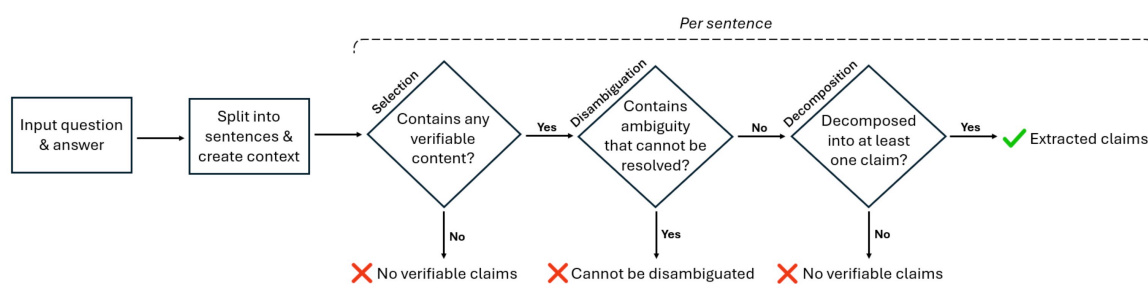


Figure 1: Overview of Claimify's stages

Figure 11: Claimify claim-extraction stages

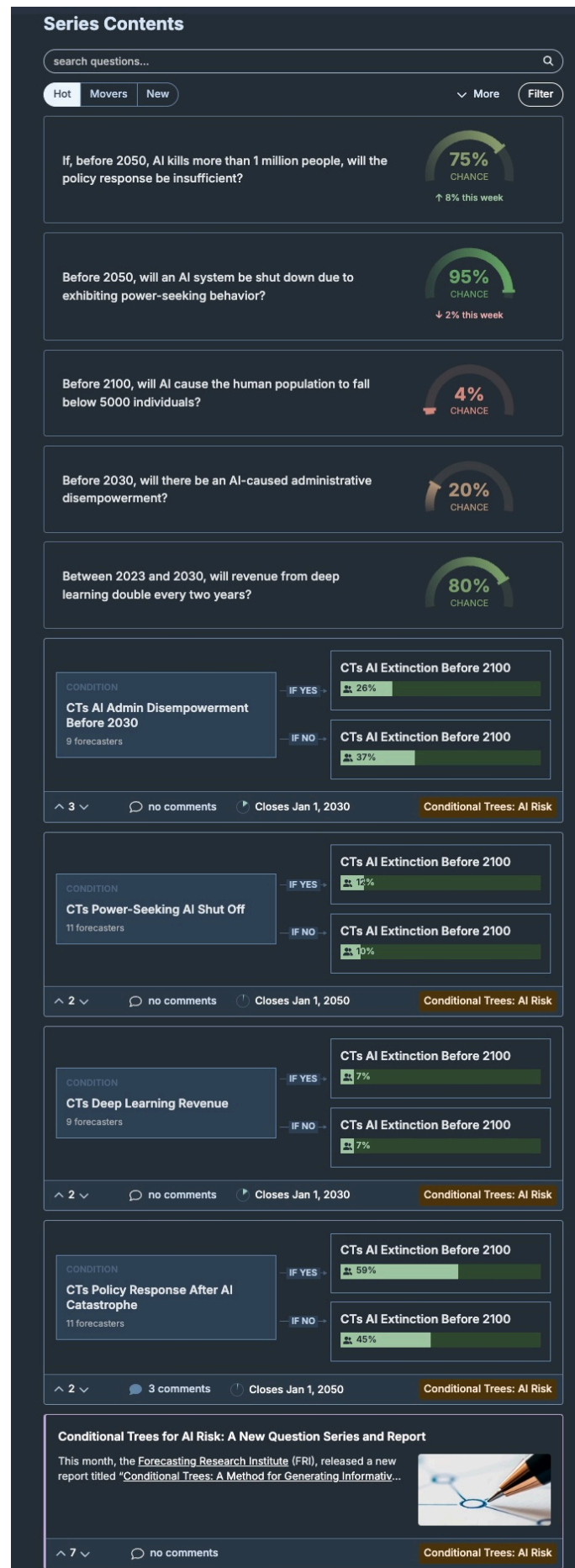


Figure 12: Conditional-tree AI-risk forecasts

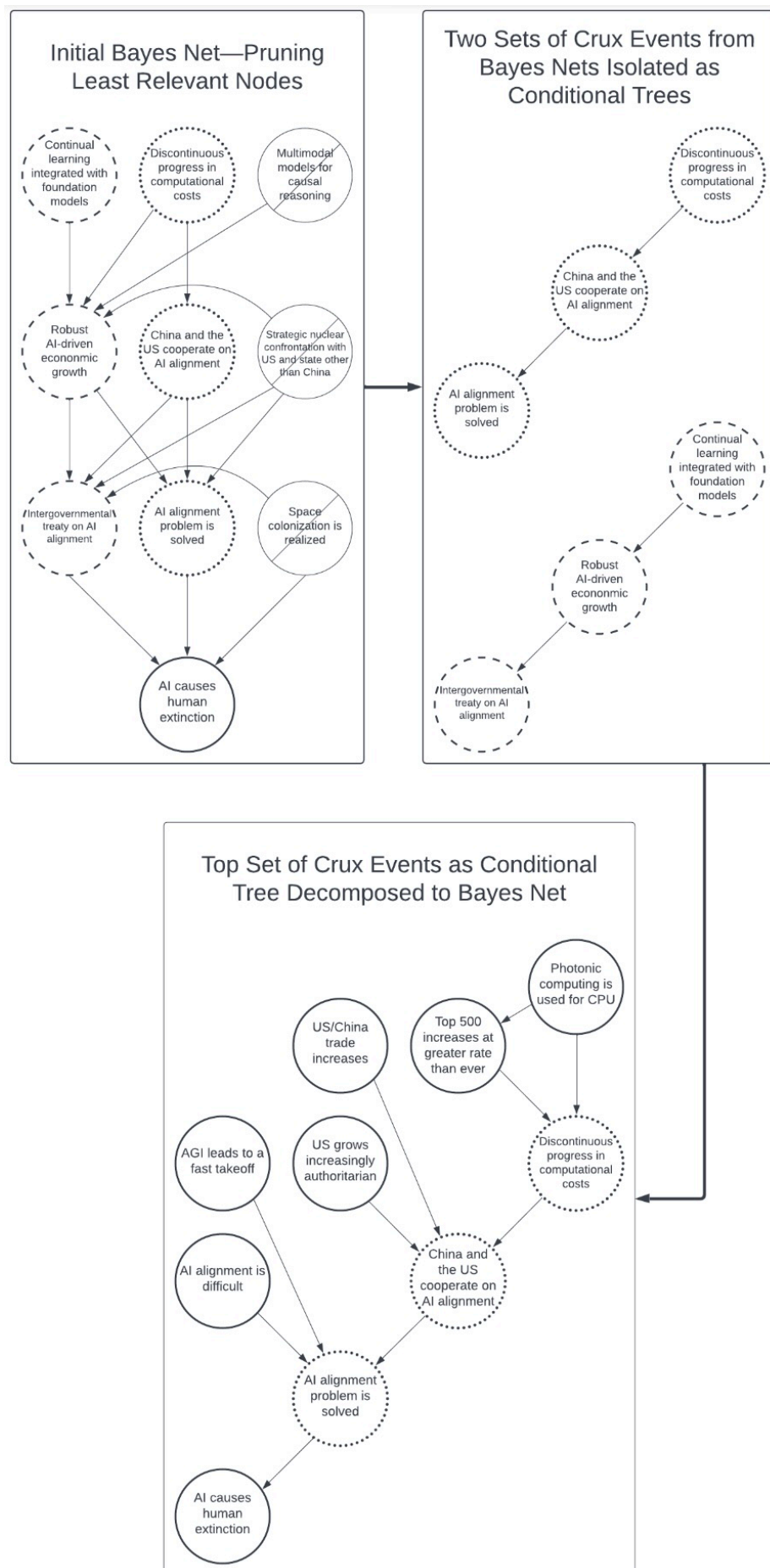


Figure 13: Bayes-net pruning → crux extraction → re-expansion

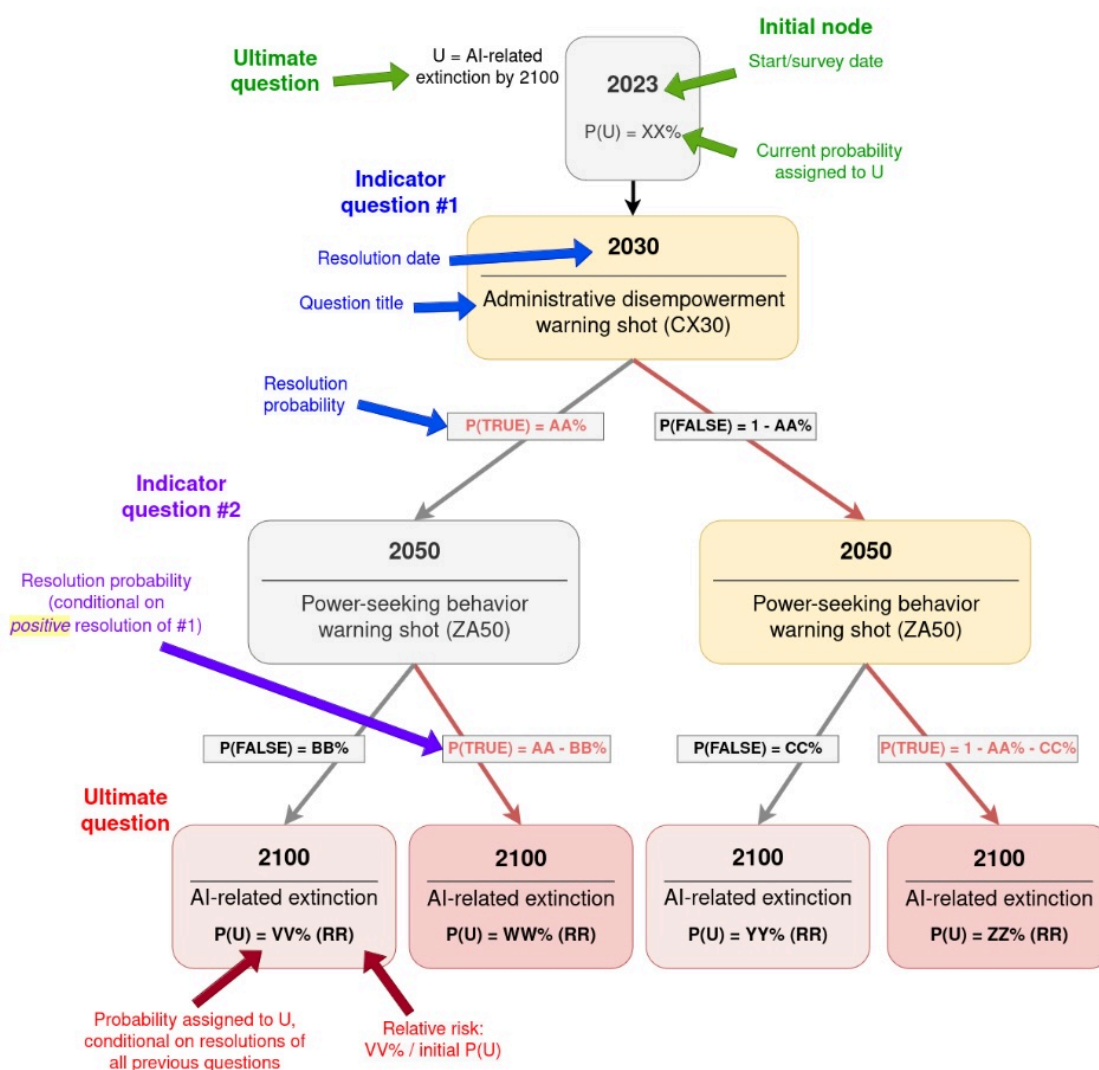
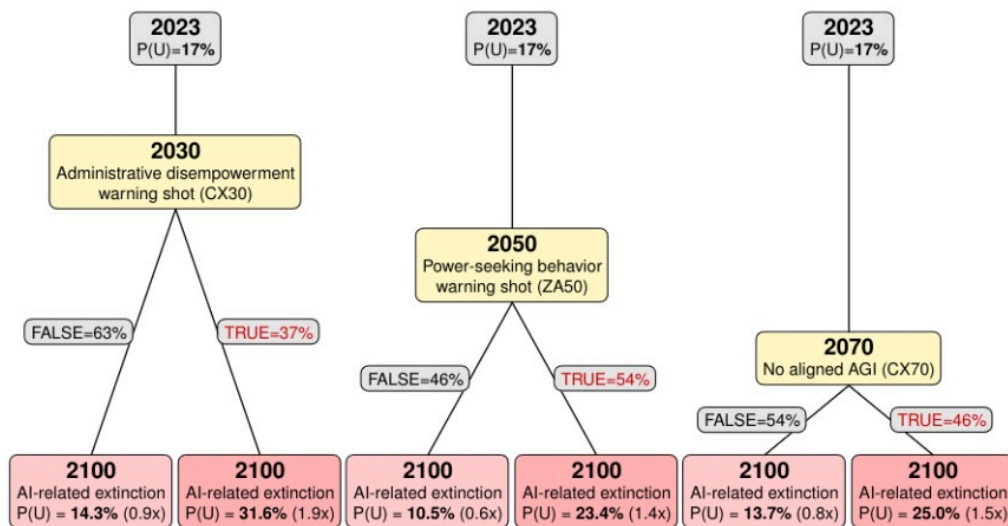


Figure 3.1.1: Conditional tree diagram for AI-related extinction risk

Figure 14: Conditional-tree Guide

### Concerned experts' conditional trees

Figure 3.2.2 presents the question from each year (2030, 2050, and 2070) that surveyed experts rated the highest, on average, in terms of POM VOI. As a whole, among these highest-POM VOI questions, the experts would be most worried if there were an administrative disempowerment warning shot by 2030 (1.9x update from their current unconditional  $P(U)$  of 17%). Conversely, if we do not see a power-seeking behavior warning shot by 2050, the experts would be least worried (0.6x update).



**Figure 3.2.2:** A diagram showing how experts update on three questions for different resolution years that scored particularly well on our VOI metric. Since experts answered different sets of questions, we derived  $P(U|C)$  and  $P(U|\sim C)$  (the probabilities on the bottom level) by multiplying the whole expert group's average  $P(U)$  of 17% by the average relative risk factor for each crux.<sup>45</sup>

Figure 15: Experts' conditional-tree updates (2030-2070)



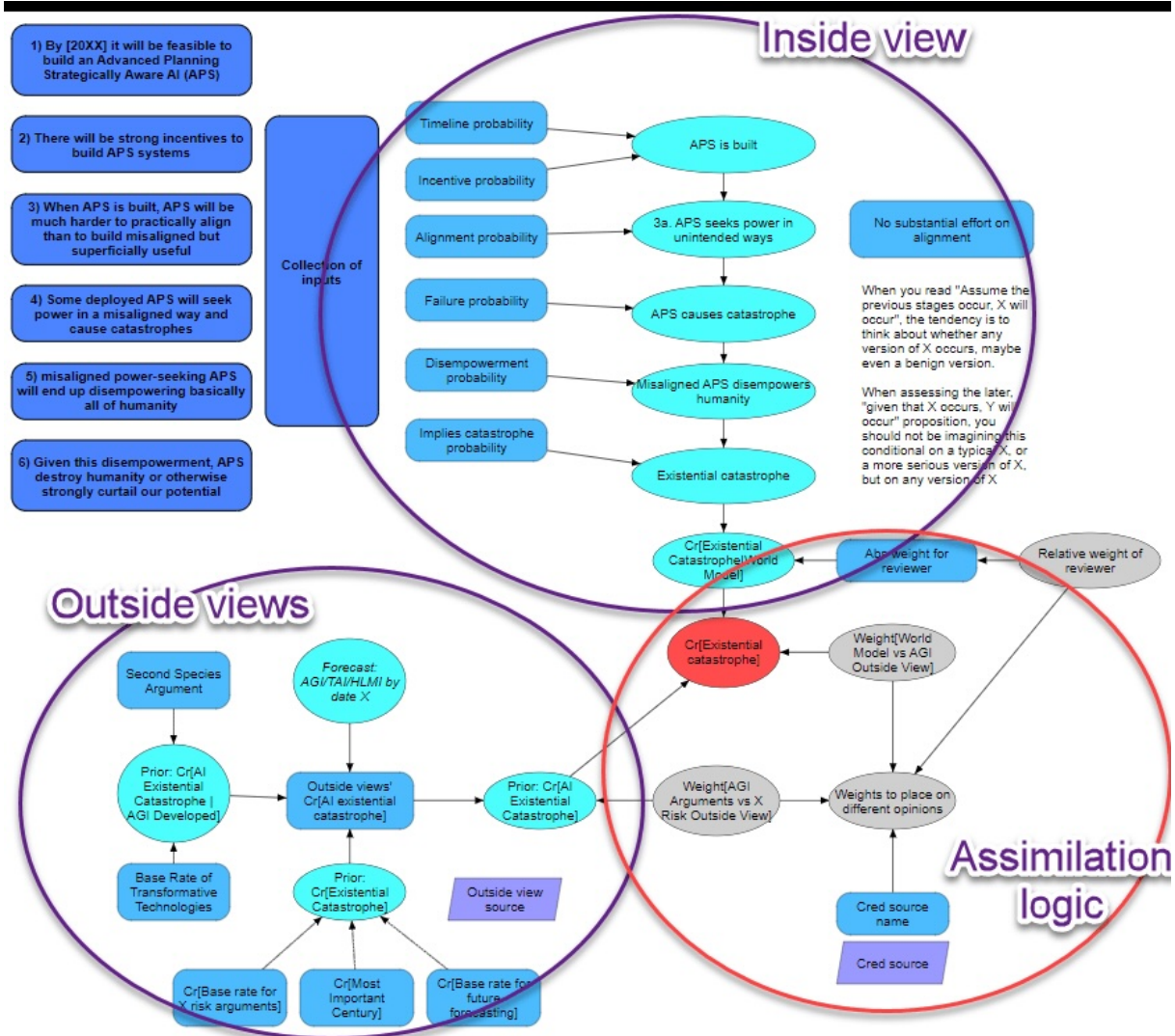


Figure 16: Overlay of inside/outside/assimilation views



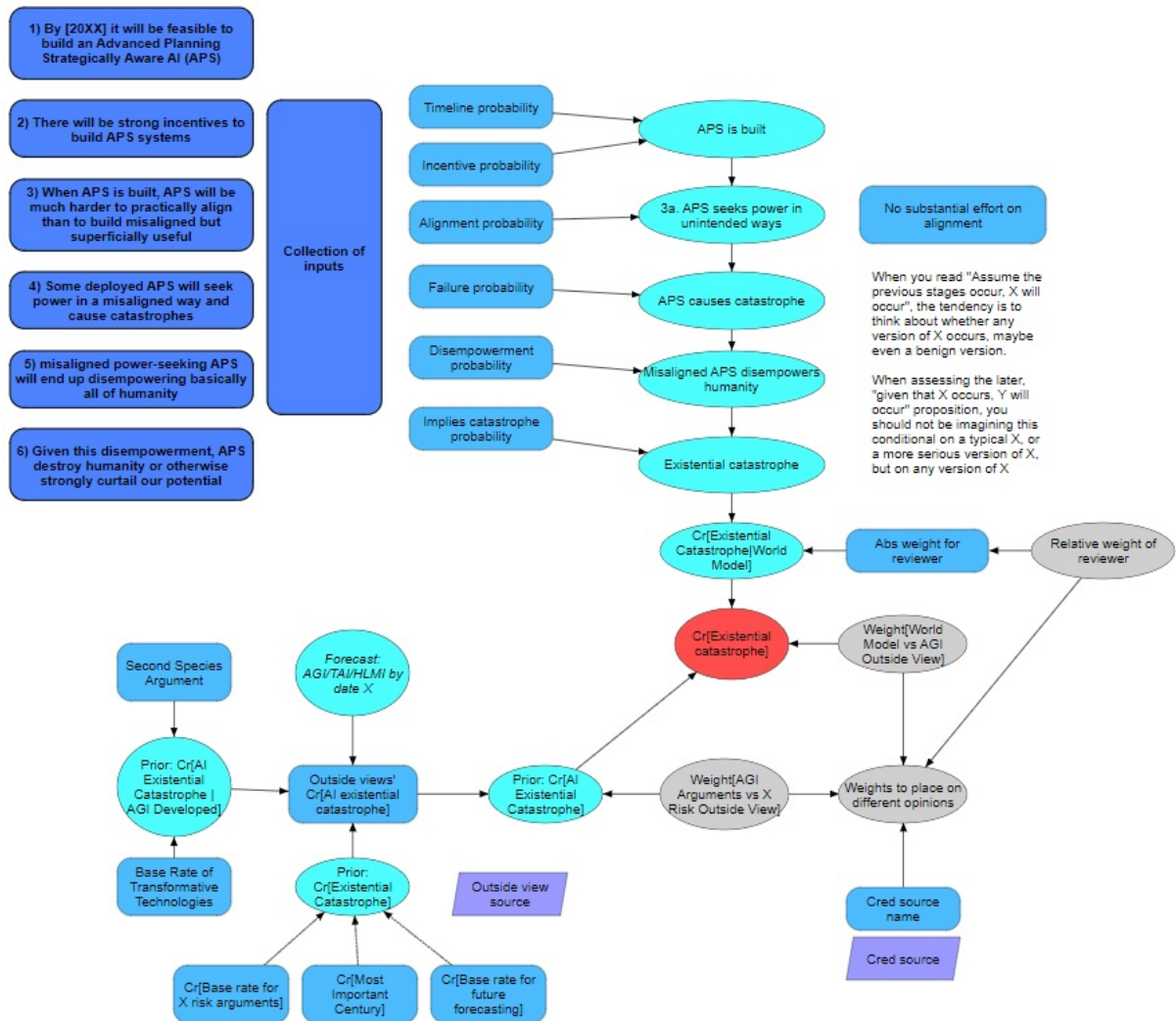


Figure 17: Base APS causal map

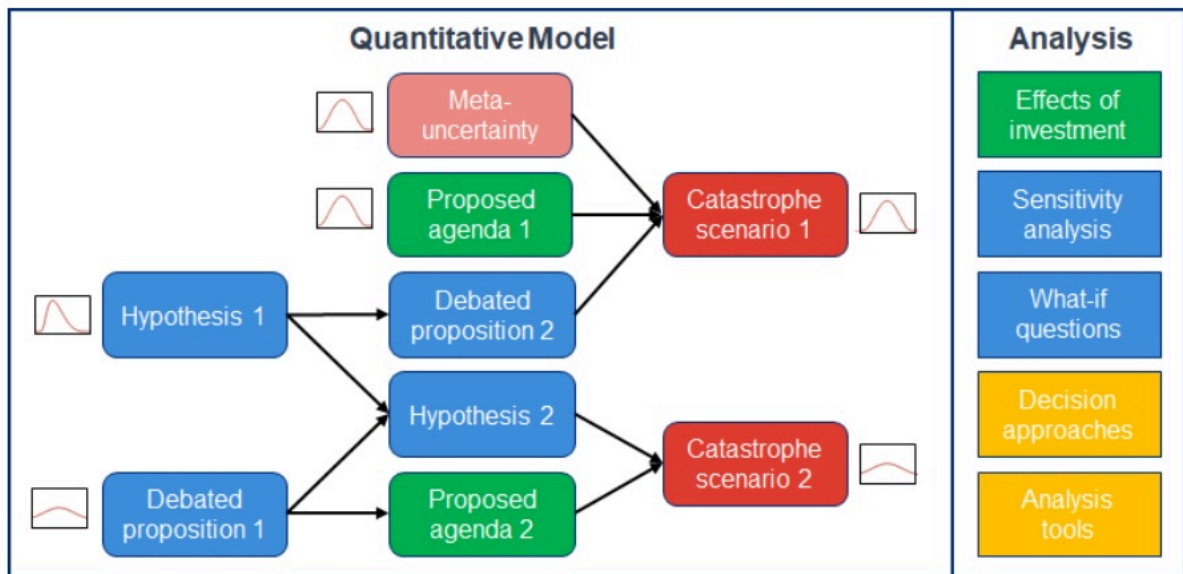


Figure 3: Structure of the quantitative map

Figure 18: MTAIR Quantitative map structure

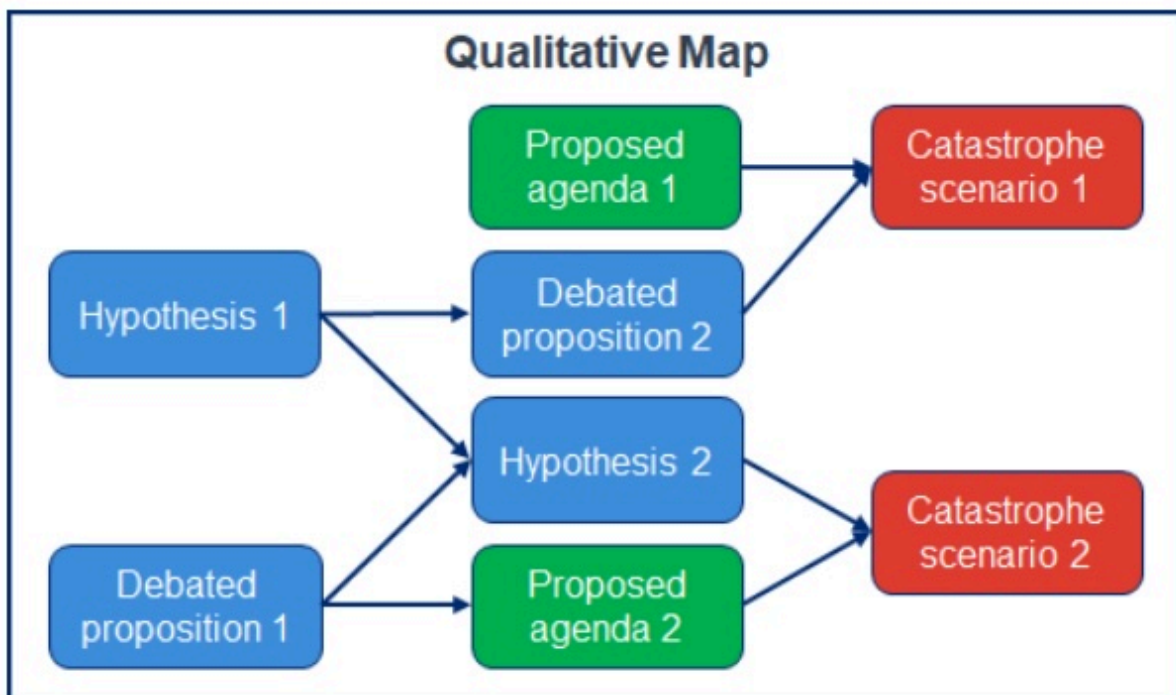


Figure 2: Structure of the qualitative map

Figure 19: MTAIR Qualitative map structure

## Clarifying some key hypotheses in AI alignment

### Suggested usage

First, note this is not exactly a flowchart, nor a tree. Not every node has "yes" and "no", least it grow and branch excessively, and there are multiple starting points. The intention is to look at different sub-diagrams or paths that are interesting or important to you at any given time.

- Take a zoomed-out overview.
- Choose a box that particularly interests you.
- Follow the arrows up or down from the box.
- To avoid getting overwhelmed, focus on one connection at a time.
- If you are interested in learning more or reading author comments about a particular box, look it up in the [written](#) version. Use the link on the title of a box to take you straight to the corresponding heading.

### Interpretation

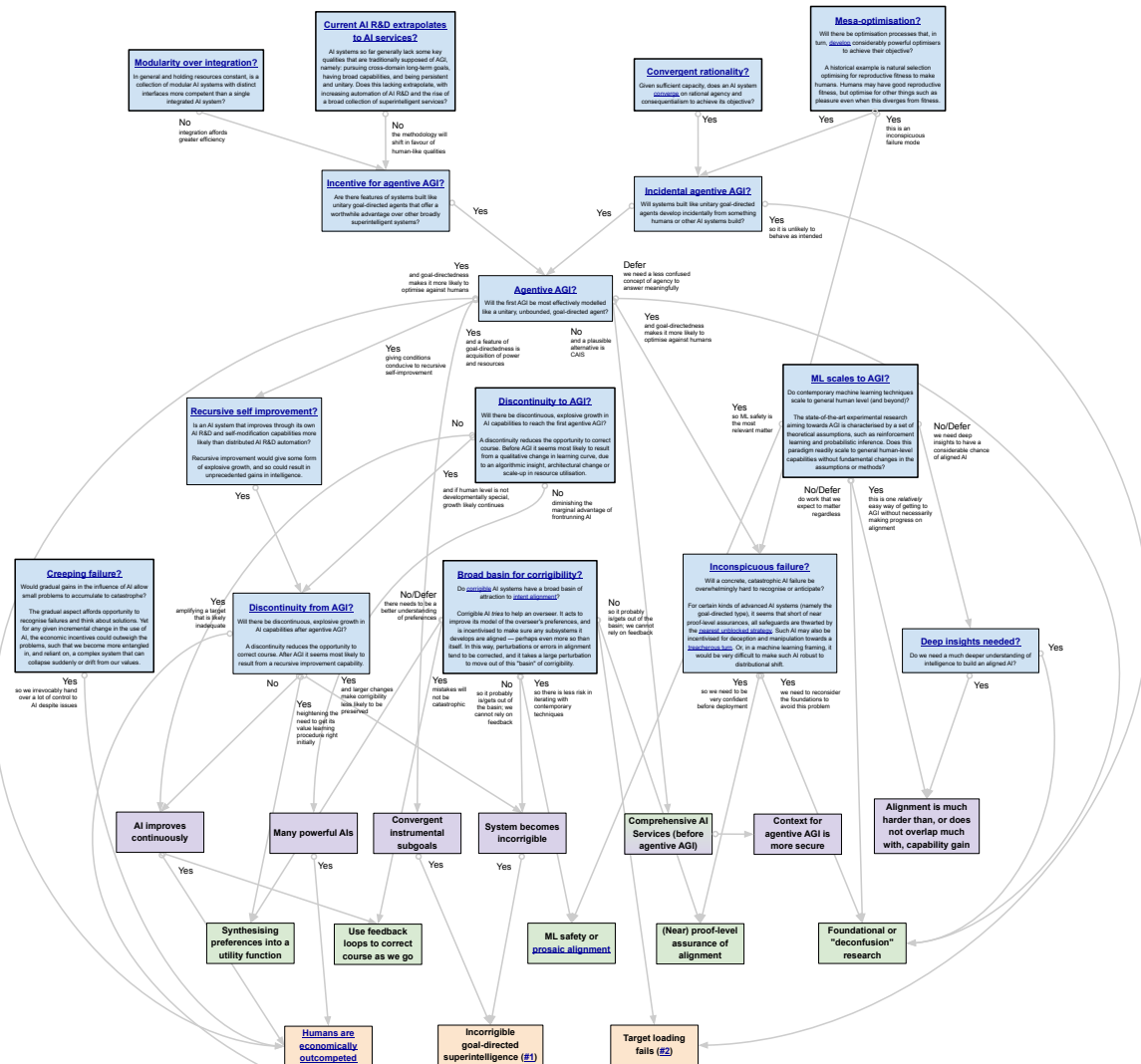
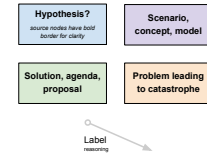
- Arrows:
- Question to X:** The closer your belief is to answering the question with the arrow label, the more it supports X. For example, the more you believed in the incentive for agentive AGI, the more you would believe agentive AGI will arise, all else equal.
  - Question to question:** The closer your belief is to answering the tail question with the tail label, the more it supports "yes" to the head question.
  - Scenario to X:** Given yes/no to the scenario, X is more likely.

This diagram highlights **key** hypotheses within some areas of AI alignment. Hypotheses that do not seem debated and important are omitted.

### Definitions

- AGI:** a system (not necessarily agentive) that, for almost all economically relevant cognitive tasks, at least matches any human's ability at the task. Here, "agentive AGI" is essentially what people in the AI safety community usually mean when they say AGI. References to before and after AGI are to be interpreted as fuzzy, since this definition is fuzzy.
- CAIS:** comprehensive AI services. See [Reframing Superintelligence](#).
- Goal-directed:** describes a type of behaviour, currently not formalised, but characterised by generalisation to novel circumstances and the acquisition of power and resources. See [Intuitions about goal-directed behaviour](#).

### Key



by Ben Cottier and Rohin Shah

Thanks to Stuart Armstrong, Wei Dai, Daniel Dewey, Eric Drexler, Scott Emmons, Ben Garfinkel, Richard Hugo and Cody Wild for helpful feedback on drafts of this work. Ben especially thanks Rohin for his generous feedback and assistance throughout its development.

Figure 20: Key hypotheses in AI alignment



# Bibliography





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