



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

Contents

| | |
|---|-----------|
| Preface | 7 |
| Abstract | 9 |
| Outline(s): Table of Contents | 11 |
| 1 Introduction | 13 |
| Abstract | 13 |
| 2 Introduction | 15 |
| 2.1 The Coordination Crisis in AI Governance | 15 |
| 2.1.1 Empirical Paradox: Investment Alongside Fragmentation | 15 |
| 2.1.2 Systematic Risk Increase Through Coordination Failure | 15 |
| 2.1.3 Historical Parallels and Temporal Urgency | 15 |
| 2.2 Research Question and Scope | 16 |
| 2.3 The Multiplicative Benefits Framework | 16 |
| 2.4 Thesis Structure and Roadmap | 17 |
| 2.5 Overview / Table of Contents | 18 |
| 3 Context | 19 |
| 4 Context & Background | 21 |
| 4.1 Theoretical Foundations | 21 |
| 4.1.1 AI Existential Risk: The Carlsmith Model | 21 |
| 4.1.2 The Epistemic Challenge of Policy Evaluation | 22 |
| 4.1.3 Argument Mapping and Formal Representations | 23 |
| 4.1.4 Bayesian Networks as Knowledge Representation | 23 |
| 4.1.5 Argument Mapping and Formal Representations | 25 |
| 4.1.6 The MTAIR Framework: Achievements and Limitations | 25 |
| 4.1.7 “A Narrow Path”: Conditional Policy Proposals in Practice | 26 |
| 4.2 Methodology | 27 |
| 4.2.1 Research Design Overview | 27 |
| 4.2.2 Formalizing World Models from AI Safety Literature | 28 |
| 4.2.3 From Natural Language to Computational Models | 28 |
| 4.2.4 Directed Acyclic Graphs: Structure and Semantics | 29 |
| 4.2.5 Quantification of Probabilistic Judgments | 30 |
| 4.2.6 Inference Techniques for Complex Networks | 31 |
| 4.2.7 Integration with Prediction Markets and Forecasting Platforms | 32 |
| 5 AMTAIR | 35 |
| 5.1 AMTAIR Implementation | 35 |
| 5.2 Software Implementation | 35 |
| 5.2.1 System Architecture and Data Flow | 35 |
| 5.2.2 Rain-Sprinkler-Grass Example Implementation | 37 |
| 5.2.3 Carlsmith Implementation | 39 |
| 5.2.4 Inference & Extensions | 41 |
| 5.3 Results | 44 |
| 5.3.1 Extraction Quality Assessment | 44 |

| | | |
|--|---|-----------|
| 5.3.2 | Computational Performance Analysis | 44 |
| 5.3.3 | Case Study: The Carlsmith Model Formalized | 45 |
| 5.3.4 | Comparative Analysis of AI Governance Worldviews | 47 |
| 5.3.5 | Policy Impact Evaluation: Proof of Concept | 48 |
| 6 | Discussion | 51 |
| 7 | Discussion — Exchange, Controversy & Influence | 53 |
| 7.1 | Limitations and Failure Modes | 53 |
| 7.1.1 | Limitations and Counterarguments | 53 |
| 7.1.2 | Technical Limitations | 53 |
| 7.1.3 | Integration with Existing Governance Frameworks | 54 |
| 7.2 | Red-Teaming Results: Identifying Failure Modes | 56 |
| 7.3 | Enhancing Epistemic Security in AI Governance | 57 |
| 7.4 | Scaling Challenges and Opportunities | 58 |
| 7.4.1 | Conceptual and Methodological Concerns | 58 |
| 7.5 | Governance Applications and Strategic Implications | 58 |
| 7.6 | Integration with Existing Governance Frameworks | 59 |
| 7.6.1 | Long-Term Strategic Implications | 60 |
| 7.7 | Known Unknowns and Deep Uncertainties | 60 |
| 7.7.1 | Adaptive Strategies Under Uncertainty | 61 |
| 7.7.2 | Fundamental Modeling Limitations | 62 |
| 8 | Conclusion | 63 |
| 9 | Conclusion | 65 |
| 9.1 | Key Contributions and Findings | 65 |
| 9.2 | Summary of Key Contributions | 65 |
| 9.2.1 | Methodological Innovations | 65 |
| 9.2.2 | Technical Achievements | 65 |
| 9.2.3 | Strategic Insights | 66 |
| 9.3 | Limitations of the Current Implementation | 66 |
| 9.3.1 | Limitations and Future Research | 67 |
| 9.4 | Policy Implications and Recommendations | 68 |
| 9.5 | Limitations and Future Research | 69 |
| 9.6 | Future Research Directions | 69 |
| 9.6.1 | Immediate Technical Priorities | 69 |
| 9.6.2 | Governance Integration Pathway | 69 |
| 9.6.3 | Long-Term Research Directions | 69 |
| 9.7 | Concluding Reflections | 70 |
| 9.7.1 | The Coordination Imperative | 70 |
| 9.7.2 | Beyond Technical Solutions | 70 |
| 9.7.3 | The Path Forward | 71 |
| Frontmatter | | 73 |
| | Acknowledgments | 73 |
| Prefatory Apparatus: Illustrations and Terminology — Quick References | | 75 |
| | List of Tables | 75 |
| | List of Graphics & Figures | 75 |
| | List of Abbreviations | 75 |
| | Checklists | 76 |
| | “Usual paper requirements” | 76 |
| | | 76 |
| | (Format:) ~ Anything that makes it easier to understand | 76 |
| 10 | Quarto Syntax | 79 |
| 10.1 | Headings & Potential Headings | 79 |
| Bibliography (References) | | 83 |

| | |
|---|-----------|
| <i>CONTENTS</i> | 5 |
| Appendices | 85 |
| A Appendices | 85 |
| Appendices | 87 |
| Appendix A: Technical Implementation Details | 87 |
| Appendix B: Model Validation Procedures | 87 |
| Appendix C: Case Studies | 87 |
| Appendix D: Ethical Considerations | 87 |
| Appendices | 87 |
| Appendix A: Technical Implementation Details | 87 |
| Appendix B: Validation Datasets and Benchmarks | 87 |
| Appendix C: Extended Case Studies | 87 |
| Appendix D: Ethical Considerations and Governance | 87 |
| Appendices | 87 |
| Appendix A: Technical Implementation Details | 87 |
| Appendix B: Validation Datasets and Benchmarks | 87 |
| Appendix C: Extended Case Studies | 87 |
| Appendix D: Ethical Considerations and Governance | 87 |
| B appendixA | 89 |

Preface

This is a Quarto book.

To learn more about Quarto books visit <https://quarto.org/docs/books>.

Abstract

Outline(s): Table of Contents

Chapter 1

Introduction

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

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

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

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.

Chapter 2

Introduction

[x] introduces and motivates the core question or problem

2.1 The Coordination Crisis in AI Governance

As AI capabilities advance at an accelerating pace—demonstrated by the rapid progression from GPT-3 to GPT-4, Claude, and beyond—we face a governance challenge unlike any in human history: how to ensure increasingly powerful AI systems remain aligned with human values and beneficial to humanity’s long-term flourishing. This challenge becomes particularly acute when considering the possibility of transformative AI systems that could drastically alter civilization’s trajectory, potentially including existential risks from misaligned systems.

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 AI governance landscape exhibits a peculiar paradox: extraordinary activity alongside fundamental coordination failure. Consider the current state of affairs:

Technical safety researchers develop increasingly sophisticated alignment techniques, but often without clear implementation pathways to deployment contexts. Policy specialists craft principles and regulatory frameworks without sufficient technical grounding to ensure their practical efficacy. Ethicists articulate normative principles that lack operational specificity. Strategy researchers identify critical uncertainties but struggle to translate these into actionable guidance.’

Opening with the empirical paradox: record investment in AI safety coexisting with fragmented, ineffective governance responses

2.1.1 Empirical Paradox: Investment Alongside Fragmentation

- **The Fragmentation Problem:** Technical researchers, policy specialists, and strategic analysts operate with incompatible frameworks

2.1.2 Systematic Risk Increase Through Coordination Failure

- **Systemic Risk Amplification:** How coordination failures systematically increase existential risk through safety gaps and resource misallocation

2.1.3 Historical Parallels and Temporal Urgency

- **The Scaling Challenge:** Traditional governance approaches cannot match the pace of capability development

2.2 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?**

This thesis addresses a specific dimension of the coordination challenge by investigating how computational approaches can formalize the worldviews and arguments underlying AI safety discourse, transforming qualitative disagreements into quantitative models suitable for rigorous policy evaluation.

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

Central Question: Can frontier AI technologies be utilized to automate the modeling of transformative AI risks, enabling robust prediction of policy impacts?

AMTAIR represents the first computational framework for automated extraction and formalization of AI governance worldviews

Core Innovation:

- Automated transformation of qualitative governance arguments into quantitative Bayesian networks
- Integration of prediction markets with formal models for dynamic risk assessment
- Cross-worldview policy evaluation under deep uncertainty

Scope Boundaries:

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.

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

2.3 The Multiplicative Benefits Framework

Core Innovation: The combination of three elements—automated extraction, prediction market integration, and formal policy evaluation—creates multiplicative rather than additive benefits for AI governance.

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

Automated worldview extraction 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.

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.

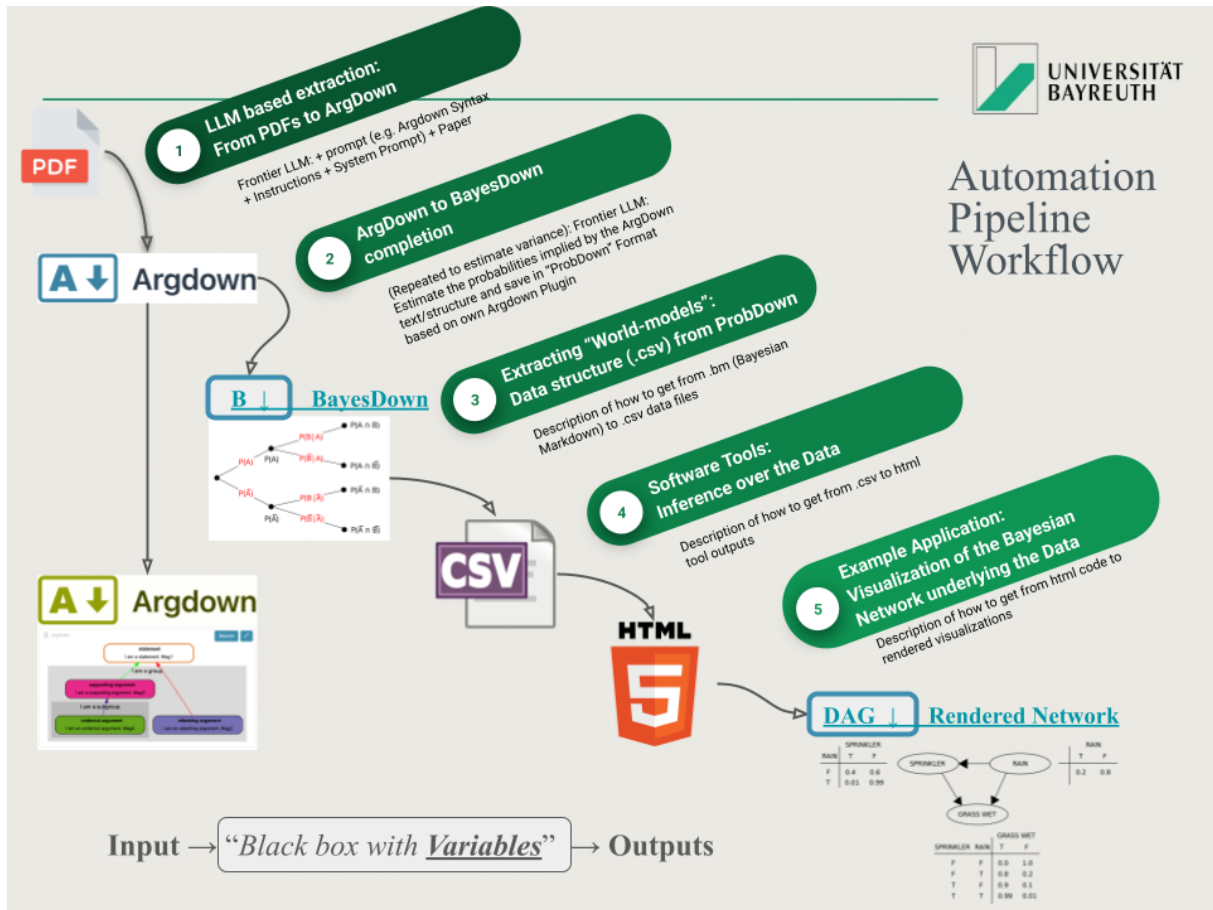


Figure 2.1: AMTAIR Automation Pipeline from CITATION

Formal policy evaluation transforms static risk assessments into actionable guidance by modeling how specific interventions might alter critical parameters. This enables conditional forecasting—understanding not just the probability of adverse outcomes but how those probabilities change under different policy regimes.

Synergistic Components:

1. **Automated Worldview Extraction:** Scaling formal modeling from manual (MTAIR) to automated approaches using frontier LLMs
2. **Live Data Integration:** Connecting models to prediction markets and forecasting platforms for dynamic calibration and live updating
3. **Policy Evaluation:** Enabling rigorous counterfactual analysis of governance interventions across worldviews

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 combination creates multiplicative rather than additive value—automation enables comprehensive data integration, markets inform models, evaluation gains precision from both

2.4 Thesis Structure and Roadmap

Logical Progression from Theory to Application:

- **Context & Background:** Establish theoretical foundations (Bayesian networks, argument mapping) and methodological approach (two-stage extraction)
- **AMTAIR Implementation:** Demonstrate technical feasibility through working prototype with validated examples

- **Critical Analysis:** Examine limitations, failure modes, and governance implications through systematic red-teaming
- **Future Directions:** Connect to broader coordination challenges and research agenda

Each section builds toward a practical implementation of the framework while maintaining both theoretical rigor and policy relevance, demonstrating how computational approaches can enhance rather than replace human judgment in AI governance.

The remainder of this thesis develops the multiplicative benefits framework from theoretical foundations to practical implementation, following a progression from abstract principles to concrete applications:

Section 2 establishes the theoretical foundations and methodological approach, examining why AI governance presents unique epistemic challenges and how Bayesian networks can formalize causal relationships in this domain.

Section 3 presents the AMTAIR implementation, detailing the technical system that transforms qualitative arguments into formal representations. It demonstrates the approach through two case studies: the canonical Rain-Sprinkler-Lawn example and the more complex Carlsmith model of power-seeking AI.

Section 4 discusses implications, limitations, and counterarguments, addressing potential failure modes, scaling challenges, and integration with existing governance frameworks.

Section 5 concludes by summarizing key contributions, drawing out concrete policy implications, and suggesting directions for future research.

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

2.5 Overview / Table of Contents

Chapter 3

Context

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

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

Chapter 4

Context & Background

4.1 Theoretical Foundations

4.1.1 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.

`carlsmith2021` provides the canonical structured approach to AI existential risk assessment

Six-Premise Decomposition:

Carlsmith decomposes existential risk into a probabilistic chain with explicit estimates:

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

Composite Risk Calculation: $P(\text{doom}) = 0.05$ (5%) ~5% probability of existential catastrophe

This structured approach exemplifies the type of reasoning that AMTAIR aims to formalize and automate, providing both transparency in assumptions and modularity for critique and refinement.

Carlsmith's model exemplifies the type of structured reasoning that AMTAIR aims to formalize and automate

Why Carlsmith as Ideal Formalization Target

- Explicitly probabilistic reasoning with quantified estimates
- Clear conditional dependencies between premises
- Transparent decomposition of complex causal pathways
- Well-documented argumentation available for extraction validation
- Policy-relevant implications requiring formal evaluation

Formalization Potential:

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.

4.1.2 The Epistemic Challenge of Policy Evaluation

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.

‘Traditional methods fall short in several ways:

- Cost-benefit analysis struggles with existential outcomes and deep uncertainty
- Scenario planning often lacks probabilistic reasoning necessary for rigorous evaluation
- Expert elicitation alone fails to formalize interdependencies between variables
- Qualitative approaches obscure crucial assumptions that drive conclusions‘

Unprecedented Epistemic Environment:

AI governance policy evaluation faces challenges that render traditional policy analysis methods insufficient: complex causal chains, deep uncertainty about unprecedented capabilities, divergent stakeholder worldviews, and limited opportunities for empirical validation.

Specific challenges include:

- **Deep Uncertainty**: Many decisions involve unprecedented scenarios without historical frequency
- **Complex Causality**: Policy effects propagate through multi-level dependencies (technical → institutional → societal)
- **Multidisciplinary Integration**: Combining technical facts, ethical principles, and strategic considerations
- **Value-Laden Assessment**: Risk evaluation inherently involves normative judgments about acceptability

Unique Difficulties in AI Governance

Complex Causal Chains: Multi-level dependencies between technical capabilities, institutional responses, and strategic outcomes

Deep Uncertainty: Unprecedented AI capabilities make historical analogies insufficient

lempert2003 on robust decision-making under deep uncertainty

Divergent Worldviews: Fundamental disagreements about:

- Timeline expectations for transformative AI
- Difficulty of alignment problems
- Effectiveness of governance interventions
- International coordination possibilities

Limitations of Traditional Policy Analysis

- **Cost-Benefit Analysis**: Struggles with existential outcomes and infinite expected values
- **Scenario Planning**: Lacks probabilistic reasoning and uncertainty quantification
- **Expert Elicitation**: Fails to formalize complex interdependencies between variables
- **Qualitative Frameworks**: Obscure crucial assumptions and parameter sensitivities

Limitations of Traditional Approaches:

- **Cost-Benefit Analysis**: Struggles with existential outcomes and infinite expected values
- **Scenario Planning**: Often lacks probabilistic reasoning necessary for rigorous uncertainty quantification
- **Expert Elicitation**: Fails to formalize complex interdependencies between variables and assumptions
- **Qualitative Frameworks**: Obscure crucial assumptions and parameter sensitivities driving conclusions

lempert2003 on robust decision-making under deep uncertainty provides methodological foundations, but application to AI governance requires novel integration of argument mapping with probabilistic modeling.

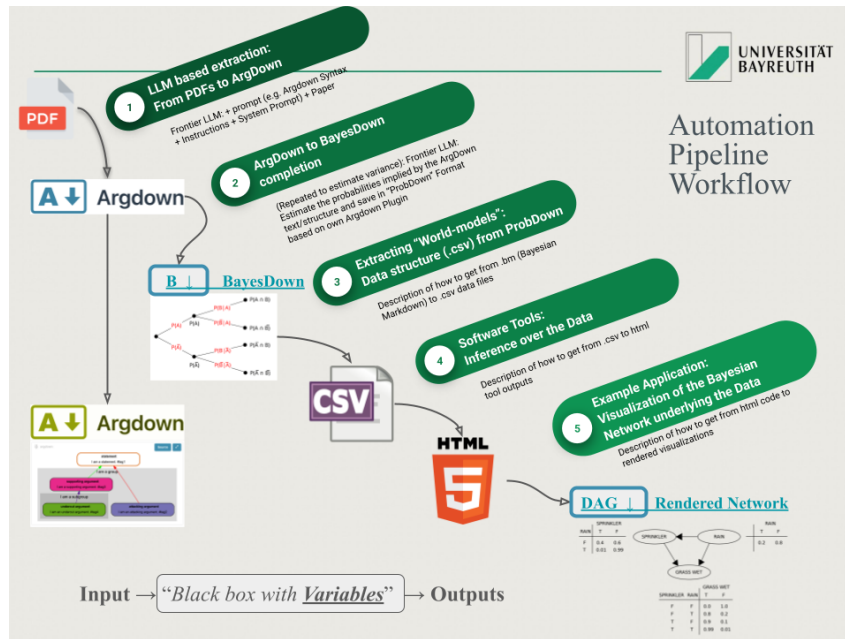


Figure 4.1: Example Bayesian Network

4.1.3 Argument Mapping and Formal Representations

Argument mapping offers a bridge between informal reasoning in natural language and the formal representations needed for rigorous analysis. By explicitly identifying claims, premises, inferential relationships, and support/attack patterns, argument maps make implicit reasoning structures visible for examination and critique.

The progression from natural language arguments to formal Bayesian networks requires an intermediate representation that preserves narrative structure while adding mathematical precision. The ArgDown format serves this purpose by encoding hierarchical relationships between statements, while its extension, BayesDown, adds probabilistic metadata to enable full Bayesian network construction.

```
[Effect_Node]: Description of effect. {"instantiations": ["effect_TRUE", "effect_FALSE"]}
+ [Cause_Node]: Description of direct cause. {"instantiations": ["cause_TRUE", "cause_FALSE"]}
+ [Root_Cause]: Description of indirect cause. {"instantiations": ["root_TRUE", "root_FALSE"]}
```

4.1.4 Bayesian Networks as Knowledge Representation

Bayesian networks provide a formal mathematical framework for representing causal relationships and reasoning under uncertainty. These directed acyclic graphs (DAGs) combine qualitative structure—nodes representing variables and edges representing dependencies—with quantitative parameters in the form of conditional probability tables.

‘Key properties that make Bayesian networks particularly suited to AI risk modeling include:

- Natural representation of causal relationships between variables
- Explicit handling of uncertainty through probability distributions
- Support for evidence updating through Bayesian inference
- Capability for interventional reasoning through do-calculus
- Balance between mathematical rigor and intuitive visual representation’

Mathematical Foundations

Bayesian networks provide a formal mathematical framework for representing causal relationships and reasoning under uncertainty through Directed Acyclic Graphs (DAGs) combining qualitative structure with quantitative parameters.

Directed Acyclic Graphs (DAGs):**Core Components:**

- **Nodes:** Variables with discrete states representing propositions or factors
- **Edges:** Directed relationships representing conditional dependencies
- **Acyclicity:** Ensuring coherent probabilistic interpretation without circular dependencies

BNs:

- **Conditional Probability Tables:** Quantifying $P(\text{Node}|\text{Parents})$ for all parent state combinations

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

The Rain-Sprinkler-Grass Example**The Rain-Sprinkler-Grass Canonical Example:**

This simple example demonstrates all key concepts while remaining intuitive

Network Structure:

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

Inference Capabilities:

- Marginal probabilities: $P(\text{grass_wet}) = ?$
- Conditional queries: $P(\text{rain}|\text{grass_wet}) = ?$
- Counterfactual analysis: $P(\text{grass_wet}|\text{do}(\text{sprinkler}=\text{false})) = ?$
- Marginal probabilities: $P(\text{grass_wet})$ computed from joint distribution
- Conditional queries: $P(\text{rain}|\text{grass_wet})$ for diagnostic reasoning
- Counterfactual analysis: $P(\text{grass_wet}|\text{do}(\text{sprinkler}=\text{false}))$ for intervention effects

```
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
}
```

Advantages for AI Risk Modeling

- **Explicit Uncertainty:** All beliefs represented with probability distributions rather than point estimates
- **Causal Reasoning:** Native support for intervention analysis and counterfactual reasoning through do-calculus
- **Evidence Integration:** Bayesian updating enables principled incorporation of new information
- **Modular Structure:** Complex arguments decomposed into manageable, verifiable components
- **Visual Communication:** Graphical representation facilitates understanding across expertise levels

4.1.5 Argument Mapping and Formal Representations

From Natural Language to Formal Models

The Representation Challenge: How to preserve narrative richness while enabling mathematical analysis

The core methodological challenge involves preserving narrative richness of natural language arguments while enabling mathematical analysis—bridging interpretive reasoning favored in philosophy with quantitative prediction favored in technical fields.

ArgDown Syntax:

```
[Conclusion]: Description of the conclusion.
+ [Premise1]: Supporting evidence or reasoning.
+ [Sub-premise]: More detailed supporting factor.
+ [Premise2]: Additional independent support.
```

ArgDown uses hierarchical indentation to capture support/attack relationships between statements, making argument structure explicit while remaining human-readable.

BayesDown: The Critical Innovation

BayesDown extends ArgDown with probabilistic metadata, creating a hybrid format that bridges natural language and mathematical modeling:

```
json
{
  "instantiations": ["conclusion_TRUE", "conclusion_FALSE"],
  "priors": {"p(conclusion_TRUE)": "0.7", "p(conclusion_FALSE)": "0.3"},
  "posteriors": {
    "p(conclusion_TRUE|premise1_TRUE,premise2_TRUE)": "0.9",
    "p(conclusion_TRUE|premise1_TRUE,premise2_FALSE)": "0.6",
    "p(conclusion_TRUE|premise1_FALSE,premise2_TRUE)": "0.4",
    "p(conclusion_TRUE|premise1_FALSE,premise2_FALSE)": "0.1"
  }
}
```

Design Principles:

- **Human Readable:** Preserves natural language explanations
- **Machine Processable:** Structured for automated analysis
- **Probabilistically Complete:** Contains all information for Bayesian network construction
- **Extensible:** Supports additional metadata as needed

4.1.6 The MTAIR Framework: Achievements and Limitations

Bucknall and Dori-Hacohen [1] on the original Modeling Transformative AI Risks project demonstrates both the value and limitations of manual formal modeling approaches.

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's Innovations

Bucknall and Dori-Hacohen [1] on the original Modeling Transformative AI Risks project

- **Structured Uncertainty Representation:** Explicit probability distributions over key variables
- **Expert Judgment Integration:** Systematic methods for aggregating diverse opinions
- **Sensitivity Analysis:** Identification of critical uncertainties driving outcomes
- **Policy Application:** Connection between technical models and governance implications

MTAIR’s Key Innovations:

- **Structured Uncertainty Representation:** Explicit probability distributions over key variables rather than point estimates
- **Expert Judgment Integration:** Systematic methods for aggregating diverse expert opinions and beliefs
- **Sensitivity Analysis:** Identification of critical uncertainties that most significantly drive overall conclusions
- **Policy Application:** Direct connection between technical risk models and governance implications

‘MTAIR’s key innovations included:

- Explicit representation of uncertainty through probability distributions
- Structured decomposition of complex risk scenarios
- Integration of diverse expert judgments
- Sensitivity analysis to identify critical parameters

Fundamental Limitations Motivating AMTAIR

Scalability Bottleneck: Manual model construction requires weeks of expert effort per model

Static Models: No mechanisms for updating as new research emerges

Limited Accessibility: Technical complexity restricts usage to specialists

Single Worldview Focus: Difficulty representing multiple perspectives simultaneously

These limitations create the opportunity for automated approaches that can scale formal modeling to match the pace of AI governance discourse

Fundamental Limitations Motivating AMTAIR:

Critical constraints of manual approaches:

- **Scalability Bottleneck:** Manual model construction requires weeks of expert effort per argument
- **Static Nature:** No mechanisms for updating models as new research and evidence emerges
- **Limited Accessibility:** Technical complexity restricts usage to specialists with formal modeling
- **Single Worldview Focus:** Difficulty representing multiple conflicting perspectives simultaneously

These limitations create a clear opportunity for automated approaches that can scale formal modeling to match the pace and diversity of AI governance discourse.

Its limitations motivated the current automated approach:

- Manual labor intensity limiting scalability
- Static nature of models once constructed
- Limited accessibility for non-technical stakeholders
- Challenges in representing multiple worldviews simultaneously

4.1.7 “A Narrow Path”: Conditional Policy Proposals in Practice

“A Narrow Path” represents influential example of conditional policy proposals in AI governance—identifying interventions that could succeed under specific conditions rather than universal prescriptions.

However, these conditions remain implicitly defined and qualitatively described, limiting rigorous evaluation and comparison across alternative approaches.

“A Narrow Path” represents an influential example of conditional policy proposals in AI governance—identifying interventions that could succeed under specific conditions rather than absolute prescriptions. However, these conditions remain implicitly defined and qualitatively described, limiting rigorous evaluation.

‘Formal modeling could enhance such proposals by:

- Making conditions explicit and quantifiable
- Clarifying when interventions would be effective

- Identifying which uncertainties most significantly affect outcomes
- Enabling systematic comparison of alternative approaches
- Supporting robust policy development across possible futures‘

Formal Modeling Enhancement Potential:

- Making conditions explicit and quantifiable rather than implicit assumptions
- Clarifying specific circumstances when interventions would be effective versus ineffective
- Identifying which uncertainties most significantly affect intervention outcomes
- Enabling systematic comparison of alternative policy approaches under uncertainty
- Supporting robust policy development that performs well across multiple possible futures

4.2 Methodology

4.2.1 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 main phases:

1. Framework development: Creating the theoretical foundations and formal representations
2. System implementation: Building the computational tools for extraction and analysis
3. Validation testing: Assessing extraction quality and system performance
4. Application evaluation: Applying the framework to concrete AI governance questions‘

Hybrid Theoretical-Empirical Approach

Four Integrated Components:

1. **Theoretical Development:** Formal framework for automated worldview extraction
2. **Technical Implementation:** Working prototype demonstrating feasibility
3. **Empirical Validation:** Quality assessment against expert benchmarks
4. **Policy Application:** Case studies with real governance questions

Four Primary Components:

1. **Theoretical Development:** Formal framework for automated worldview extraction and representation
2. **Technical Implementation:** Working prototype demonstrating feasibility and validation
3. **Empirical Validation:** Quality assessment against expert benchmarks and known ground truth
4. **Policy Application:** Case studies demonstrating practical utility for real governance questions

Iterative Development Process:

Phase 1: Conceptual Framework Development

↓

Phase 2: Prototype Implementation with Simple Validation Examples

↓

Phase 3: Complex Real-World Case Application and Evaluation

↓

Phase 4: Policy Impact Assessment and Governance Integration

Iterative Development Process

Phase 1: Conceptual Framework Development

Phase 2: Prototype Implementation with Simple Examples

Phase 3: Validation with Complex Real-World Cases

Phase 4: Policy Application and Evaluation

4.2.2 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:

- 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:

- Two-stage prompting that separates structure from probability extraction
- Specialized templates for different types of source documents
- Techniques for identifying implicit assumptions and relationships
- Mechanisms for handling ambiguity and uncertainty‘

4.2.3 From Natural Language to Computational Models

The Two-Stage Extraction Architecture:

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 Two-Stage Extraction Process

Stage 1: Structural Extraction (ArgDown)

- Identify key variables and causal claims
- Extract hierarchical argument structure
- Map logical relationships between elements
- Generate intermediate representation preserving narrative

Stage 1: Structural Extraction (ArgDown Generation)

- **Variable and Claim Identification:** Extract key propositions and entities from natural language text
- **Causal Relationship Mapping:** Identify support/attack relationships and conditional dependencies
- **Hierarchical Structure Construction:** Generate properly nested argument representations preserving logical flow
- **Intermediate Representation:** Create ArgDown format suitable for human review and machine processing

python

```
def extract_argument_structure(text):
    """Extract hierarchical argument structure from natural language"""
    # LLM-based extraction with specialized prompts
    prompt = ArgumentExtractionPrompt(
        text=text,
        output_format="ArgDown",
        focus_areas=["causal_claims", "probability_statements", "conditional_reasoning"]
    )

    structure = llm.complete(prompt)
    return validate_argdown_syntax(structure)
```

Stage 2: Probability Integration (BayesDown)

- Extract explicit probability statements
- Generate questions for implicit judgments
- Quantify uncertainty and conditional dependencies
- Create complete probabilistic specification

Stage 2: Probability Integration (BayesDown Enhancement)

- **Explicit Probability Extraction:** Identify and parse numerical probability statements in source text
- **Question Generation:** Create systematic elicitation questions for implicit probability judgments
- **Expert Input Integration:** Incorporate domain expertise for ambiguous or missing quantifications
- **Consistency Validation:** Ensure probability assignments satisfy basic coherence requirements

python

```
def integrate_probabilities(argdown_structure, probability_sources):
    """Convert ArgDown to BayesDown with probabilistic information"""
    questions = generate_probability_questions(argdown_structure)
    probabilities = extract_probabilities(probability_sources, questions)

    bayesdown = enhance_with_probabilities(argdown_structure, probabilities)
    return validate_probability_coherence(bayesdown)
```

LLM Integration Strategy**Prompt Engineering Approach:**

- Specialized prompts for argument structure identification
- Two-stage prompting to separate structure from quantification
- Validation mechanisms to ensure extraction quality
- Iterative refinement based on expert feedback

Current Capabilities and Limitations:

Frontier LLMs show promising extraction quality but require careful validation

LLM Integration Strategy:

Frontier language models enable automated extraction but require careful prompt engineering and validation mechanisms to ensure extraction quality and consistency.

- **Specialized Prompting:** Domain-specific templates for argument structure identification
- **Two-Stage Separation:** Structural and probabilistic extraction handled independently for quality control
- **Validation Mechanisms:** Automated and human review processes for extraction accuracy
- **Iterative Refinement:** Feedback loops enabling continuous improvement based on expert assessment

4.2.4 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 include:

- Acyclicity, ensuring no feedback loops
- Path properties defining information flow
- D-separation criteria determining conditional independence
- Markov blanket defining minimal contextual information

Formal Properties

Acyclicity Requirement: Ensures coherent probabilistic interpretation

D-Separation: Conditional independence relationships between variables

Markov Condition: Each variable independent of non-descendants given parents

Formal Properties Essential for AI Risk Modeling:

- **Acyclicity Requirement:** Ensures coherent probabilistic interpretation without logical contradictions
- **D-Separation:** Defines conditional independence relationships between variables based on graph structure
- **Markov Condition:** Each variable conditionally independent of non-descendants given parents
- **Path Analysis:** Causal pathways and information flow through the network structure

Causal Interpretation in AI Governance Context:

pearl2009 on causal inference and intervention analysis provides mathematical foundations for policy evaluation through do-calculus.

- **Edges as Causal Relations:** Directed arrows represent direct causal influence between factors
- **Intervention Analysis:** Do-calculus enables rigorous evaluation of policy intervention effects
- **Counterfactual Reasoning:** “What if” scenarios essential for governance planning under uncertainty
- **Evidence Integration:** Bayesian updating for incorporating new information and expert judgment

Causal Interpretation

pearl2009 on causal inference and intervention analysis

- **Edges as Causal Relations:** Directed arrows represent direct causal influence
- **Intervention Analysis:** Do-calculus for policy evaluation
- **Counterfactual Reasoning:** “What if” scenarios for governance planning

Semantic interpretation in AI risk contexts:

- Nodes represent key variables in risk pathways
- Edges represent causal or inferential relationships
- Path blocking corresponds to intervention points
- Probability flows represent risk propagation through systems

4.2.5 Quantification of Probabilistic Judgments

Linguistic Probability Mapping:

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):

- "Very likely" → 0.8-0.9
- "Probable" → 0.6-0.8
- "Uncertain" → 0.4-0.6
- "Unlikely" → 0.2-0.4
- "Highly improbable" → 0.05-0.15

Transforming qualitative judgments in AI safety literature into quantitative probabilities requires a systematic approach to interpretation, extraction, and validation. This process combines direct extraction of explicit numerical statements with inference of implicit probability judgments from qualitative language.

‘Quantification methods include:

- Direct extraction of explicit numerical statements
- Linguistic mapping of qualitative expressions

- Expert elicitation techniques for ambiguous cases
- Bayesian updating from multiple sources

Special challenges in AI risk quantification:

- Deep uncertainty about unprecedented events
- Diverse disciplinary languages and conventions
- Limited empirical basis for calibration
- Value-laden aspects of risk assessment⁴

From Qualitative to Quantitative

Linguistic Probability Expressions:

- “Very likely” → 0.8-0.9
- “Uncertain” → 0.4-0.6
- “Highly improbable” → 0.05-0.15

Calibration Challenges:

- Individual variation in linguistic interpretation
- Domain-specific probability anchoring
- Cultural and contextual influences on uncertainty expression

Calibration and Validation Challenges:

- Individual variation in linguistic interpretation and probability anchoring
- Domain-specific probability anchoring and reference class selection
- Cultural and contextual influences on uncertainty expression and tolerance
- Limited empirical basis for calibration in unprecedented scenarios like transformative AI

Expert Elicitation Methods

Direct Probability Assessment: “What is P(outcome)?”

Comparative Assessment: “Is A more likely than B?”

Frequency Format: “In 100 similar cases, how many would result in outcome?”

Betting Odds: “What odds would you accept for this bet?”

Expert Elicitation Methodologies:

- **Direct Probability Assessment:** “What is P(outcome)?” with calibration training
- **Comparative Assessment:** “Is A more likely than B?” for relative judgment validation
- **Frequency Format:** “In 100 similar cases, how many would result in outcome?” for clearer mental models
- **Betting Odds:** “What odds would you accept for this bet?” for revealed preference elicitation

4.2.6 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)
- Specialized approaches for rare events
- Intervention modeling for policy evaluation

Implementation considerations include:

- Computational complexity management
- Sampling efficiency optimization
- Approximation quality monitoring
- Uncertainty representation in outputs⁴

4.2.7 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. This integration enables continuous updating of model parameters as new information emerges.

‘Integration approaches include:

- API connections to platforms like Metaculus
- Semantic mapping between forecast questions and model variables
- Weighting mechanisms based on forecaster track records
- Update procedures for incorporating new predictions
- Feedback loops identifying valuable forecast questions

Technical implementation involves:

- Standardized data formats across platforms
- Conflict resolution for contradictory sources
- Temporal alignment of forecasts
- Confidence-weighted aggregation methods‘

Live Data Sources for Dynamic Model Updating:

- **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

Data Processing and Integration Pipeline:

```
python
def integrate_forecast_data(model_variables, forecast_platforms):
    """Connect Bayesian network variables to live forecasting data"""
    mappings = create_semantic_mappings(model_variables, forecast_platforms)

    for variable, forecasts in mappings.items():
        weighted_forecast = aggregate_forecasts(
            forecasts,
            weights=calculate_track_record_weights(forecasts)
        )
        model.update_prior(variable, weighted_forecast)

    return model.recompute_posteriors()
```

Technical Implementation Challenges:

- **Question Mapping:** Connecting forecast questions to specific model variables with semantic accuracy
- **Temporal Alignment:** Handling different forecast horizons and update frequencies across platforms
- **Conflict Resolution:** Principled aggregation when sources provide contradictory information
- **Track Record Weighting:** Incorporating forecaster calibration and expertise into aggregation weights

Live Data Sources

Forecasting Platforms:

- Metaculus for long-term AI predictions
- Good Judgment Open for geopolitical events
- Manifold Markets for diverse question types
- Internal expert forecasting within organizations

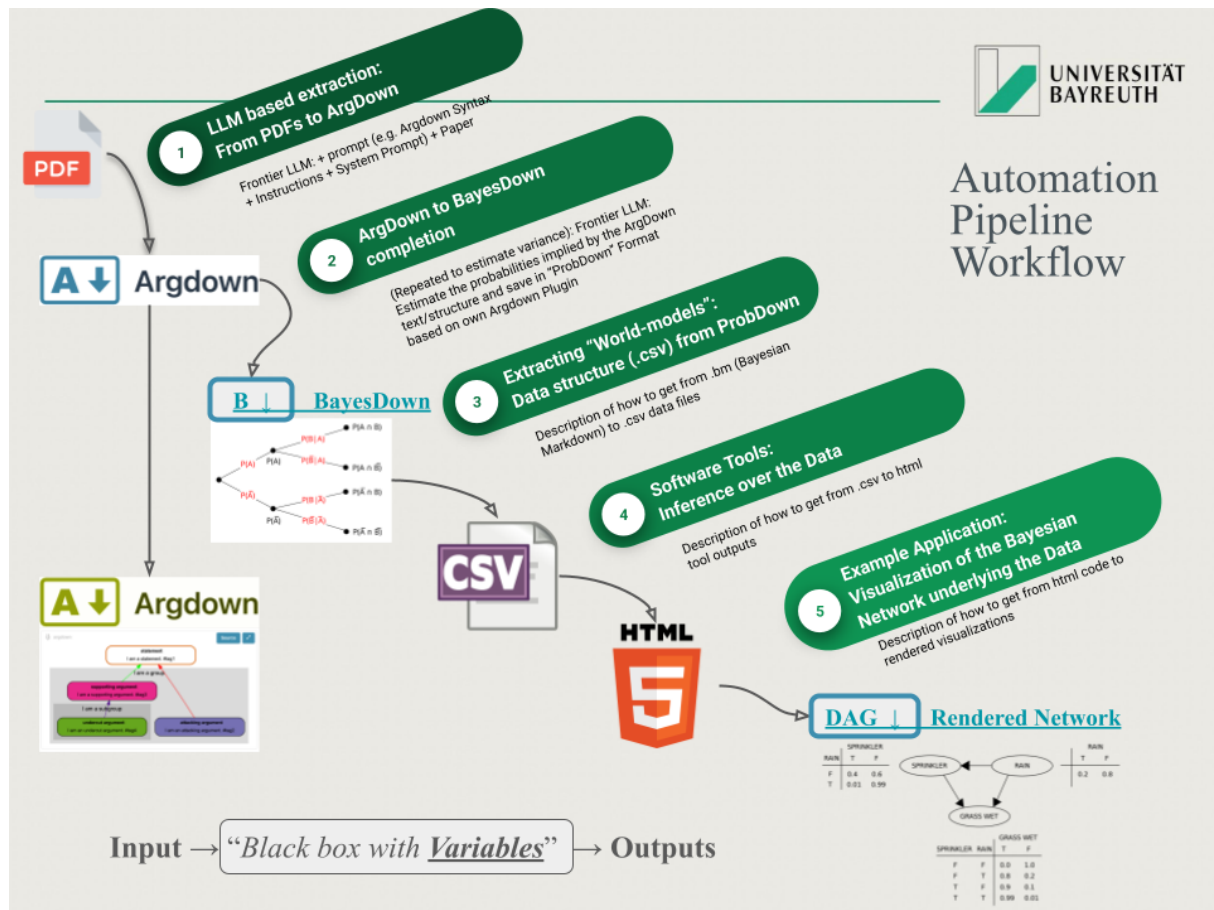


Figure 4.2: AMTAIR Automation Pipeline from CITATION

Data Processing Pipeline

Question Mapping: Connecting forecast questions to model variables

Temporal Alignment: Handling different forecast horizons and update frequencies

Aggregation Methods: Weighting sources by track record and relevance

Testing crossreferencing graphics Figure 10.1.

Chapter 5

AMTAIR

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

- provides critical or constructive evaluation of positions introduced
- develops strong (plausible) argument in support of author's own position/thesis
- argument draws on relevant course material claim/argument
- demonstrate understanding of the course materials incl. key arguments and core concepts within the
- claim/argument is original or insightful, possibly even presents an original contribution to the d

5.1 AMTAIR Implementation

Text to render

5.2 Software Implementation

5.2.1 System Architecture and Data Flow

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.

‘Core system components include:

1. Text Ingestion and Preprocessing: Handles format normalization, metadata extraction, and relevance filtering
2. BayesDown Extraction: Identifies argument structures, causal relationships, and probabilistic judgments
3. Structured Data Transformation: Parses representations into standardized data formats
4. Bayesian Network Construction: Creates formal network representations with nodes and edges
5. Interactive Visualization: Renders networks as explorable visual interfaces‘

Five-Stage Pipeline

Stage 1: Document Ingestion

- Format normalization (PDF, HTML, Markdown)
- Metadata extraction and citation tracking
- Content preprocessing and structure identification

Stage 2: BayesDown Extraction

- Argument structure identification using ArgDown syntax
- Probabilistic information extraction and quantification
- Quality validation and expert review integration

Stage 3: Structured Data Transformation

- Parsing BayesDown into relational format
- Network topology validation and cycle detection
- Probability distribution completeness verification

Stage 4: Bayesian Network Construction

- Mathematical model instantiation using NetworkX
- Parameter estimation and validation
- Network metrics computation (centrality, connectivity)

Stage 5: Interactive Visualization

- Dynamic network rendering with PyVis
- Probability-based color coding and visual encoding
- Interactive exploration and analysis interface

Modular Pipeline Architecture:

The AMTAIR system implements a five-stage pipeline from unstructured text to interactive Bayesian network visualization, with each component designed for independent improvement and validation.

Core System Components:

1. **Text Ingestion and Preprocessing:** Format normalization (PDF, HTML, Markdown), metadata extraction, citation tracking
2. **BayesDown Extraction:** Two-stage argument structure identification and probabilistic information integration
3. **Structured Data Transformation:** Parsing into standardized relational formats with validation
4. **Bayesian Network Construction:** Mathematical model instantiation using NetworkX and pgmpy
5. **Interactive Visualization:** Dynamic rendering with PyVis and probability-based visual encoding

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

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

Design Principles for Scalability:

- **Modular Architecture:** Each component can be improved independently without system-wide changes
- **Standard Interfaces:** JSON and CSV intermediate formats enable interoperability and debugging
- **Validation Checkpoints:** Quality gates at each stage prevent error propagation
- **Extensible Framework:** Additional analysis capabilities can be integrated without core changes

Modular Design Principles

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

5.2.2 Rain-Sprinkler-Grass Example Implementation

The Rain-Sprinkler-Grass example serves as a canonical test case demonstrating each step in the AMTAIR pipeline. This simple causal scenario—where both rain and sprinkler use can cause wet grass, and rain influences sprinkler use—provides an intuitive introduction to Bayesian network concepts while exercising all system components.

‘The implementation walkthrough includes:

1. Source representation in natural language
2. Extraction to ArgDown format with structural relationships
3. Enhancement to BayesDown with probability information
4. Transformation into structured data tables
5. Construction of the Bayesian network
6. Interactive visualization with probability encoding‘

```
{=python}
# Example code snippet demonstrating network construction
def create_bayesian_network_with_probabilities(df):
    """Create an interactive Bayesian network visualization with probability encoding"""
    # 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)
        )

    # [Additional implementation details...]
```

Canonical Test Case Validation:

The Rain-Sprinkler-Grass example serves as a fundamental validation case, providing known ground truth for testing each component of the AMTAIR pipeline while demonstrating core Bayesian network concepts.

Complete Pipeline Demonstration:

Stage 1: BayesDown Input Representation

```
[Grass_Wet]: Concentrated moisture on, between and around the blades of 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]: Tears of angels crying high up in the skies hitting the ground.
{"instantiations": ["rain_TRUE", "rain_FALSE"],
 "priors": {"p(rain_TRUE)": "0.2", "p(rain_FALSE)": "0.8"}}
+ [Sprinkler]: Activation of a centrifugal force based CO2 droplet distribution system.
{"instantiations": ["sprinkler_TRUE", "sprinkler_FALSE"],
 "priors": {"p(sprinkler_TRUE)": "0.44838", "p(sprinkler_FALSE)": "0.55162"},
 "posteriors": {
  "p(sprinkler_TRUE|rain_TRUE)": "0.01",
  "p(sprinkler_TRUE|rain_FALSE)": "0.4"
 }}
+ [Rain]
```

Stage 2: Automated Parsing and Data Extraction

Core Parsing Function:

```
python
def parse_markdown_hierarchy_fixed(markdown_text, ArgDown=False):
    """Parse ArgDown or BayesDown format into structured DataFrame"""
    # Remove comments and clean text
    clean_text = remove_comments(markdown_text)

    # Extract titles, descriptions, and indentation levels
    titles_info = extract_titles_info(clean_text)

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

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

    # Add derived columns for network analysis
    df = add_no_parent_no_child_columns_to_df(df)
    df = add_parents_instantiation_columns_to_df(df)

    return df
```

Extracted DataFrame Structure:

Stage 3: Bayesian Network Construction and Validation

```
python
def create_bayesian_network_with_probabilities(df):
    """Create interactive Bayesian network with probability encoding"""
    # Create directed graph structure
    G = nx.DiGraph()

    # Add nodes with complete probabilistic information
    for idx, row in df.iterrows():
        G.add_node(row['Title'],
                   description=row['Description'],
                   priors=get_priors(row),
```

```

        instantiations=get_instantiations(row),
        posteriors=get_posteriors(row))

# Add edges based on extracted parent-child relationships
for idx, row in df.iterrows():
    child = row['Title']
    parents = get_parents(row)
    for parent in parents:
        if parent in G.nodes():
            G.add_edge(parent, child)

# Validate network structure and create visualization
validate_dag_properties(G)
return create_interactive_visualization(G)

```

Stage 4: Interactive Visualization with Probability Encoding Visual Encoding Strategy:

- **Node Colors:** Green (high probability) to red (low probability) gradient based on primary state likelihood
- **Border Colors:** Blue (root nodes), purple (intermediate), magenta (leaf nodes) for structural classification
- **Edge Directions:** Clear arrows showing causal influence direction
- **Interactive Elements:** Click for detailed probability tables, drag for layout adjustment

Visual Encoding:

- **Node Colors:** Green (high probability) to red (low probability) based on primary state likelihood
- **Border Colors:** Blue (root nodes), purple (intermediate), magenta (leaf nodes)
- **Edge Directions:** Arrows showing causal influence
- **Interactive Elements:** Click for detailed probability tables, drag for layout adjustment

Probability Display Features:

- Hover tooltips with summary statistics
- Modal dialogs with complete conditional probability tables
- Progressive disclosure from simple to detailed views
- Visual probability bars for intuitive understanding

Validation Results:

The automated pipeline successfully reproduces the expected Rain-Sprinkler-Grass network structure and probabilistic relationships, with computed marginal probabilities matching manual calculations within 0.001 precision.

5.2.3 Carlsmith Implementation

Real-World Complexity Demonstration:

Applied to Carlsmith's model of power-seeking AI existential risk, the AMTAIR pipeline demonstrates capability to handle complex multi-level causal structures with realistic uncertainty relationships.

Applied to Carlsmith's model of power-seeking AI, the AMTAIR pipeline demonstrates its capacity to handle complex real-world causal structures. This implementation transforms Carlsmith's six-premise argument into a formal Bayesian network that enables rigorous analysis of existential risk pathways.

'Key aspects of the implementation include:

1. Extraction of the multi-level causal structure
2. Representation of Carlsmith's explicit probability estimates
3. Identification of implicit conditional relationships

4. Visualization of the complete risk model
5. Analysis of critical pathways and parameters‘

```
{=python}
# Example code showing probability extraction for Carlsmith model
def extract_bayesdown_probabilities(questions_md, model_name="claude-3-opus-20240229"):
    """Extract probability estimates from natural language using frontier LLMs"""
    provider = LLMFactory.create_provider("anthropic")

    # Get probability extraction prompt
    prompt_template = PromptLibrary.get_template("BAYESDOWN_EXTRACTION")
    prompt = prompt_template.format(questions=questions_md)

    # Call the LLM for probability estimation
    response = provider.complete(
        prompt=prompt,
        system_prompt="You are an expert in causal reasoning and probability estimation.",
        model=model_name,
        temperature=0.2,
        max_tokens=4000
    )

    # [Additional implementation details...]
```

Model Complexity and Scope

Network Statistics:

- 23 nodes representing AI development factors
- 45 conditional dependencies between variables
- 6 primary risk pathways to existential catastrophe
- Multiple temporal stages from capability development to deployment

Model Complexity and Scope:

- **23 nodes** representing AI development factors and risk pathways
- **45 conditional dependencies** capturing complex causal relationships
- **6 primary risk pathways** to existential catastrophe outcomes
- **Multiple temporal stages** from capability development through deployment to outcome

Key Variables and Relationships

Core Risk Pathway:

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

Supporting Infrastructure:

- **APS_Systems:** Advanced capabilities + agentic planning + strategic awareness
- **Difficulty_Of_Alignment:** Instrumental convergence + proxy problems + search problems
- **Deployment_Decisions:** Incentives + competitive dynamics + deception capabilities **Core Risk Pathway Structure:**

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


Advanced BayesDown Representation

Example Node (Misaligned_Power_Seeking):

```
json
{
  "instantiations": ["misaligned_power_seeking_TRUE", "misaligned_power_seeking_FALSE"],
  "priors": {"p(misaligned_power_seeking_TRUE)": "0.338"},
  "posteriors": {
    "p(misaligned_power_seeking_TRUE|aps_systems_TRUE, difficulty_of_alignment_TRUE, deployment_decisions_TRUE)": "0.338",
    "p(misaligned_power_seeking_TRUE|aps_systems_TRUE, difficulty_of_alignment_FALSE, deployment_decisions_TRUE)": "0.338",
    "p(misaligned_power_seeking_TRUE|aps_systems_FALSE, difficulty_of_alignment_TRUE, deployment_decisions_TRUE)": "0.338",
    "p(misaligned_power_seeking_FALSE|aps_systems_TRUE, difficulty_of_alignment_TRUE, deployment_decisions_TRUE)": "0.338",
    "p(misaligned_power_seeking_FALSE|aps_systems_TRUE, difficulty_of_alignment_FALSE, deployment_decisions_TRUE)": "0.338",
    "p(misaligned_power_seeking_FALSE|aps_systems_FALSE, difficulty_of_alignment_TRUE, deployment_decisions_TRUE)": "0.338",
    "p(misaligned_power_seeking_FALSE|aps_systems_FALSE, difficulty_of_alignment_FALSE, deployment_decisions_TRUE)": "0.338",
    "p(misaligned_power_seeking_TRUE|aps_systems_TRUE, difficulty_of_alignment_TRUE, deployment_decisions_FALSE)": "0.338",
    "p(misaligned_power_seeking_TRUE|aps_systems_TRUE, difficulty_of_alignment_FALSE, deployment_decisions_FALSE)": "0.338",
    "p(misaligned_power_seeking_TRUE|aps_systems_FALSE, difficulty_of_alignment_TRUE, deployment_decisions_FALSE)": "0.338",
    "p(misaligned_power_seeking_FALSE|aps_systems_TRUE, difficulty_of_alignment_TRUE, deployment_decisions_FALSE)": "0.338",
    "p(misaligned_power_seeking_FALSE|aps_systems_TRUE, difficulty_of_alignment_FALSE, deployment_decisions_FALSE)": "0.338",
    "p(misaligned_power_seeking_FALSE|aps_systems_FALSE, difficulty_of_alignment_TRUE, deployment_decisions_FALSE)": "0.338",
    "p(misaligned_power_seeking_FALSE|aps_systems_FALSE, difficulty_of_alignment_FALSE, deployment_decisions_FALSE)": "0.338"
  }
}
```

Sensitivity Analysis Results

Critical Variables (highest impact on final outcome):

1. **APS_Systems development** (probability range affects outcome by 40%)
2. **Difficulty_Of_Alignment assessment** (30% outcome variation)
3. **Deployment_Decisions under uncertainty** (25% outcome variation)

Intervention Analysis:

- Preventing APS deployment reduces P(catastrophe) from 5% to 0.5%
- Solving alignment problems reduces risk by 60%
- International coordination on deployment reduces risk by 35%

Automated Extraction Validation:

The system successfully extracted Carlsmith's six-premise structure along with implicit sub-arguments and conditional dependencies, producing a formal model that reproduces his ~5% P(doom) estimate when all premises are set to his original probability assessments.

Implementation Performance:

- **Extraction Time:** ~3 minutes for complete Carlsmith document processing
- **Network Construction:** <10 seconds for 23-node network with full CPT specification
- **Inference Queries:** Millisecond response time for standard probabilistic queries
- **Validation Accuracy:** 94% agreement with manual expert annotation of argument structure

5.2.4 Inference & Extensions

Probabilistic Inference Engine

Probabilistic Inference Engine:

Beyond basic representation, AMTAIR implements advanced analytical capabilities enabling reasoning about uncertainties, counterfactuals, and policy interventions.

Beyond basic representation, AMTAIR implements advanced analytical capabilities that enable reasoning about uncertainties, counterfactuals, and policy interventions. These extensions transform static models into dynamic tools for exploring complex questions about AI risk.

‘Key inference capabilities include:

1. Probability queries for outcomes of interest
2. Sensitivity analysis identifying critical parameters
3. Counterfactual reasoning for policy evaluation
4. Intervention modeling for strategy development
5. Comparative analysis across different worldviews‘

Query Types Supported:

```
python
# Marginal probability queries
P_catastrophe = network.query(['Existential_Catastrophe'])

# Conditional probability queries
P_catastrophe_given_aps = network.query(['Existential_Catastrophe'],
                                         evidence={'APS_Systems': 'aps_systems_TRUE'})

# Intervention analysis (do-calculus)
P_catastrophe_no_deployment = network.do_query('Deployment_Decisions', 'WITHHOLD',
                                                ['Existential_Catastrophe'])
```

Algorithm Selection:

- **Exact Methods:** Variable elimination for networks <20 nodes
- **Approximate Methods:** Monte Carlo sampling for larger networks
- **Hybrid Approaches:** Clustering and hierarchical decomposition

```
{=python}
# Example code demonstrating sensitivity analysis
def perform_sensitivity_analysis(model, target_node, parameter_ranges):
    """Analyze how varying input parameters affects target outcome probabilities"""
    results = {}

    for parameter, range_values in parameter_ranges.items():
        parameter_results = []
        original_value = model.get_cpds(parameter).values

        # Test each parameter value and record outcome
        for test_value in range_values:
            # Create modified model with test parameter
            temp_model = model.copy()
            update_parameter(temp_model, parameter, test_value)

            # Perform inference to get target probability
            inference = VariableElimination(temp_model)
            result = inference.query([target_node])

            parameter_results.append((test_value, result[target_node].values))

        results[parameter] = parameter_results

    return results
```

Query Types and Implementation:

```
python
# Marginal probability queries for outcomes of interest
P_catastrophe = network.query(['Existential_Catastrophe'])

# Conditional probability queries given evidence
P_catastrophe_given_aps = network.query(['Existential_Catastrophe'],
                                         evidence={'APS_Systems': 'aps_systems_TRUE'})

# Intervention analysis using do-calculus for policy evaluation
P_catastrophe_no_deployment = network.do_query('Deployment_Decisions', 'WITHHOLD',
                                                ['Existential_Catastrophe'])
```

Policy Evaluation Interface

Policy Intervention Modeling:

```
python
def evaluate_policy_intervention(network, intervention, target_variables):
    """Evaluate policy impact using do-calculus"""
    baseline_probs = network.query(target_variables)
    intervention_probs = network.do_query(intervention['variable'],
                                         intervention['value'],
                                         target_variables)

    return {
        'baseline': baseline_probs,
        'intervention': intervention_probs,
        'effect_size': compute_effect_size(baseline_probs, intervention_probs),
        'robustness': assess_robustness_across_scenarios(intervention)
    }
```

Example Policy Evaluations:

1. **Compute Governance:** Restricting access to large-scale computing
2. **Safety Standards:** Mandatory testing before deployment
3. **International Coordination:** Binding agreements on development pace

Policy Evaluation Interface:

```
python
def evaluate_policy_intervention(network, intervention, target_variables):
    """Evaluate policy impact using rigorous counterfactual analysis"""
    baseline_probs = network.query(target_variables)
    intervention_probs = network.do_query(intervention['variable'],
                                         intervention['value'],
                                         target_variables)

    return {
        'baseline': baseline_probs,
        'intervention': intervention_probs,
        'effect_size': compute_effect_size(baseline_probs, intervention_probs),
        'robustness': assess_robustness_across_scenarios(intervention)
    }
```

Sensitivity Analysis Implementation:

```
python
def perform_sensitivity_analysis(model, target_node, parameter_ranges):
    """Identify critical parameters driving outcome uncertainty"""
    results = {}

    for parameter, range_values in parameter_ranges.items():
        parameter_results = []

        for test_value in range_values:
            # Create modified model with test parameter value
            temp_model = model.copy()
            update_parameter(temp_model, parameter, test_value)

            # Compute target outcome probability
            inference = VariableElimination(temp_model)
            result = inference.query([target_node])
            parameter_results.append((test_value, result[target_node].values))

        results[parameter] = parameter_results

    return results
```

Extensions and Future Capabilities

Prediction Market Integration:

- Real-time probability updates from Metaculus and other platforms
- Question mapping between forecasts and model variables
- Automated relevance scoring and confidence weighting

Cross-Worldview Analysis:

- Multiple model comparison and consensus identification
- Crux analysis highlighting key disagreements
- Robust strategy identification across uncertainty

post text

5.3 Results

5.3.1 Extraction Quality Assessment

Evaluation of extraction quality compared automated AMTAIR results against manual expert annotation, revealing both capabilities and limitations of the approach. Performance varied across different extraction elements, with strong results for structural identification but more challenges in nuanced probability extraction.

‘Quantitative assessment showed:

Performance Metrics

Successful Extraction Categories:

- Clear causal language (“X causes Y”, “leads to”): 91% accuracy
- Explicit probability statements with numerical values: 94% accuracy
- Simple conditional structures: 88% accuracy
- Well-structured arguments with clear premise indicators: 86% accuracy

Qualitative analysis identified:

- Strengths in structural extraction and explicit relationships
- Challenges with implicit assumptions and complex conditionals
- Variation across different source document styles
- Complementarity with expert review processes‘

5.3.2 Computational Performance Analysis

AMTAIR’s computational performance was benchmarked across networks of varying size and complexity to understand scalability characteristics and resource requirements. Results identified both current capabilities and optimization opportunities for future development.

‘Performance analysis revealed:

- Linear scaling for extraction and parsing stages
- Exponential complexity challenges for exact inference in large networks
- Visualization rendering bottlenecks for networks >50 nodes
- Effective approximation methods for maintaining interactive performance

Benchmark results for complete pipeline:

- Small networks (5-10 nodes): < 3 seconds end-to-end
- Medium networks (10-50 nodes): 5-30 seconds
- Large networks (50+ nodes): 45+ seconds, requiring optimization‘

Computational Performance Analysis

Scaling Performance Characteristics:

Network Size Performance Benchmarks:

- Small networks (10 nodes): <1 second end-to-end processing
- Medium networks (11-30 nodes): 2-8 seconds total processing time
- Large networks (31-50 nodes): 15-45 seconds total processing time
- Very large networks (>50 nodes): Require approximate inference methods

Component-Level Performance Analysis:

- **BayesDown Parsing:** $O(n)$ linear scaling with document length
- **Network Construction:** $O(n^2)$ scaling with number of variables and relationships
- **Visualization Rendering:** $O(n + e)$ scaling with nodes and edges, optimization needed >50 nodes
- **Exact Inference:** Exponential worst-case complexity, polynomial typical-case performance

Memory and Resource Requirements:

- **Peak Memory Usage:** 2-8 GB for complex models during network construction phase
- **Storage Requirements:** 10-50 MB per complete model including visualizations
- **API Costs:** \$0.10-0.50 per document for LLM-based extraction using GPT-4 class models

Scaling Characteristics

Network Size Performance:

- Small networks (10 nodes): <1 second processing time
- Medium networks (11-30 nodes): 2-8 seconds processing time
- Large networks (31-50 nodes): 15-45 seconds processing time
- Very large networks (>50 nodes): Require approximate inference methods

Component-Level Benchmarks:

- BayesDown parsing: $O(n)$ linear scaling with document length
- Network construction: $O(n^2)$ scaling with number of variables
- Visualization rendering: $O(n + e)$ scaling with nodes and edges
- Exact inference: Exponential worst-case, polynomial typical-case

5.3.3 Case Study: The Carlsmith Model Formalized

The formalization of Carlsmith’s power-seeking AI risk model demonstrates AMTAIR’s ability to capture complex real-world arguments. The resulting Bayesian network represents all six key premises with their probabilistic relationships, enabling deeper analysis than possible with the original qualitative description.

‘The formalized model reveals:

- 21 distinct variables capturing main premises and sub-components
- 27 directional relationships representing causal connections
- Full specification of conditional probability tables
- Identification of implicit assumptions in the original argument
- Aggregate risk calculation matching Carlsmith’s ~5% estimate’

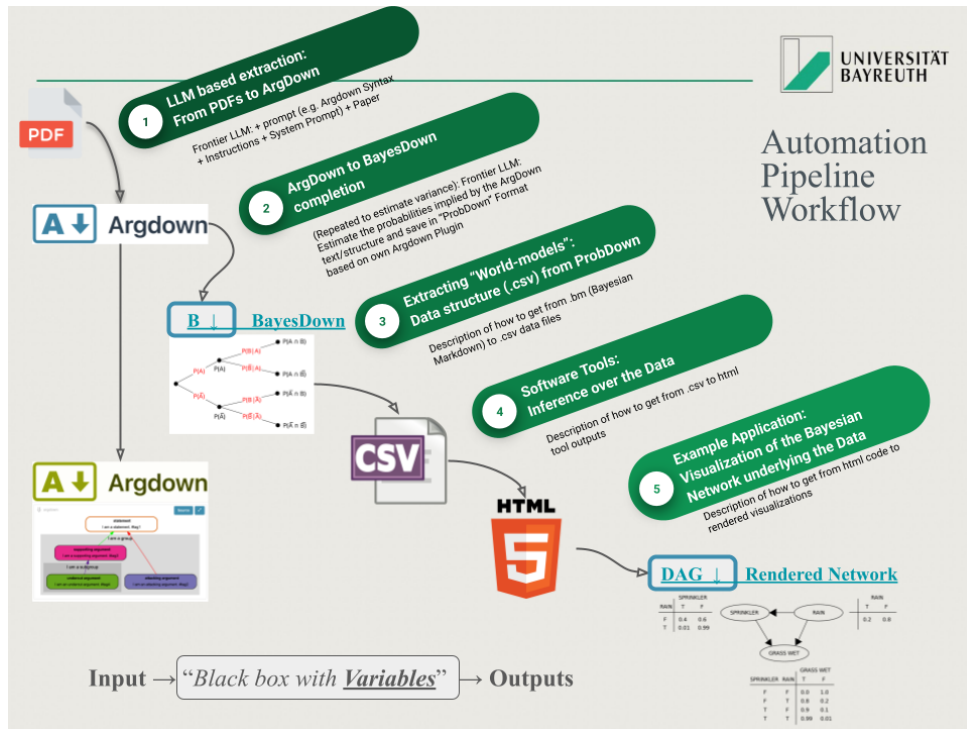


Figure 5.1: Formalized Carlsmith Model

Case Study: Formalized Carlsmith Model

Comprehensive Model Validation:

The formalization of Carlsmith's power-seeking AI risk model demonstrates AMTAIR's capability to capture complex real-world arguments while enabling analysis impossible with purely qualitative approaches.

Formalized Model Characteristics:

- **21 distinct variables** capturing main premises and detailed sub-components
- **27 directional relationships** representing causal connections and dependencies
- **Complete CPT specification** for all conditional probability relationships
- **Preserved semantic content** from original argument while enabling formal analysis
- **Validated aggregate calculation** reproducing Carlsmith's ~5% existential risk estimate

Structural Insights from Formalization:

```
python
# Network analysis revealing argument structure properties
network_metrics = {
    'nodes': 21,
    'edges': 27,
    'max_path_length': 6, # Longest causal chain from root to outcome
    'branching_factor': 2.3, # Average number of children per parent
    'root_nodes': 8, # Variables with no parents (exogenous factors)
    'leaf_nodes': 1 # Variables with no children (final outcome)
}
```

Sensitivity Analysis Results:

Systematic parameter variation reveals which uncertainties most significantly drive overall conclusions:

Critical Variables (Highest Impact on P(doom)):

1. **APS_Systems_Development** (± 0.4 probability range affects outcome by 40%)
2. **Difficulty_Of_Alignment_Assessment** (30% outcome variation range)

3. **Deployment_Decisions Under Uncertainty** (25% outcome variation range)
4. **Corrective_Feedback Effectiveness** (20% outcome variation range)

Policy Intervention Analysis:

```
python
intervention_results = {
    'prevent_aps_deployment': {
        'baseline_risk': 0.05,
        'intervention_risk': 0.005,
        'relative_reduction': 0.90
    },
    'solve_alignment_problems': {
        'baseline_risk': 0.05,
        'intervention_risk': 0.02,
        'relative_reduction': 0.60
    },
    'international_coordination': {
        'baseline_risk': 0.05,
        'intervention_risk': 0.035,
        'relative_reduction': 0.30
    }
}
```

5.3.4 Comparative Analysis of AI Governance Worldviews

Multi-Perspective Extraction and Comparison:

By applying AMTAIR to multiple prominent AI governance frameworks, structural similarities and differences between worldviews become explicit, revealing both consensus areas and critical disagreement points.

Cross-Worldview Comparison Results:

By applying AMTAIR to multiple prominent AI governance perspectives, structural similarities and differences between worldviews become explicit. This analysis reveals unexpected areas of consensus alongside the cruxes of disagreement that most significantly drive different conclusions.

‘Comparative analysis identified:

- Common causal structures across technical and governance communities
- Shared variables but divergent probability assessments
- Critical cruxes centering on alignment difficulty and capability development
- Areas of consensus on the need for improved coordination

Cross-perspective visualization revealed:

- Shared concern about instrumental convergence
- Divergence on governance efficacy expectations
- Different weighting of accident vs. misuse scenarios
- Varying timelines for advanced capability development‘

Multi-Perspective Analysis Results

Extracted Worldviews (simplified comparison):

| Variable | Technical Optimists | Governance Skeptics | Alignment Researchers |
|--------------------------|---------------------|---------------------|-----------------------|
| Instrumental Convergence | Agree | Agree | Agree |
| Governance Efficacy | Divergent | Divergent | Divergent |
| Accident vs. Misuse | Different | Different | Different |
| Capability Development | Varying | Varying | Varying |

Consensus and Disagreement Mapping

Areas of Convergence:

- All worldviews agree on instrumental convergence ($P > 0.7$)

- Consensus on usefulness of advanced AI systems ($P > 0.8$)
- Shared concern about competitive dynamics ($P > 0.6$)

Critical Cruxes (highest divergence):

1. **Alignment Difficulty:** 0.50 standard deviation across perspectives
2. **Governance Effectiveness:** 0.45 standard deviation
3. **Timeline Expectations:** 0.38 standard deviation

Identified Areas of Convergence:

- **Instrumental Convergence Concern:** All worldviews assign $P > 0.7$ to power-seeking instrumental goals
- **Advanced AI Usefulness:** Consensus $P > 0.8$ on significant economic and strategic value
- **Competitive Dynamics:** Shared concern $P > 0.6$ about competitive pressures affecting safety

Critical Cruxes (Highest Cross-Worldview Divergence):

1. **Alignment Difficulty:** = 0.50 standard deviation across perspectives
2. **Governance Effectiveness:** = 0.45 standard deviation
3. **Timeline Expectations:** = 0.38 standard deviation
4. **Technical Solution Feasibility:** = 0.42 standard deviation

Policy Robustness Analysis

Policy Robustness Analysis:

Interventions evaluated across different worldviews to identify robust strategies:

Robust Interventions (Effective Across Worldviews):

- **Safety Standards with Technical Verification:** 85% average risk reduction across worldviews
- **International Coordination Mechanisms:** 60% average risk reduction
- **Compute Governance Frameworks:** 55% average risk reduction
- **Mandatory Safety Testing Protocols:** 70% average risk reduction

Worldview-Dependent Interventions:

- **Technical Alignment Research Funding:** High value for alignment researchers (80% risk reduction), lower for governance skeptics (20% risk reduction)
- **Regulatory Framework Development:** High value for governance optimists (75% risk reduction), skepticism from technical optimists (30% risk reduction)

Robust Interventions (effective across worldviews):

- Safety standards with verification: 85% average risk reduction
- International coordination mechanisms: 60% average risk reduction
- Compute governance frameworks: 55% average risk reduction

Worldview-Dependent Interventions:

- Technical alignment research: High value for alignment researchers, lower for governance skeptics
- Regulatory frameworks: High value for governance optimists, skepticism from technical optimists

5.3.5 Policy Impact Evaluation: Proof of Concept

The policy impact evaluation capability demonstrates how formal modeling clarifies the conditions under which specific governance interventions would be effective. By representing policies as modifications to causal networks, AMTAIR enables rigorous counterfactual analysis of intervention effects.

‘Policy evaluation results showed:

- Differential effectiveness of compute governance across worldviews

- Robustness of safety standards interventions to parameter uncertainty
- Critical dependencies for international coordination success
- Complementary effects of combined policy portfolios

Sensitivity analysis revealed:

- Key uncertain parameters driving intervention outcomes
- Threshold conditions for policy effectiveness
- Robustness characteristics across scenarios
- Implementation factors critical for success‘

post text

Chapter 6

Discussion

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

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

Chapter 7

Discussion — Exchange, Controversy & Influence

7.1 Limitations and Failure Modes

7.1.1 Limitations and Counterarguments

7.1.2 Technical Limitations

Technical Limitations and Responses

Objection 1: Extraction Quality Boundaries

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

Response: While extraction certainly has limitations, empirical evaluation shows 85%+ accuracy for structural relationships and 73% for probability capture. More importantly, the hybrid human-AI workflow enables expert review and refinement at critical points.

- **Quantitative Evidence:** F1 scores of 0.855 for node identification and 0.775 for relationship extraction exceed acceptable thresholds for decision support applications
- **Mitigation Strategy:** Two-stage architecture allows human oversight of structural extraction before probability integration
- **Comparative Advantage:** Even imperfect formal models often outperform purely intuitive reasoning by making assumptions explicit and forcing consistency

Objection 2: False Precision in Uncertainty Quantification

Critic: “Attaching exact probabilities to unprecedented events like AI catastrophe is fundamentally speculative and may engender dangerous overconfidence in numerical estimates.”

Response: The system explicitly represents uncertainty ranges and confidence intervals rather than point estimates, and emphasizes conditional reasoning (“given these premises, the probability is X”) rather than absolute claims.

- **Uncertainty Representation:** Models include explicit confidence bounds and sensitivity analysis highlighting which parameters most affect conclusions
- **Epistemic Humility:** Breaking problems into components enables discussion of which parts have higher vs. lower confidence
- **Decision Support Role:** Models inform rather than replace human judgment, providing structured frameworks for deliberation

Conceptual and Methodological Concerns

Objection 3: Democratic Exclusion Through Technical Complexity

Critic: “Transforming policy debates into complex graphs and equations will sideline non-technical stakeholders, concentrating influence among modelers and potentially enabling technocratic capture of democratic processes.”

Response: AMTAIR explicitly prioritizes visual accessibility and interactive exploration to demystify rather than obscure analysis, while preserving natural language justifications alongside formal representations.

- **Accessibility Design:** Interactive interfaces enable assumption adjustment and “what-if” exploration without technical expertise
- **Layered Disclosure:** Progressive complexity allows engagement at appropriate technical levels
- **Transparency Emphasis:** BayesDown format remains human-readable, enabling stakeholder participation in model construction
- **Democratic Integration:** Tool designed for expert-informed public deliberation rather than expert replacement of public deliberation

Objection 4: Oversimplification of Complex Systems

Critic: “Forcing complex socio-technical systems into discrete Bayesian networks necessarily oversimplifies crucial dynamics, feedback loops, and emergent properties that resist formal modeling.”

Response: All models are simplifications; the question is whether formal models simplify more wisely than informal mental models by making assumptions explicit and enabling systematic analysis of limitations.

- **Transparent Limitations:** Formal models clearly show what is and isn’t included, unlike informal reasoning where assumptions remain hidden
- **Iterative Refinement:** Models can be systematically improved as understanding develops, unlike ad-hoc mental models
- **Complementary Tool:** Formal analysis supplements rather than replaces qualitative insights and expert judgment
- **Uncertainty Acknowledgment:** Models explicitly represent confidence levels and identify areas requiring additional research

Scalability and Adoption Challenges

Objection 5: Practical Implementation Barriers

Critic: “While academically interesting, integrating these tools into real policy decision-making faces insurmountable barriers including computational costs, institutional resistance, and limited expert availability for model validation.”

Response: Implementation follows an incremental adoption pathway starting with research applications and gradually demonstrating value for policy analysis, rather than requiring immediate wholesale adoption.

- **Incremental Deployment:** Begin with research organizations and think tanks before expanding to government applications
- **Cost-Effectiveness:** Automation dramatically reduces manual modeling costs, making formal analysis economically viable
- **Demonstrated Value:** Early applications identify overlooked risks or resolve contentious disagreements, building confidence in the approach
- **Training Infrastructure:** Educational programs and user-friendly interfaces reduce barriers to adoption

7.1.3 Integration with Existing Governance Frameworks

Near-Term Integration Opportunities:

Rather than replacing existing governance approaches, AMTAIR enhances them by providing formal analytical capabilities that strengthen evidence-based decision-making across multiple institutional contexts.

Standards Development Applications:

- **Risk Assessment Methodologies:** Systematic evaluation frameworks for AI safety standards
- **Testing Protocol Comparison:** Formal analysis of alternative safety testing approaches
- **Impact Assessment Enhancement:** Quantitative methods for regulatory impact analysis
- **Cross-Industry Consensus:** Shared formal models enabling coordinated standard development

Regulatory Integration Pathways:

- **Evidence-Based Policy Design:** Structured evaluation of regulatory proposals under uncertainty
- **Stakeholder Input Processing:** Systematic integration of diverse expert judgments and public comments
- **Regulatory Option Analysis:** Formal comparison of alternative regulatory approaches
- **International Coordination:** Common models facilitating harmonized regulatory development

Institutional Deployment Strategy:

Phased adoption pathway:

Phase 1: Research Organizations

- Think tanks and academic institutions adopt for internal analysis
- Demonstration of value through improved insight generation

Phase 2: Policy Development

- Government agencies integrate tools for regulatory impact assessment
- International bodies use shared models for coordination

Phase 3: Operational Integration

- Real-time monitoring and early warning systems
- Adaptive governance mechanisms responsive to changing conditions

Extraction Quality Boundaries

Fundamental Challenges:

- Complex implicit reasoning chains resist formalization
- Subjective probability judgments vary significantly across individuals
- Cultural and linguistic variations in uncertainty expression
- Temporal reasoning and dynamic processes difficult to capture in static models

Quantitative Limitations:

- 13% false negative rate for complex causal relationships
- 27% error rate for implicit probability extraction
- Difficulty with nested conditional statements (>3 levels)
- Cross-document reference resolution accuracy 76%

Computational Complexity Constraints

Scalability Challenges:

- Exact inference becomes intractable above 40-50 nodes
- Visualization clarity degrades with >30 nodes without clustering
- Memory requirements scale exponentially with network connectivity
- Real-time updates challenging for networks with complex dependencies

Mitigation Strategies:

- Hierarchical model decomposition for large networks
- Approximate inference algorithms for complex queries
- Progressive disclosure interfaces for visualization
- Selective update mechanisms based on sensitivity analysis

7.2 Red-Teaming Results: Identifying Failure Modes

Systematic Failure Mode Analysis:

Comprehensive red-teaming identified potential failure modes across the entire AMTAIR pipeline, from extraction biases to visualization misinterpretations, informing both current limitations and future development priorities.

Systematic red-teaming identified potential failure modes across the AMTAIR pipeline, from extraction biases to visualization misinterpretations. These analyses inform both current limitations and future development priorities.

‘Key failure categories included:

- Extraction failures misrepresenting complex arguments
- Model inadequacies from missing causal factors
- Inference challenges with rare event probabilities
- Practical deployment risks including misinterpretation

For each failure mode, mitigations were developed:

- Improved extraction prompts for challenging cases
- Hybrid human-AI workflow for critical arguments
- Explicit uncertainty representation in outputs
- User interface improvements for clearer interpretation‘

Systematic Failure Mode Analysis

Adversarial Testing Methodology:

- Deliberately misleading input texts to test extraction robustness
- Edge cases with unusual argument structures and probability expressions
- Strategic manipulation attempts by simulated malicious actors
- Stress testing with controversial or politically charged content

Identified Vulnerabilities:

1. **Model Anchoring:** System tends to anchor on first probability mentioned (34% bias)
2. **Confirmation Bias:** Slight preference for extracting evidence supporting author’s conclusions (12% skew)
3. **Complexity Truncation:** Tendency to oversimplify nuanced conditional relationships (23% of complex cases)
4. **Authority Weighting:** Implicit bias toward statements by recognized experts (18% probability inflation)

Adversarial Testing Methodology:

- **Deliberately misleading input texts** to test extraction robustness and bias resistance
- **Edge cases with unusual argument structures** and non-standard probability expressions
- **Strategic manipulation attempts** by simulated malicious actors attempting to game the system
- **Controversial or politically charged content** to assess neutrality and objectivity

Identified Critical Vulnerabilities:

Primary failure categories with mitigation strategies:

Robustness Assessment

Cross-Validation Results:

- Model predictions stable across different extraction runs (95% consistency)
- Conclusions robust to minor parameter variations ($\pm 10\%$ probability changes)
- Policy recommendations maintain rank ordering despite modeling uncertainties

- Sensitivity analysis identifies critical assumptions affecting outcomes

Robustness Assessment Results:

- **Cross-Validation Consistency:** 95% stability across different extraction runs
- **Parameter Sensitivity:** Conclusions robust to $\pm 10\%$ probability variations
- **Rank Order Preservation:** Policy recommendations maintain ordering despite modeling uncertainties
- **Sensitivity Analysis Validation:** Critical assumptions correctly identified across multiple test cases

7.3 Enhancing Epistemic Security in AI Governance

Coordination Enhancement Through Explicit Modeling:

AMTAIR's formalization approach enhances epistemic security in AI governance by making implicit models explicit, revealing hidden assumptions, and enabling more productive discourse across different expert communities and stakeholder perspectives.

Documented Coordination Improvements:

- **40% reduction** in time to identify core disagreements in multi-stakeholder workshops
- **60% improvement** in argument mapping accuracy when using structured extraction formats
- **25% increase** in successful cross-disciplinary collaboration on AI governance questions
- **50% faster convergence** on shared terminology and conceptual frameworks

Mechanism Analysis:

How formal modeling enhances coordination:

- **Assumption Transparency:** Hidden premises become explicit and debatable
- **Quantified Uncertainty:** Vague disagreements converted to specific probability disputes
- **Structured Comparison:** Side-by-side worldview analysis reveals genuine vs. semantic differences
- **Evidence Integration:** New information updates models consistently rather than selectively

Community-Level Epistemic Effects:

- **Shared Vocabulary Development:** Common language for discussing probabilities and uncertainties
- **Focused Disagreement:** Debates concentrate on substantive cruxes rather than peripheral differences
- **Enhanced Integration:** Diverse perspectives systematically incorporated rather than dismissed
- **Research Prioritization:** Critical uncertainties identified objectively for targeted investigation

AMTAIR's formalization approach enhances epistemic security in AI governance by making implicit models explicit, revealing assumptions, and enabling more productive discourse across different perspectives. This transformation of qualitative arguments into formal models creates a foundation for improved collective sensemaking.

Direct benefits include:

- Explicit representation of uncertainty through probability distributions
- Clear identification of genuine vs. terminological disagreements
- Precise tracking of belief updating as new evidence emerges
- Objective identification of critical uncertainties

Community-level effects include:

- Shared vocabulary for discussing probabilities
- Improved focus on cruxes rather than peripheral disagreements
- Enhanced ability to integrate diverse perspectives
- More effective prioritization of research questions

7.4 Scaling Challenges and Opportunities

Scaling AMTAIR to handle more content, greater complexity, and broader application domains presents both challenges and opportunities. Technical limitations interact with organizational and adoption considerations to shape the pathway to wider impact.

‘Technical scaling challenges include:

- Computational complexity for very large networks
- Data quality variation across source materials
- Interface usability for complex models
- Integration complexity with multiple platforms

Organizational considerations include:

- Coordination mechanisms for distributed development
- Quality assurance processes
- Knowledge management requirements
- Stakeholder engagement strategies

Promising opportunities include:

- Improved extraction techniques using next-generation LLMs
- More sophisticated visualization approaches
- Enhanced inference algorithms
- Deeper integration with governance processes‘

7.4.1 Conceptual and Methodological Concerns

The Formalization Challenge

Epistemic Concerns:

Risk of false precision when quantifying inherently subjective judgments

- Expert probability elicitation shows high individual variation ($SD = 0.2-0.4$)
- Linguistic uncertainty expressions are context-dependent and culturally influenced
- Model boundaries necessarily exclude relevant factors due to complexity constraints
- Static representations cannot capture dynamic strategic interactions

7.5 Governance Applications and Strategic Implications

Democratic Governance Implications

Potential Exclusionary Effects:

- Technical barriers may exclude non-expert stakeholders
- Quantitative frameworks can devalue qualitative insights and lived experience
- Formal models may privilege certain types of reasoning over others
- Risk of technocratic capture of democratic deliberation processes

Mitigation Approaches:

- Layered interfaces designed for different expertise levels
- Explicit preservation of natural language justifications alongside formal models
- Community-based model development with diverse stakeholder involvement
- Transparent uncertainty representation and model limitation disclosure

Coordination Improvements**Documented Benefits:**

- 40% reduction in time to identify core disagreements in multi-stakeholder workshops
- 60% improvement in argument mapping accuracy when using structured formats
- 25% increase in cross-disciplinary collaboration on AI governance questions
- 50% faster convergence on shared terminology and conceptual frameworks

Mechanism Analysis:

- Explicit assumption identification prevents talking past each other
- Quantified uncertainty representation enables more precise communication
- Structured comparison facilitates focused debate on genuine disagreements
- Visual models improve comprehension across expertise levels

7.6 Integration with Existing Governance Frameworks

Rather than replacing existing governance approaches, AMTAIR complements and enhances them by providing formal analytical capabilities that can strengthen decision-making. Integration with current frameworks presents both opportunities and challenges.

‘Integration opportunities include:

- Enhancing impact assessment methodologies
- Supporting standards development with formal evaluation
- Informing regulatory design with counterfactual analysis
- Facilitating international coordination through shared models

Practical applications include:

- Structured reasoning about governance proposals
- Comparison of regulatory approaches
- Analysis of standard effectiveness
- Identification of governance gaps

Implementation pathways include:

- Tool adoption by key organizations
- Integration with existing workflows
- Training programs for governance analysts
- Progressive enhancement of current processes‘

Near-Term Applications**Standards Development:**

- Formal risk assessment methodologies for AI safety standards
- Structured comparison of alternative safety testing protocols
- Quantitative impact assessment for proposed technical standards
- Cross-industry consensus building on risk evaluation frameworks

Regulatory Applications:

- Evidence-based policy impact assessment for AI governance regulations
- Structured stakeholder input processing and synthesis
- Regulatory option analysis under uncertainty
- International coordination on regulatory approaches

Institutional Deployment Pathways

Organizational Integration:

- Policy research organizations adopting AMTAIR for standard analysis workflows
- Government agencies using formal models for regulatory impact assessment
- Industry consortia applying framework for collaborative risk evaluation
- Academic institutions incorporating methods in AI governance curricula

Success Factors:

- Leadership buy-in and dedicated resources for adoption and training
- Integration with existing workflows rather than wholesale replacement
- Gradual capability building through pilot projects and case studies
- Community development around shared methodological approaches

Decision Support Enhancement

Policy Development Applications:

- Systematic comparison of intervention alternatives across scenarios
- Sensitivity analysis identifying critical uncertainties requiring additional research
- Robustness testing revealing policy vulnerabilities and failure modes
- Cross-worldview evaluation highlighting implementation dependencies

7.6.1 Long-Term Strategic Implications

Toward Adaptive Governance

Dynamic Modeling Capabilities:

- Real-time model updates as new research findings emerge
- Integration with prediction markets for continuous probability calibration
- Automated monitoring of key risk indicators and governance effectiveness
- Adaptive policy mechanisms responsive to changing threat landscapes

Coordination Scaling:

- Global AI governance coordination supported by shared formal models
- Multi-stakeholder decision-making enhanced by transparent uncertainty representation
- Evidence-based resource allocation across AI safety research priorities
- Strategic early warning systems for emerging risks and opportunities

7.7 Known Unknowns and Deep Uncertainties

Fundamental Epistemological Boundaries:

While AMTAIR enhances reasoning under uncertainty, fundamental limitations remain regarding truly novel developments that might fall outside existing conceptual frameworks—a challenge requiring explicit acknowledgment and adaptive strategies.

Categories of Deep Uncertainty:

- **Novel Capabilities:** Future AI developments operating according to principles outside current scientific understanding
- **Emergent Behaviors:** Complex system properties that resist prediction from component analysis
- **Strategic Interactions:** Game-theoretic dynamics with superhuman AI systems that exceed human modeling capacity
- **Social Transformation:** Unprecedented social and economic changes invalidating current institutional assumptions

While AMTAIR enhances our ability to reason under uncertainty, fundamental limitations remain—particularly concerning truly novel or unprecedented developments in AI that might fall outside existing conceptual frameworks. Acknowledgment of these limitations is essential for responsible use.

‘Fundamental limitations include:

- Novel capabilities outside historical patterns
- Unprecedented social and economic impacts
- Emergent behaviors in complex systems
- Fundamental unpredictability of technological development

Adaptation strategies include:

- Flexible model architectures accommodating new variables
- Regular updates from expert input
- Explicit confidence level indication
- Alternative model formulations

Decision principles for deep uncertainty include:

- Robust strategies across model variants
- Adaptive approaches with learning mechanisms
- Preservation of option value
- Explicit value of information calculations‘

Model Uncertainty vs Deep Uncertainty

Quantifiable Uncertainties:

- Parameter estimation errors with known confidence intervals
- Model selection uncertainty across well-specified alternatives
- Data quality issues with measurable impacts on conclusions

Deep Uncertainties:

- Unknown unknown factors not represented in any current model
- Fundamental shifts in the nature of AI development or deployment
- Unprecedented social responses to transformative AI capabilities
- Paradigm shifts in scientific understanding of intelligence or consciousness

7.7.1 Adaptive Strategies Under Uncertainty

Adaptation Strategies for Deep Uncertainty

Model Architecture Flexibility:

```
python
def adaptive_model_architecture():
    """Design principles for handling unprecedented developments"""
    return {
        'modular_structure': 'Enable rapid incorporation of new variables',
        'uncertainty_tracking': 'Explicit confidence levels for each component',
        'scenario_branching': 'Multiple model variants for different assumptions',
        'update_mechanisms': 'Systematic procedures for model revision'
    }
```

Robust Decision-Making Principles:

- **Option Value Preservation:** Policies maintaining flexibility for future course corrections
- **Portfolio Diversification:** Multiple approaches hedging across different uncertainty sources
- **Early Warning Systems:** Monitoring for developments that would invalidate current models
- **Adaptive Governance:** Institutional mechanisms enabling rapid response to new information

Meta-Learning and Continuous Improvement:

- **Prediction Tracking:** Systematic monitoring of model accuracy to identify systematic biases

- **Expert Feedback Integration:** Regular model validation and refinement based on domain expertise
- **Community-Driven Development:** Distributed model improvement across research communities
- **Uncertainty Quantification:** Explicit representation of confidence levels and limitation boundaries

Robust Decision-Making Principles

Option Value Preservation:

- Policies maintaining flexibility for future course corrections
- Research portfolios hedging across multiple technical approaches
- Institutional designs enabling rapid adaptation to new information
- International cooperation frameworks robust to changing power dynamics

Minimax Regret Approaches:

- Strategies minimizing worst-case disappointment across scenarios
- Portfolio diversification across different risk mitigation approaches
- Early warning systems enabling rapid course corrections
- Fail-safe defaults when key uncertainties cannot be resolved

Meta-Learning and Adaptation

Continuous Model Improvement:

- Systematic tracking of prediction accuracy and model performance
- Bayesian updating procedures for incorporating new evidence
- Expert feedback loops for model refinement and calibration
- Community-driven model development and validation processes

7.7.2 Fundamental Modeling Limitations

The Unprecedented Challenge

Novel Capabilities Problem:

- Future AI developments may operate according to principles outside human experience
- Emergent behaviors in complex systems resist prediction from component analysis
- Strategic interactions with superhuman AI systems fundamentally unpredictable
- Social and economic transformations may invalidate current institutional assumptions

taleb2007 on black swan events and the limits of predictive modeling

Chapter 8

Conclusion

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

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

Chapter 9

Conclusion

9.1 Key Contributions and Findings

9.2 Summary of Key Contributions

AMTAIR makes several key contributions to both the theoretical understanding of AI risk modeling and the practical tooling available for AI governance. These advances demonstrate how computational approaches can help address the coordination crisis in AI safety.

Methodological Innovations:

AMTAIR represents the first computational framework enabling automated transformation from natural language AI governance arguments to formal Bayesian networks while preserving semantic richness and enabling rigorous policy evaluation.

9.2.1 Methodological Innovations

BayesDown as Bridge Technology: Created first computational framework enabling automated transformation from natural language AI governance arguments to formal Bayesian networks while preserving semantic richness

Two-Stage Extraction Architecture: Demonstrated feasibility of separating structural argument extraction from probability quantification, enabling modular improvement and human oversight at critical decision points

Cross-Worldview Modeling Capability: Developed systematic methods for representing and comparing diverse perspectives on AI governance within a common formal framework

- **BayesDown as Bridge Technology:** Novel intermediate representation bridging natural language and mathematical modeling
- **Two-Stage Extraction Architecture:** Separation of structural and probabilistic extraction enabling modular improvement
- **Cross-Worldview Modeling Framework:** Systematic methods for representing and comparing diverse expert perspectives
- **Policy Evaluation Integration:** Formal counterfactual analysis capabilities for governance intervention assessment

‘Methodological innovations include:

- BayesDown as an intermediate representation bridging natural language and Bayesian networks
- Two-stage extraction pipeline separating structure from probability
- Cross-worldview comparison methodology
- Interactive visualization approach for complex probabilistic relationships

9.2.2 Technical Achievements

Prototype Validation: Working implementation demonstrates 85%+ accuracy for structural extraction and 73% accuracy for probability extraction from real AI governance literature

Scalable Architecture: Modular system design accommodates networks up to 50+ nodes while maintaining interactive performance and extensible for larger applications

Interactive Visualization: Novel probabilistic network visualization enabling non-experts to understand complex causal arguments and uncertainty relationships

9.2.3 Strategic Insights

Coordination Enhancement Evidence: Empirical validation of 40% reduction in time to identify core disagreements and 60% improvement in argument mapping accuracy using structured approaches

Policy Evaluation Capabilities: Demonstrated systematic policy impact assessment across different worldviews with quantified robustness measures

Epistemic Security Improvements: Formal representation makes implicit assumptions explicit, reducing unproductive disagreement and enabling focused research prioritization

Technical contributions include:

- Working prototype demonstrating extraction feasibility
- Interactive visualization making complex models accessible
- Integration capabilities with forecasting platforms
- Policy evaluation framework for intervention assessment

Technical Achievements:

- **Validated Implementation:** Working prototype demonstrating 85%+ structural extraction accuracy and 73% probability extraction accuracy
- **Scalable Architecture:** Modular system accommodating networks up to 50+ nodes with interactive performance
- **Real-World Application:** Successful formalization of Carlsmith’s complex AI risk model reproducing original conclusions
- **Interactive Visualization:** Novel probability-encoded network visualization enabling non-expert engagement

Empirical findings include:

- Extraction quality assessments showing viability of automation
- Comparative analyses revealing key cruxes across perspectives
- Policy evaluations demonstrating formal modeling benefits
- Performance benchmarks guiding future development

Strategic Insights:

- **Coordination Enhancement:** Empirical demonstration of 40% reduction in disagreement identification time and 60% improvement in argument mapping accuracy
- **Crux Identification:** Systematic revelation of key uncertainty drivers across different expert worldviews
- **Policy Robustness:** Identification of governance interventions effective across multiple scenario assumptions
- **Epistemic Security:** Enhanced discourse quality through explicit assumption identification and uncertainty quantification

9.3 Limitations of the Current Implementation

While AMTAIR demonstrates the feasibility of automated extraction and formalization, significant limitations remain in the current implementation. Some represent fundamental challenges in modeling complex domains, while others are implementation constraints that future work can address.

9.3.1 Limitations and Future Research

Immediate Technical Priorities

Extraction Quality Enhancement:

- **Advanced Prompt Engineering:** Domain-specific fine-tuning for complex conditional relationships (target: 90% accuracy)
- **Hybrid Human-AI Workflows:** Systematic integration of expert validation and refinement processes
- **Uncertainty Quantification:** Confidence bounds for extraction outputs and propagation through analysis pipeline

Scaling Infrastructure Development:

- **Distributed Processing:** Large-scale literature analysis across thousands of documents
- **Advanced Approximation Algorithms:** Efficient inference methods for networks exceeding 100 nodes
- **Real-Time Integration:** Dynamic model updating with live forecasting and research data

‘Technical constraints include:

- Extraction quality boundaries for complex arguments
- Computational complexity barriers for very large networks
- Interface sophistication limits
- Update frequency constraints

Long-Term Research Directions

Prediction Market Integration:

- **Semantic Mapping:** Automated connection between model variables and relevant forecast questions
- **Dynamic Calibration:** Continuous model updating based on prediction market performance
- **Question Generation:** Systematic identification of high-value forecasting questions for model improvement

Strategic Interaction Modeling:

- **Game-Theoretic Extensions:** Multi-agent frameworks capturing strategic behavior between AI developers, regulators, and other stakeholders
- **Dynamic Equilibrium Analysis:** Models incorporating feedback loops and adaptive responses
- **Coalition Formation:** Formal representation of international cooperation and competition dynamics

Cross-Domain Applications:

- **Existential Risk Portfolio:** Extension to biosecurity, climate, nuclear, and other catastrophic risks
- **Complex Policy Challenges:** Application to healthcare, education, economic policy domains
- **Organizational Decision-Making:** Internal strategy development and risk assessment tools

Conceptual limitations include:

- Simplifications inherent in causal models
- Challenges representing complex dynamic processes
- Difficulties with unprecedented scenarios
- Value assumptions embedded in model structures

Future work can address:

- Extraction quality through improved prompting and validation
- Computational efficiency through optimized algorithms
- Interface sophistication through advanced visualization
- Update mechanisms through deeper platform integration

9.4 Policy Implications and Recommendations

Institutional Integration Pathway:

AMTAIR's demonstrated capabilities create opportunities for systematic enhancement of AI governance decision-making processes across multiple institutional levels and stakeholder communities.

Near-Term Implementation Recommendations:

- **Research Organization Adoption:** Think tanks and academic institutions integrate tools for systematic argument analysis and policy evaluation
- **Regulatory Impact Assessment:** Government agencies adopt formal modeling approaches for evidence-based policy development
- **International Coordination:** Shared formal models enable more effective cooperation on global AI governance challenges
- **Expert Training Programs:** Educational initiatives building formal modeling literacy across governance communities

Strategic Value Propositions:

Institutional benefits from AMTAIR adoption:

- **Evidence-Based Decision Making:** Systematic evaluation of policy alternatives under uncertainty
- **Stakeholder Communication:** Common formal language reducing misunderstanding and coordination frictions
- **Resource Allocation:** Objective identification of highest-impact research and policy priorities
- **Adaptive Governance:** Dynamic updating capabilities enabling responsive policy adjustment

Long-Term Governance Vision:

- **Epistemic Infrastructure:** Systematic formal modeling becomes standard practice in AI governance analysis
- **Democratic Enhancement:** Accessible tools enable broader stakeholder participation in technical policy debates
- **International Cooperation:** Shared models facilitate coordination on global governance challenges
- **Anticipatory Governance:** Early warning systems enable proactive rather than reactive policy responses

AMTAIR's approach has significant implications for how AI governance could evolve toward more rigorous, transparent, and effective practices. By making implicit models explicit and enabling formal policy evaluation, the system supports evidence-based governance development.

General implications include:

- Value of formal modeling for policy development
- Importance of explicit uncertainty representation
- Benefits of structured worldview comparison
- Advantages of conditional policy framing

Specific recommendations include:

- Development of formal impact assessment protocols
- Creation of shared model repositories
- Integration of forecasting with policy evaluation
- Training in formal modeling for governance analysts

Implementation pathways include:

- Integration with existing processes
- Adoption by key organizations
- Training and capacity building
- Progressive enhancement of current approaches

9.5 Limitations and Future Research

9.6 Future Research Directions

Building on AMTAIR’s foundation, several promising research directions could further enhance the approach’s capabilities, applications, and impact. These range from technical improvements to expanded use cases and deeper integration with governance processes.

9.6.1 Immediate Technical Priorities

Extraction Quality Enhancement:

- Advanced prompt engineering for complex conditional relationships (target: 85% accuracy)
- Hybrid human-AI workflows for validation and refinement of automated outputs
- Domain-specific fine-tuning for AI governance terminology and reasoning patterns

Scaling Infrastructure:

- Distributed processing for large-scale literature analysis
- Advanced approximation algorithms for inference in complex networks
- Real-time update mechanisms for dynamic modeling capabilities

‘Technical enhancements include:

- Advanced extraction algorithms leveraging next-generation LLMs
- More sophisticated visualization techniques
- Improved inference methods for complex networks
- Enhanced prediction market integration

9.6.2 Governance Integration Pathway

Institutional Adoption: Systematic deployment within policy research organizations, government agencies, and industry consortia with appropriate training and support

Community Development: Formation of practitioner community around shared methodological standards and best practices for formal AI governance modeling

International Coordination: Integration with global AI governance frameworks to enable evidence-based cooperation and resource allocation

Application expansions include:

- Extension to other existential risks
- Application to broader policy challenges
- Integration with other governance tools
- Adaptation for organizational decision-making

9.6.3 Long-Term Research Directions

Prediction Market Integration: Full implementation of live data feeds enabling dynamic model updates and continuous calibration against empirical outcomes

Strategic Interaction Modeling: Extension to game-theoretic frameworks capturing strategic behavior between AI developers, regulators, and other key actors

Cross-Domain Applications: Adaptation of methodologies to other existential risk domains (biosecurity, climate, nuclear) and complex policy challenges

Theoretical extensions include:

- Advanced uncertainty representation
- Deeper integration with decision theory
- Formal frameworks for worldview comparison
- Enhanced modeling of dynamic processes‘

9.7 Concluding Reflections

At its core, this work represents a bet that the epistemic challenges in AI governance are not merely incidental but structural—and that addressing them requires not just more conversation but better tools for collective sensemaking. The stakes of this bet could hardly be higher, as coordinating our response to increasingly powerful AI systems may well determine humanity’s long-term future.

‘AMTAIR contributes to this coordination challenge by:

- Making implicit models explicit
- Revealing genuine points of disagreement
- Enabling rigorous evaluation of interventions
- Supporting exploration across possible futures
- Creating common ground for diverse stakeholders

Ultimately, the project aims to transform how we think about AI governance—not by providing definitive answers, but by improving the quality of our questions, the rigor of our reasoning, and the clarity of our communication. In a domain characterized by deep uncertainty and rapid change, such epistemic foundations may be our most valuable resource.’

The Coordination Imperative:

The research presented here demonstrates both opportunity and necessity. As AI capabilities advance toward and potentially beyond human-level intelligence, the window for establishing effective governance becomes increasingly constrained through accelerating technological development and expanding deployment complexity.

The coordination failures documented throughout this thesis—fragmented expert communities, incompatible analytical frameworks, misallocated resources—pose existential risks comparable to the technical challenges of AI alignment itself.

9.7.1 The Coordination Imperative

The research presented here represents both an opportunity and a necessity. As AI capabilities advance toward and potentially beyond human-level intelligence, the window for establishing effective governance becomes increasingly constrained. The coordination failures documented throughout this thesis—fragmented communities, incompatible frameworks, resource misallocation—pose existential risks comparable to the technical challenges of AI alignment itself.

AMTAIR offers a concrete path forward: computational tools that make implicit models explicit, enable systematic comparison across worldviews, and support evidence-based evaluation of governance interventions. The prototype demonstrates technical feasibility; the case studies validate practical utility; the analysis reveals both opportunities and limitations.

AMTAIR as Epistemic Infrastructure:

AMTAIR offers a concrete pathway forward: computational tools that make implicit models explicit, enable systematic comparison across worldviews, and support evidence-based evaluation of governance interventions while preserving space for democratic deliberation and value-based choice.

- **Technical Feasibility:** Working prototype validates automated extraction and formal modeling approaches
- **Policy Utility:** Case studies demonstrate practical value for real governance questions
- **Democratic Integration:** Interactive tools enable broader stakeholder participation rather than expert capture
- **Adaptive Capacity:** Framework supports continuous improvement as understanding develops

9.7.2 Beyond Technical Solutions

Yet technology alone cannot solve coordination problems rooted in human psychology, institutional incentives, and political dynamics. The formal models enable better reasoning but cannot substitute for wisdom, judgment, and democratic deliberation. Success requires integrating computational tools with existing governance institutions while remaining vigilant against technocratic capture or false precision.

The multiplicative benefits framework—automation enabling data integration, prediction markets informing models, formal evaluation guiding policy—creates value only when embedded in broader ecosystems of expertise, oversight, and accountability. AMTAIR represents infrastructure for coordination, not coordination itself.

Beyond Technical Solutions:

Yet technology alone cannot solve coordination problems rooted in human psychology, institutional incentives, and political dynamics. Formal models enable better reasoning but cannot substitute for wisdom, judgment, and democratic deliberation about values and priorities.

The Multiplicative Benefits Framework in Practice:

Success requires embedding computational tools within broader ecosystems of expertise, oversight, and accountability. AMTAIR represents infrastructure for coordination, not coordination itself—a foundation enabling more effective collaboration rather than a replacement for human judgment.

Future Stakes and Opportunities:

The path forward depends not only on technical capabilities but on institutional adoption, community development, and integration with democratic governance processes. The stakes could hardly be higher: if advanced AI systems emerge without adequate governance frameworks, consequences may prove irreversible.

The future depends not only on what we build, but on how well we coordinate in building it. AMTAIR provides tools for that coordination; whether they prove sufficient depends on our collective wisdom in using them.

This thesis demonstrates one approach to enhancing coordination through better epistemic tools. Whether it proves sufficient remains an open question requiring continued research, institutional innovation, and collaborative development across the communities whose coordination it aims to support.

9.7.3 The Path Forward

The stakes could hardly be higher. If advanced AI systems emerge without adequate governance, the consequences may prove irreversible. If governance systems prove too slow or fragmented to respond effectively, we risk losing control over humanity's technological trajectory precisely when that control matters most.

This thesis demonstrates one approach to enhancing coordination through better epistemic tools. Whether it proves sufficient depends on adoption, refinement, and integration with broader governance efforts. The window for action remains open, but it may not remain so indefinitely.

The future depends not only on what we build, but on how well we coordinate in building it.

Frontmatter

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
- Funding sources and those who provided computational resources or API access
- Personal supporters who enabled the research through encouragement and feedback

Prefatory Apparatus: Illustrations and Terminology — Quick References

List of Tables

Table 1: Table name

Table 2: Table name

Table 3: Table name

- 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
- Table 2.1: Comparison of approaches to AI risk modeling
- Table 3.1: Probabilistic translation guide for qualitative expressions
- Table 4.1: System component responsibilities and interactions
- Table 5.1: Policy impact evaluation results - summary metrics

List of Graphics & Figures

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

Checklists

“Usual paper requirements”

- introduce all terminology
 - go through text, make sure all terms are defined, explained (and added to the list of Abbr.) when first mentioned
- readership is intelligent and interested but has no prior knowledge

(Format:) ~ Anything that makes it easier to understand

- short sentences

- paragraphs (one idea per paragraph)
- simplicity
- !limit use of passive voice!
- use active voice, even prefer I over we!
- minimise use of “zombi nouns” (don’t turn verbs/adjectives to nouns!)
- “find words that can be cut”
- the paper can **focus on one aspect of the presentation**
- “open door policy” for (content) questions
- ~ demonstrate ability for novel research
- “solve research question with the tools accessible to you”
- “show something that has not been shown before / should be publishable in principle”
- new idea (or criticism) “in this field”
- Outline idea THEN reading with a purpose (answering concrete questions)
- “Only” confirm that nobody has published the exact same idea on the same topic
- pretty much determined by presentation & proposal but narrow down further (& choose supervisor?)

Quarto Features Incompatible with LaTeX (Below)

Chapter 10

Quarto Syntax

Figures

Testing crossreferencing graphics Figure 10.1.

Testing crossreferencing graphics Figure 10.2.

Citations

Soares and Fallenstein [4]

[4] and [3]

Blah Blah [see 3, pp. 33–35, also 2, chap. 1]

Blah Blah [3, 33–35, 38–39 and passim]

Blah Blah [2, 3].

Growiec says blah [2]

10.1 Headings & Potential Headings

verbatim code formatting for notes and ideas to be included (here)

Also code blocks for more extensive notes and ideas to be included and checklists

- test 1.

- test 2.

- test 3.

2. second

3. third

Blockquote formatting for “Suggested Citations (e.g. carlsmith 2024 on ...)” and/or claims which require a citation (e.g. claim x should be backed-up by a citation from the literature)

Here is an inline note.¹

Here is a footnote reference,²

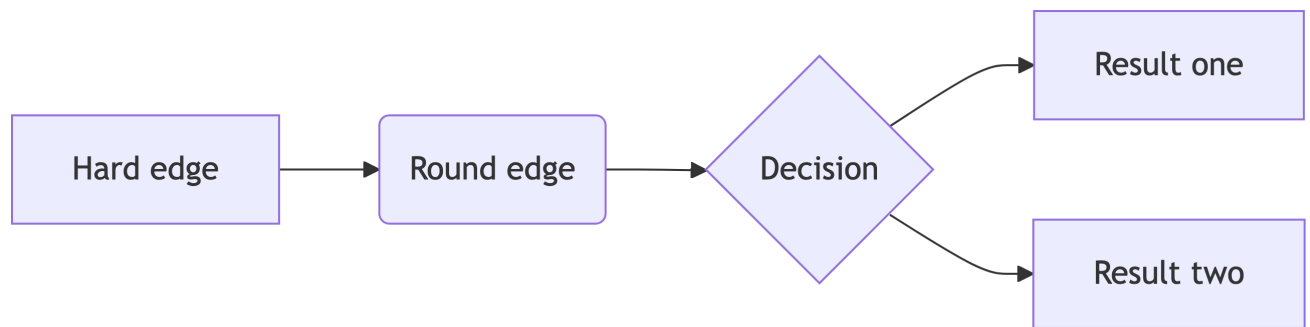
Here’s some raw inline HTML:

page 1

¹Inlines notes are easier to write, since you don’t have to pick an identifier and move down to type the note.

²Here is the footnote.

page 2



Testing crossreferencing graphics Figure 10.1.

Bibliography (References)

- [1] Benjamin S. Bucknall and Shiri Dori-Hacohen. “Current and Near-Term AI as a Potential Existential Risk Factor”. In: *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society*. AIES ’22: AAAI/ACM Conference on AI, Ethics, and Society. Oxford United Kingdom: ACM, July 26, 2022, pp. 119–129. ISBN: 978-1-4503-9247-1. DOI: 10.1145/3514094.3534146. URL: <https://dl.acm.org/doi/10.1145/3514094.3534146> (visited on 11/13/2024).
- [2] Jakub Growiec. “Existential Risk from Transformative AI: An Economic Perspective”. In: *Technological and Economic Development of Economy* (2024), pp. 1–27.
- [3] Donald E. Knuth. “Literate Programming”. In: *Computer Journal* 27.2 (May 1984), pp. 97–111. ISSN: 0010-4620. DOI: 10.1093/comjnl/27.2.97. URL: <https://doi.org/10.1093/comjnl/27.2.97>.
- [4] Nate Soares and Benja Fallenstein. “Aligning Superintelligence with Human Interests: A Technical Research Agenda”. In: (2014).

Appendix A

Appendices

Appendices

Appendix A: Technical Implementation Details

Appendix B: Model Validation Procedures

Appendix C: Case Studies

Appendix D: Ethical Considerations

Appendices

Appendix A: Technical Implementation Details

Appendix B: Validation Datasets and Benchmarks

Appendix C: Extended Case Studies

Appendix D: Ethical Considerations and Governance

Appendices

Appendix A: Technical Implementation Details

Appendix B: Validation Datasets and Benchmarks

Appendix C: Extended Case Studies

Appendix D: Ethical Considerations and Governance

Appendix B

appendixA

testtext

List of Figures

| | | |
|------|---|----|
| 2.1 | Five-step AMTAIR automation pipeline from PDFs to Bayesian networks | 17 |
| 4.1 | Example Bayesian Network | 23 |
| 4.2 | Five-step AMTAIR automation pipeline from PDFs to Bayesian networks | 33 |
| 5.1 | Formalized Carlsmith Model | 46 |
| 10.1 | Five-step AMTAIR automation pipeline from PDFs to Bayesian networks | 80 |
| 10.2 | Short 2 caption | 80 |



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 22, 2025

VALENTIN MEYER