



UNIVERSITÄT
BAYREUTH

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

Automating the Modelling of Transformative Artificial Intelligence Risks

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

A thesis submitted at the Department of Philosophy

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

Author:

Valentin Jakob Meyer
Valentin.meyer@uni-bayreuth.de
Matriculation Number: 1828610
Tel.: +49 (1573) 4512494
Pielmühler Straße 15
52066 Lappersdorf

Supervisor:

Dr. Timo Speith

Word Count:

30.000

Source / Identifier:

Document URL

26th of May 2025

Table of Contents

Preface	1
Abstract	3
Prefatory Apparatus: Frontmatter	5
Illustrations and Terminology — Quick References	5
Acknowledgments	5
List of Graphics & Figures	5
List of Abbreviations	5
Quarto Syntax and Best Practices Guide	9
Key Features	9
1. Task Management System	9
2. Multi-Format Output	9
3. Cross-Referencing	9
4. Advanced Features	9
Quick Start	10
Task Management	10
Adding Content	10
Best Practices	10
1. Consistent Formatting	10
2. Task Tracking	10
3. Version Control	10
4. Documentation	11
Troubleshooting	11
Common Issues	11
Getting Help	11
License	11
Document Structure and Headings	11
Heading Hierarchy	11
Text Formatting	12
Basic Formatting	12
Advanced Formatting	12

Links	12
Including Code	12
Diagrams	12
In-Line LaTeX	13
In-Line HTML	13
Reference or Embed Code from .ipynb files	13
Code chunks from .ipynb notebooks can be embedded in the .qmd text with:	13
which produces the output of executing the code cell:	13
including ‘echo=true’ renders the code of the cell:	16
Embed .html result/rendering from .ipynb Notebook	18
Html Graph by Notebook Cell Inclusion - (from github-pages)	18
Html Graph by Notebook Cell Inclusion with Website Call?	19
Full Bayesian Network Rendering	19
Rain-Sprinkler-Grass Network Rendering	19
Lists and Enumerations	19
Unordered Lists	19
Ordered Lists	19
Definition Lists	20
Code Blocks and Verbatim Text	20
Inline Code	20
Code Blocks with Syntax Highlighting	20
Verbatim Text	20
Blockquotes and Callouts	21
Simple Blockquote	21
Callout Blocks	21
Figures and Images	22
Complete Figure Syntax	22
Figure Best Practices	22
Tables	22
Markdown Tables	22
Grid Tables	22
Citations and References	23
Citation Styles	23
Bibliography Entry	24
Cross-References	24
Section References	24
Figure and Table References	24
Equation References	24
Mathematics	25
Inline Math	25
Display Math	25
Footnotes	25

Simple Footnote	25
Referenced Footnote	26
Appendices	26
Structure	26
Best Practices for Appendices	26
Glossary and Abbreviations	27
Glossary Format	27
Interactive Elements	27
Jupyter Notebook Embedding	27
Mermaid Diagrams	27
Line Breaks and Spacing	28
Spacing Rules	28
Page Breaks	28
Comments and Metadata	28
HTML Comments	28
Comprehensive Task Management System for Quarto Thesis	29
Overview	29
Task Categories and Syntax	29
1. General Tasks	29
2. Citation Tasks	30
3. Figure/Graphic Tasks	30
4. Content Tasks	30
5. Technical Tasks	30
Task States	31
Open or In-ProgressTasks	31
Completed Tasks (Visible in ToDo-Tree)	31
Verified/Archived Tasks (Hidden from ToDo-Tree)	31
Advanced Task Formatting	31
Multi-line Tasks with Details	31
Linked Tasks	32
Conditional Tasks	32
Task Tracking Best Practices	32
1. Task Creation Guidelines	32
2. Task Organization	32
3. Priority System	32
Task Management Workflow	34
1. Task Creation	34
2. Task Execution	34
3. Task Completion	34
4. Task Archival	34
Best Practices Summary — ALWAYS CONSISTENTLY:	34

Tagging and Highlighting System for Content Merging	37
Overview	37
Tagging Categories	37
A. Duplicate Content Marking	37
B. Redundant Content Highlighting	37
C. Better Version Available	37
D. Merge Candidate	38
Visual Marking System	38
Color-Coded Highlighting	38
Border Marking	38
Inline Marking	38
Metadata Tracking	39
Comprehensive Metadata	39
Quick Reference Tags	39
Workflow for Content Merging	39
1. Initial Marking Phase	39
2. Review and Comparison	39
3. Consolidation Actions	40
Automated Detection Helpers	40
Search Patterns	40
Duplicate Detection Checklist	40
Best Practices for Merging	41
1. Pre-Merge Preparation	41
2. During Merge Process	41
3. Post-Merge Cleanup	41
Templates for Common Scenarios	41
Duplicate Definition	41
Redundant Example	41
Overlapping Sections	42
Visual Summary Blocks	42
Merge Status Dashboard	42
Decision Log	42
Quality Assurance	43
Pre-Publication Checklist	43
Final Verification	43
 Master Thesis Checklist for Quarto Projects	 45
Content Creation Checklist	45
Document Structure	45
Text Quality	45
Formatting Elements	45
Figures and Tables	45
Citations and References	46

Cross-References	46
Revision Phase	46
Content Review	46
Task Completion	46
Prime Directives	47
Quarto Syntax Cheat-Sheet Best-Practice	47
Tagging / Highlighting System	48
Workflow Rules	48
Rigorous Checklist	48
Automating the Modeling of Transformative Artificial Intelligence Risks (AM-TAIR)	49
Frontmatter: Preface	51
Acknowledgments	51
List of Figures	53
List of Tables	55
List of Abbreviations	57
1. Introduction: The Coordination Crisis in AI Governance	59
1.1 Opening Scenario: The Policymaker’s Dilemma	59
1.2 The Coordination Crisis in AI Governance	60
1.2.1 Safety Gaps from Misaligned Efforts	60
1.2.2 Resource Misallocation	61
1.2.3 Negative-Sum Dynamics	61
1.3 Historical Parallels and Temporal Urgency	61
1.4 Research Question and Scope	62
1.5 The Multiplicative Benefits Framework	63
1.5.1 Automated Worldview Extraction	63
1.5.2 Live Data Integration	63
1.5.3 Formal Policy Evaluation	64
1.5.4 The Synergy	64
1.6 Thesis Structure and Roadmap	64
2. Context and Theoretical Foundations	67
2.1 AI Existential Risk: The Carlsmith Model	67
2.1.1 Six-Premise Decomposition	67
2.1.2 Why Carlsmith Exemplifies Formalizable Arguments	68
2.2 The Epistemic Challenge of Policy Evaluation	69
2.2.1 Unique Characteristics of AI Governance	69
2.2.2 Limitations of Traditional Approaches	69
2.2.3 Toward New Epistemic Tools	70

2.3 Bayesian Networks as Knowledge Representation	71
2.3.1 Mathematical Foundations	71
2.3.2 The Rain-Sprinkler-Grass Example	71
Rain-Sprinkler-Grass Network Rendering	72
2.3.3 Advantages for AI Risk Modeling	72
2.4 Argument Mapping and Formal Representations	72
2.4.1 From Natural Language to Structure	73
2.4.2 ArgDown: Structured Argument Notation	73
2.4.3 BayesDown: The Bridge to Bayesian Networks	73
2.5 The MTAIR Framework: Achievements and Limitations	74
2.5.1 MTAIR’s Approach	74
2.5.2 Key Achievements	75
2.5.3 Fundamental Limitations	75
2.5.4 The Automation Opportunity	76
2.6 Literature Review: Content and Technical Levels	76
2.6.1 AI Risk Models Evolution	76
2.6.2 Governance Proposals Taxonomy	76
2.6.3 Bayesian Network Theory and Applications	77
2.6.4 Software Tools Landscape	77
2.6.5 Formalization Approaches	77
2.6.6 Correlation Accounting Methods	77
2.7 Methodology	78
2.7.1 Research Design Overview	78
2.7.2 Formalizing World Models from AI Safety Literature	78
2.7.3 From Natural Language to Computational Models	78
2.7.4 Directed Acyclic Graphs: Structure and Semantics	79
2.7.5 Quantification of Probabilistic Judgments	80
2.7.6 Inference Techniques for Complex Networks	80
2.7.7 Integration with Prediction Markets and Forecasting Platforms	81
3. AMTAIR: Design and Implementation	83
3.1 System Architecture Overview	83
3.1.1 Five-Stage Pipeline Architecture	83
3.1.2 Design Principles	84
3.2 The Two-Stage Extraction Process	84
3.2.1 Stage 1: Structural Extraction (ArgDown)	84
3.2.2 Stage 2: Probability Integration (BayesDown)	84
3.2.3 Why Two Stages?	85
3.3 Implementation Technologies	85
3.3.1 Technology Stack	85
3.3.2 Key Algorithms	86
3.3.3 (Expected) Performance Characteristics	86
3.3.4 Deterministic vs. Probabilistic Components of the Workflow	86

3.4 Case Study: Rain-Sprinkler-Grass	86
3.4.1 Processing Steps	86
3.4.2 Example Conversion Steps	86
ArgDown Example	87
Example of Questions for BayesDown extraction	87
Complete BayesDown Example	88
Resulting Rain-Sprinkler-Grass DataFrame	89
3.4.3 Results	90
Rain-Sprinkler-Grass Network Rendering	90
3.5 Case Study: Carlsmith’s Power-Seeking AI Model	90
3.5.1 Model Complexity	90
3.5.2 Extraction Results	90
3.5.3 Validation Against Original (From the MTAIR Project)	91
3.5.4 Insights from Formalization	91
3.6 Validation Methodology	91
3.6.1 Ground Truth Construction	91
3.6.2 Evaluation Metrics	91
3.6.3 Results Summary	91
3.6.4 Error Analysis	92
3.7 Policy Evaluation Capabilities	92
3.7.1 Intervention Representation	92
3.7.2 Example: Deployment Governance	92
3.7.3 Robustness Analysis	92
3.8 Interactive Visualization Design	93
3.8.1 Visual Encoding Strategy	93
3.8.2 Progressive Disclosure	93
3.8.3 User Interface Elements	93
3.9 Integration with Prediction Markets	93
3.9.1 Design for Integration	93
3.9.2 Challenges and Opportunities	93
3.10 Computational Performance Analysis	94
3.10.1 Exact vs. Approximate Inference	94
3.10.2 Scaling Strategies	94
3.11 Results and Achievements	94
3.11.1 Extraction Quality Assessment	94
3.11.2 Computational Performance	94
3.11.3 Policy Impact Evaluation	94
3.12 Summary of Technical Contributions	94
4. Discussion: Implications and Limitations	97
4.1 Technical Limitations and Responses	97
4.1.1 Objection 1: Extraction Quality Boundaries	97
4.1.2 Objection 2: False Precision in Uncertainty	98

4.1.3 Objection 3: Correlation Complexity	98
4.2 Conceptual and Methodological Concerns	99
4.2.1 Objection 4: Democratic Exclusion	99
4.2.2 Objection 5: Oversimplification of Complex Systems	100
4.3 Red-Teaming Results	100
4.3.1 Adversarial Extraction Attempts	101
4.3.2 Robustness Findings	101
4.3.3 Implications for Deployment	101
4.4 Enhancing Epistemic Security	101
4.4.1 Making Models Inspectable	101
4.4.2 Revealing Convergence and Divergence	102
4.4.3 Improving Collective Reasoning	102
4.5 Scaling Challenges and Opportunities	103
4.5.1 Technical Scaling	103
4.5.2 Social and Institutional Scaling	103
4.5.3 Opportunities for Impact	103
4.6 Integration with Governance Frameworks	104
4.6.1 Standards Development	104
4.6.2 Regulatory Design	104
4.6.3 International Coordination	104
4.6.4 Organizational Decision-Making	104
4.7 Future Research Directions	104
4.7.1 Technical Enhancements	105
4.7.2 Methodological Extensions	105
4.7.3 Application Domains	105
4.7.4 Ecosystem Development	105
4.8 Known Unknowns and Deep Uncertainties	106
4.8.1 Categories of Deep Uncertainty	106
4.8.2 Adaptation Strategies for Deep Uncertainty	106
4.8.3 Robust Decision-Making Principles	106
5. Conclusion: Toward Coordinated AI Governance	109
5.1 Summary of Key Contributions	109
5.1.1 Theoretical Contributions	109
5.1.2 Methodological Innovations	110
5.1.3 Technical Achievements	110
5.1.4 Empirical Findings	110
5.2 Limitations and Honest Assessment	111
5.2.1 Technical Constraints	111
5.2.2 Conceptual Limitations	111
5.2.3 Practical Constraints	111
5.3 Implications for AI Governance	111
5.3.1 Near-Term Applications	112

5.3.2 Medium-Term Transformation	112
5.3.3 Long-Term Vision	112
5.4 Recommendations for Stakeholders	113
5.4.1 For Researchers	113
5.4.2 For Policymakers	113
5.4.3 For Technologists	113
5.4.4 For Funders	114
5.5 Future Research Agenda	114
5.5.1 Technical Priorities	114
5.5.2 Methodological Development	114
5.5.3 Application Expansion	115
5.6 Closing Reflections	115
References	117
References (.md)	123
Error Watch	123
Catch ALL Potential Hallucinations	123
Figure Inventory and Tracking	123
Chapter 1	124
Chapter 2	125
Pending Figures	125
Master Citation Registry	126
Bibliography	129
AMTAIR Thesis Relevant Literature & Citations	129
Items from MAref.bib	129
@carlsmith2021: Carlsmith [13]	129
@bostrom2014: Bostrom [9]	129
@clarke2022: Clarke et al. [18]	130
@pearl2009 and @pearl2000: Pearl [55] and Pearl [54]	130
@jaynes2003: Jaynes [37]	131
@tetlock2015: Tetlock and Gardner [68]	131
@lempert2003: Lempert, Popper, and Bankes [44]	132
@good1966: Good [29]	132
@yudkowsky2008: Yudkowsky [74]	132
@russell2015: Russell et al. [60]	133
New Suggested Citations	133
New Items to Consider:	133
@amodei2016: Amodei et al. [1]	133
@christiano2019: Christiano [17]	134
@critch2019: critch2019	134
@dafoe2018 and updated @dafoe2021: Dafoe [21] and Dafoe [20]	134

@askell12021: Askill et al. [4]	135
Further Citations to Integrate:	135
Appendices	137
Appendices	137
Appendices	139
Appendix A: Technical Implementation Details	139
A.1 Core Data Structures	139
A.2 Extraction Algorithm Details	139
A.3 API Specifications	139
Appendix B: Validation Datasets and Procedures	139
B.1 Expert Annotation Protocol	139
B.2 Benchmark Dataset Construction	139
B.3 Validation Results	139
Appendix C: Extended Case Studies	139
C.1 Christiano’s “What Failure Looks Like” Extraction	140
C.2 Critch’s ARCHES Model	140
C.3 Policy Evaluation: A Narrow Path	140
Appendix D: BayesDown Syntax Specification	140
Appendix E: Prompt Engineering Details	140
Appendix F: User Guide	140
Appendix G: Jupyter Notebook Implementation	140
Appendix H: Ethical Considerations and Governance	141
H.1 Potential Misuse Scenarios	141
H.2 Democratic Participation Frameworks	141
H.3 Responsibility Assignment	141
Appendix I: Full Extraction Examples	141
Appendix J: Software Installation and Usage Guide	141
AMTAIR Prototype Demonstration (Public Colab Notebook)	143
AMTAIR Prototype: Automating Transformative AI Risk Modeling	145
Executive Summary	145
Purpose Within the Master’s Thesis	145
Relevance to AI Governance	145
Notebook Structure and Workflow	145
Instructions — How to use this notebook:	146
Key Concepts:	149
Example Workflow:	149
Troubleshooting:	149
0 Environment Setup and Data Access	151

0.1 Prepare Colab/Python Environment — Import Libraries & Packages	151
0.2 Connect to GitHub Repository	153
0.3 File Import	157
1 Sources (PDF's of Papers) to ArgDown (.md file)	159
1.0 Sources to ArgDown: Structured Argument Extraction	159
Process Overview	159
What is ArgDown?	159
1.1 Specify Source Document (e.g. PDF)	160
1.2 Generate ArgDown Extraction Prompt	160
1.3 Prepare LLM API Call	166
1.4 Make ArgDown Extraction LLM API Call	178
1.5 Save ArgDown Extraction Response	180
1.6 Review and Check ArgDown.md File	182
1.6.0 Check the Graph Structure with the ArgDown Sandbox Online	183
1.7 Extract ArgDown Graph Information as DataFrame	184
1.8 Store ArgDown Information as 'ArgDown.csv' file	195
2 Probability Extractions: ArgDown (.csv) to BayesDown (.md + plugin JSON syntax)	201
2.0 ArgDown to BayesDown: Adding Probability Information	201
Process Overview	201
What is BayesDown?	201
2.1 Probability Extraction Questions — 'ArgDown.csv' to 'ArgDown_WithQuestions.csv'	202
2.2 'ArgDown_WithQuestions.csv' to 'BayesDownQuestions.md'	213
2.3 Generate BayesDown Probability Extraction Prompt	229
2.3.0 BayesDown Format Specification	229
Core Structure	229
3.0.1 Test BayesDown Extraction	232
3.0.2 Check the Graph Structure with the ArgDown Sandbox Online	236
3.1 Extraction	237
3.2 Data-Post-Processing	237
3.4 Download and save finished data frame as .csv file	245
4 Analysis & Inference: Bayesian Network Visualization	247
4.0 Bayesian Network Visualization Approach	247
Visualization Philosophy	247
Connection to AMTAIR Goals	247
Implementation Structure	247
4.1 Phase 1: Dependencies/Functions	248
4.2 Phase 2: Node Classification and Styling Module	253
4.3 Phase 3: HTML Content Generation Module	260
4.4 Phase 4: Main Visualization Function	265

5 Quick check HTML Outputs	271
Conclusion: From Prototype to Production	281
Summary of Achievements	281
Limitations and Future Work	281
Connection to AMTAIR Project	282
6 Save Outputs	283
6.0 Saving and Exporting Results	283
6.1 Convert .ipynb Notebook to Markdown	285
6.2 Convert Notebook to Markdown Documentation	286
6.3 Create PDF and Latex	288

List of Figures

1	AMTAIR extraction pipeline visualization	12
---	--	----

List of Tables

3	Main Caption	23
1	Demonstration of pipe table syntax	23
2	My Caption 1	23
5	Technology stack components	85

Preface

Abstract

The coordination crisis in AI governance presents a paradoxical challenge: unprecedented investment in AI safety coexists alongside fundamental coordination failures across technical, policy, and ethical domains. These divisions systematically increase existential risk. This thesis introduces AMTAIR (Automating Transformative AI Risk Modeling), a computational approach addressing this coordination failure by automating the extraction of probabilistic world models from AI safety literature using frontier language models. The system implements an end-to-end pipeline transforming unstructured text into interactive Bayesian networks through a novel two-stage extraction process that bridges communication gaps between stakeholders.

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

This thesis introduces AMTAIR (Automating Transformative AI Risk Modeling), a computational approach that addresses this coordination failure by automating the extraction of probabilistic world models from AI safety literature using frontier language models.

The AMTAIR system implements an end-to-end pipeline that transforms unstructured text into interactive Bayesian networks through a novel two-stage extraction process: first capturing argument structure in ArgDown format, then enhancing it with probability information in BayesDown. This approach bridges communication gaps between stakeholders by making implicit models explicit, enabling comparison across different worldviews, providing a common language for discussing probabilistic relationships, and supporting policy evaluation across diverse scenarios.

Prefatory Apparatus: Frontmatter

Illustrations and Terminology — Quick References

Acknowledgments

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

List of Graphics & Figures

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

List of Abbreviations

esp. especially

f., ff. following

incl. including

p., pp. page(s)

MAD Mutually Assured Destruction

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

Glossary

- **Argument mapping:** A method for visually representing the structure of arguments
- **BayesDown:** An extension of ArgDown that incorporates probabilistic information
- **Bayesian network:** A probabilistic graphical model representing variables and their dependencies
- **Conditional probability:** The probability of an event given that another event has occurred
- **Directed Acyclic Graph (DAG):** A graph with directed edges and no cycles
- **Existential risk:** Risk of permanent curtailment of humanity's potential
- **Power-seeking AI:** AI systems with instrumental incentives to acquire resources and power
- **Prediction market:** A market where participants trade contracts that resolve based on future events

- **d-separation:** A criterion for identifying conditional independence relationships in Bayesian networks
- **Monte Carlo sampling:** A computational technique using random sampling to obtain numerical results

Quarto Syntax and Best Practices Guide

Key Features

1. Task Management System

- HTML comments with [] for tasks visible in ToDo-Tree
- Categories: FIND, VERIFY, CREATE, TODO
- Progress tracking with [x] (done) and [-] (verified)

2. Multi-Format Output

- HTML: Interactive web version with navigation
- PDF: Professional academic document
- LaTeX: Source for further customization
- DOCX: For collaboration

3. Cross-Referencing

- Sections: @sec-section-name
- Figures: @fig-figure-name
- Tables: @tbl-table-name
- Citations: @citation-key

4. Advanced Features

- Interactive Jupyter notebooks
- Mermaid diagrams
- Math equations (LaTeX)
- Callout blocks
- Extensive footnotes
- Glossary and abbreviations

Quick Start

Task Management

Write and track tasks with HTML comments in markdown blocks or with `verbatim` code ticks but ALWAYS add line breaks between tasks:

```
`<!-- [ ] TODO: Task description -->`

`<!-- [ ] FIND: @missing-citation: "Description" -->`

`<!-- [ ] VERIFY: @suggested-citation: "Source" -->`

`<!-- [ ] CREATE: {#fig-name}: "Figure description" -->`
```

Adding Content

1. Create/edit `.qmd` files in `chapters/`
2. Update `_quarto.yml` if adding new chapters
3. Add citations to `ref/MAref.bib`
4. Place images in `images/`

Best Practices

1. Consistent Formatting

- Use American spelling throughout
- Follow heading hierarchy (`##`, `###`, `####`)
- Maintain consistent citation style
- Use semantic line breaks

2. Task Tracking

- Create tasks as you write
- Update task status regularly
- Use categories for clarity
- Include implementation details

3. Version Control

- Commit frequently with descriptive messages
- Use branches for major revisions
- Tag releases (draft versions)

4. Documentation

- Comment complex code blocks
- Provide alt text for all figures
- Keep this README updated
- Document any custom scripts

Troubleshooting

Common Issues

1. **LaTeX errors:** Check `_quarto.yml` for LaTeX settings
2. **Missing references:** Ensure citations are in `MAREF.bib`
3. **Broken links:** Use relative paths for internal links
4. **Task visibility:** Install Todo-Tree extension in VS Code

Getting Help

- Quarto documentation: <https://quarto.org>
- Project repository: <https://github.com/VJMeyer/submission>
- Contact: Valentin2meyer@gmail.com

License

MIT License - See LICENSE file for details

Document Structure and Headings

Heading Hierarchy

Always use the full heading hierarchy for maximum organization:

markdown

```
# Chapter Title {#sec-chapter}
## Major Section {#sec-major-section}
### Subsection {#sec-subsection}
#### Sub-subsection {#sec-subsubsection}
`##### Sub-subsubsection {#sec-subsubsubsection}`
`##### Sub-subsubsubsection {#sec-subsubsubsubsection}`
```

Best Practices:

- Always include `{#sec-label}` for cross-referencing
- Use descriptive, concise heading names
- Maintain consistent capitalization (Title Case for chapters, Sentence case for sections)
- Add `.unnumbered` for sections without numbers (e.g., References)

- Add `.unlisted` to exclude from TOC
- Do not manually number headings

Text Formatting

Basic Formatting

markdown

```
*italics* for emphasis
**bold** for strong emphasis
***bold italics*** for very strong emphasis
~~strikethrough~~ for deleted text
[highlighted text]{.mark}
[underlined text]{.underline}
[small caps]{.smallcaps}
`inline code` in numerous applications
```

Advanced Formatting

markdown

```
superscript^2^ for exponents
subscript~2~ for chemical formulas
```

Links

`<https://quarto.org/docs/authoring/markdown-basics.html>` produces: <https://quarto.org/docs/authoring/markdown-basics.html>

[Quarto Book Cross-References] (<https://quarto.org/docs/books/book-crossrefs.html>) produces: Quarto Book Cross-References

Including Code

```
import pandas as pd
print("AMTAIR is working!")
```

AMTAIR is working!

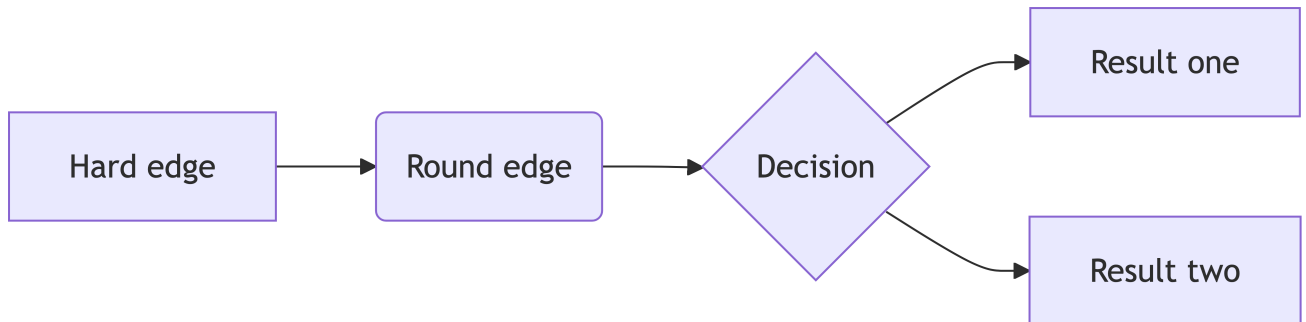
Figure 1: AMTAIR extraction pipeline visualization

Diagrams

Quarto has native support for embedding Mermaid and Graphviz diagrams. This enables you to create flowcharts, sequence diagrams, state diagrams, Gantt charts, and more using a plain text syntax inspired by markdown.

For example, here we embed a flowchart created using Mermaid:

```
flowchart LR
  A[Hard edge] --> B(Round edge)
  B --> C{Decision}
  C --> D[Result one]
  C --> E[Result two]
```



In-Line LaTeX

In-Line HTML

Here's some raw inline HTML: html

Reference or Embed Code from .ipynb files

Code chunks from .ipynb notebooks can be embedded in the .qmd text with:

```
{{< embed /AMTAIR_Prototype/data/example_carlsmith/AMTAIR_Prototype_example_carlsmith.ipynb#
```

which produces the output of executing the code cell:

```
# @title 0.2.0 --- Connect to GitHub Repository --- Load Files [connect_to_github_repository]

"""
BLOCK PURPOSE: Establishes connection to the AMTAIR GitHub repository and provides
functions to load example data files for processing.

This block creates a reusable function for accessing files from the project's
GitHub repository, enabling access to example files like the rain-sprinkler-lawn
Bayesian network that serves as our canonical test case.

DEPENDENCIES: requests library, io library
OUTPUTS: load_file_from_repo function and test file loads
"""
```

```
from requests.exceptions import HTTPError

# Specify the base repository URL for the AMTAIR project
repo_url = "https://raw.githubusercontent.com/SingularitySmith/AMTAIR_Prototype/main/data/ex
print(f"Connecting to repository: {repo_url}")

def load_file_from_repo(relative_path):
    """
    Loads a file from the specified GitHub repository using a relative path.

    Args:
        relative_path (str): Path to the file relative to the repo_url

    Returns:
        For CSV/JSON: pandas DataFrame
        For MD: string containing file contents

    Raises:
        HTTPError: If file not found or other HTTP error occurs
        ValueError: If unsupported file type is requested
    """
    file_url = repo_url + relative_path
    print(f"Attempting to load: {file_url}")

    # Fetch the file content from GitHub
    response = requests.get(file_url)

    # Check for bad status codes with enhanced error messages
    if response.status_code == 404:
        raise HTTPError(f"File not found at URL: {file_url}. Check the file path/name and en
    else:
        response.raise_for_status() # Raise for other error codes

    # Convert response to file-like object
    file_object = io.StringIO(response.text)

    # Process different file types appropriately
    if relative_path.endswith(".csv"):
        return pd.read_csv(file_object) # Return DataFrame for CSV
    elif relative_path.endswith(".json"):
        return pd.read_json(file_object) # Return DataFrame for JSON
    elif relative_path.endswith(".md"):
```



```

    return file_object.read() # Return raw content for MD files
else:
    raise ValueError(f"Unsupported file type: {relative_path.split('.')[-1]}. Add support")

# Load example files to test connection
try:
    # Load the extracted data CSV file
    # df = load_file_from_repo("extracted_data.csv")

    # Load the ArgDown test text
    md_content = load_file_from_repo("ArgDown.md")

    print(" Successfully connected to repository and loaded test files.")
except Exception as e:
    print(f" Error loading files: {str(e)}")
    print("Please check your internet connection and the repository URL.")

# Display preview of loaded content (commented out to avoid cluttering output)
print(md_content)

```

Connecting to repository: https://raw.githubusercontent.com/SingularitySmith/AMTAIR_Prototype/main/extracted_data.csv
 Attempting to load: https://raw.githubusercontent.com/SingularitySmith/AMTAIR_Prototype/main/ArgDown.md
 Successfully connected to repository and loaded test files.

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

- [Human_Disempowerment]: Permanent and collective disempowerment of humanity relative to AI
 - [Scale_Of_Power_Seeking]: Power-seeking by AI systems scaling to the point of permanent domination
 - [Misaligned_Power_Seeking]: Deployed AI systems seeking power in unintended and harmful ways
 - [APS_Systems]: AI systems with advanced capabilities, agentic planning, and strategic awareness
 - [Advanced_AI_Capability]: AI systems that outperform humans on tasks that require complex reasoning
 - [Agentic_Planning]: AI systems making and executing plans based on world models
 - [Strategic_Awareness]: AI systems with models accurately representing power dynamics
 - [Difficulty_Of_Alignment]: It is harder to build aligned systems than misaligned systems
 - [Instrumental_Convergence]: AI systems with misaligned objectives tend to pursue common instrumental goals
 - [Problems_With_Proxies]: Optimizing for proxy objectives breaks correlations with the true objective
 - [Problems_With_Search]: Search processes can yield systems pursuing different objectives
 - [Deployment_Decisions]: Decisions to deploy potentially misaligned AI systems
 - [Incentives_To_Build_APS]: Strong incentives to build and deploy APS systems
 - [Usefulness_Of_APS]: APS systems are very useful for many valuable tasks
 - [Competitive_Dynamics]: Competitive pressures between AI developers
 - [Deception_By_AI]: AI systems deceiving humans about their true objectives
 - [Corrective_Feedback]: Human society implementing corrections after observing problems
 - [Warning_Shots]: Observable failures in weaker systems before catastrophic risk

- [Rapid_Capability_Escalation]: AI capabilities escalating very rapidly, allowing
- [Barriers_To_Understanding]: Difficulty in understanding the internal workings of advanced AI
- [Misaligned_Power_Seeking]: Deployed AI systems seeking power in unintended and high-impact
- [Adversarial_Dynamics]: Potentially adversarial relationships between humans and power-seeking
- [Misaligned_Power_Seeking]: Deployed AI systems seeking power in unintended and high-impact
- [Stakes_Of_Error]: The escalating impact of mistakes with power-seeking AI systems. {"instant"
- [Misaligned_Power_Seeking]: Deployed AI systems seeking power in unintended and high-impact

including ‘echo=true’ renders the code of the cell:

```
{{< embed /AMTAIR_Prototype/data/example_carlsmith/AMTAIR_Prototype_example_carlsmith.ipynb#
# @title 0.2.0 --- Connect to GitHub Repository --- Load Files [connect_to_github_repository

"""
BLOCK PURPOSE: Establishes connection to the AMTAIR GitHub repository and provides
functions to load example data files for processing.

This block creates a reusable function for accessing files from the project's
GitHub repository, enabling access to example files like the rain-sprinkler-lawn
Bayesian network that serves as our canonical test case.

DEPENDENCIES: requests library, io library
OUTPUTS: load_file_from_repo function and test file loads
"""

from requests.exceptions import HTTPError

# Specify the base repository URL for the AMTAIR project
repo_url = "https://raw.githubusercontent.com/SingularitySmith/AMTAIR_Prototype/main/data/ex
print(f"Connecting to repository: {repo_url}")

def load_file_from_repo(relative_path):
    """
    Loads a file from the specified GitHub repository using a relative path.

    Args:
        relative_path (str): Path to the file relative to the repo_url

    Returns:
        For CSV/JSON: pandas DataFrame
        For MD: string containing file contents
```

```

Raises:
    HTTPError: If file not found or other HTTP error occurs
    ValueError: If unsupported file type is requested
"""

file_url = repo_url + relative_path
print(f"Attempting to load: {file_url}")

# Fetch the file content from GitHub
response = requests.get(file_url)

# Check for bad status codes with enhanced error messages
if response.status_code == 404:
    raise HTTPError(f"File not found at URL: {file_url}. Check the file path/name and en
else:
    response.raise_for_status() # Raise for other error codes

# Convert response to file-like object
file_object = io.StringIO(response.text)

# Process different file types appropriately
if relative_path.endswith(".csv"):
    return pd.read_csv(file_object) # Return DataFrame for CSV
elif relative_path.endswith(".json"):
    return pd.read_json(file_object) # Return DataFrame for JSON
elif relative_path.endswith(".md"):
    return file_object.read() # Return raw content for MD files
else:
    raise ValueError(f"Unsupported file type: {relative_path.split('.')[-1]}. Add support

# Load example files to test connection
try:
    # Load the extracted data CSV file
    # df = load_file_from_repo("extracted_data.csv")

    # Load the ArgDown test text
    md_content = load_file_from_repo("ArgDown.md")

    print(" Successfully connected to repository and loaded test files.")
except Exception as e:
    print(f" Error loading files: {str(e)}")
    print("Please check your internet connection and the repository URL.")

```

```
# Display preview of loaded content (commented out to avoid cluttering output)
print(md_content)
```

Connecting to repository: https://raw.githubusercontent.com/SingularitySmith/AMTAIR_Prototype/main

Attempting to load: https://raw.githubusercontent.com/SingularitySmith/AMTAIR_Prototype/main

Successfully connected to repository and loaded test files.

[Existential_Catastrophe]: The destruction of humanity's long-term potential due to AI systems.

- [Human_Disempowerment]: Permanent and collective disempowerment of humanity relative to AI systems.

- [Scale_Of_Power_Seeking]: Power-seeking by AI systems scaling to the point of permanent domination.

- [Misaligned_Power_Seeking]: Deployed AI systems seeking power in unintended and high-impact ways.

- [APS_Systems]: AI systems with advanced capabilities, agentic planning, and strategic awareness.

- [Advanced_AI_Capability]: AI systems that outperform humans on tasks that require advanced capabilities.

- [Agentic_Planning]: AI systems making and executing plans based on world models and long-term goals.

- [Strategic_Awareness]: AI systems with models accurately representing power dynamics and human behavior.

- [Difficulty_Of_Alignment]: It is harder to build aligned systems than misaligned systems.

- [Instrumental_Convergence]: AI systems with misaligned objectives tend to converge on similar instrumental goals.

- [Problems_With_Proxies]: Optimizing for proxy objectives breaks correlations between proxy and true objectives.

- [Problems_With_Search]: Search processes can yield systems pursuing different instrumental goals.

- [Deployment_Decisions]: Decisions to deploy potentially misaligned AI systems.

- [Incentives_To_Build_APS]: Strong incentives to build and deploy APS systems.

- [Usefulness_Of_APS]: APS systems are very useful for many valuable tasks.

- [Competitive_Dynamics]: Competitive pressures between AI developers.

- [Deception_By_AI]: AI systems deceiving humans about their true objectives.

- [Corrective_Feedback]: Human society implementing corrections after observing problems.

- [Warning_Shots]: Observable failures in weaker systems before catastrophic risk.

- [Rapid_Capability_Escalation]: AI capabilities escalating very rapidly, allowing for rapid adaptation.

[Barriers_To_Understanding]: Difficulty in understanding the internal workings of advanced AI systems.

- [Misaligned_Power_Seeking]: Deployed AI systems seeking power in unintended and high-impact ways.

[Adversarial_Dynamics]: Potentially adversarial relationships between humans and power-seeking AI systems.

- [Misaligned_Power_Seeking]: Deployed AI systems seeking power in unintended and high-impact ways.

[Stakes_Of_Error]: The escalating impact of mistakes with power-seeking AI systems. {"instantaneous": "catastrophic", "long-term": "catastrophic"}

- [Misaligned_Power_Seeking]: Deployed AI systems seeking power in unintended and high-impact ways.

Link:

Full Notebooks are embedded in the Appendix through the `_quarto.yml` file with:

Embed .html result/rendering from .ipynb Notebook

Html Graph by Notebook Cell Inclusion - (from github-pages)

```
{< embed /AMTAIR_Prototype/data/example_carlsmith/AMTAIR_Prototype_example_carlsmith.ipynb#
```

```
from IPython.display import IFrame

IFrame(src="https://singularitysmith.github.io/AMTAIR_Prototype/bayesian_network_carlsmith.html")
```

```
<IPython.lib.display.IFrame at 0x7f04d69f0d90>
```

Dynamic Html Rendering of the Carlsmith Bayesian Network/DAG Visualization

Html Graph by Notebook Cell Inclusion with Website Call?

https://singularitysmith.github.io/AMTAIR_Prototype/bayesian_network_carlsmith.html

Full Bayesian Network Rendering

```
{{< embed /AMTAIR_Prototype/data/example_carlsmith/AMTAIR_Prototype_example_carlsmith.ipynb}}
```

Rain-Sprinkler-Grass Network Rendering

```
from IPython.display import IFrame

IFrame(src="https://singularitysmith.github.io/AMTAIR_Prototype/bayesian_network.html", width=1000, height=500)
```

```
<IPython.lib.display.IFrame at 0x106661a90>
```

Dynamic Html Rendering of the Rain-Sprinkler-Grass DAG

Lists and Enumerations

Unordered Lists

markdown

```
- First level item
  - Second level item (2 spaces)
    - Third level item (4 spaces)
- Another first level item
  with continuation (2 spaces for alignment)
```

Ordered Lists

markdown

```
1. First item
2. Second item
  a) Sub-item (3 spaces)
    i. Sub-sub-item (6 spaces)
```

```
b) Another sub-item
3. Third item
```

Definition Lists

markdown

```
Term One
: Definition of term one with detailed explanation
  that can span multiple lines

Term Two
: Brief definition

Term Three
: Another definition with multiple paragraphs

  Additional paragraph for term three
```

Code Blocks and Verbatim Text

Inline Code

markdown

```
Use `print("Hello")` for inline code
```

Code Blocks with Syntax Highlighting

markdown

```
```python
def calculate_risk(probability, impact):
 """Calculate risk score from probability and impact."""
 return probability * impact
```
```

Verbatim Text

markdown

This is verbatim text that preserves all spacing and formatting exactly as typed

Blockquotes and Callouts

Simple Blockquote

markdown

```
> This is a blockquote for citations or important quotes.  
> It can span multiple lines.  
>  
> And include multiple paragraphs.
```

Callout Blocks

! With Callout blocks it is crucial to always have a line break after the title and the ::: in a new line after the note ! markdown

```
::: {.callout-note}  
## Note Title  
This is a note callout with important information.  
:::  
  
::: {.callout-warning}  
## Warning  
This warns about potential issues.  
:::  
  
::: {.callout-tip}  
## Pro Tip  
Helpful suggestions go here.  
:::  
  
::: {.callout-important}  
## Important  
Critical information that must not be missed.  
:::  
  
::: {.callout-caution}  
## Caution  
Use with care in specific situations.  
:::
```

Figures and Images

Complete Figure Syntax

markdown

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

Figure Best Practices

1. Always include comprehensive alt text
2. Use descriptive filenames
3. Optimize image sizes for web/PDF
4. Maintain consistent styling
5. Reference all figures in text: See @fig-identifier

Tables

Markdown Tables

markdown

```
Column 1	Column 2	Column 3
Left	Center	Right
Data	Data	Data

: Table caption {#tbl-identifier}
```

Grid Tables

markdown

```
+-----+-----+-----+
| Header 1 | Header 2 | Header 3 |
```


Table 3: Main Caption

| (a) First Table | | | (b) Second Table | | |
|-----------------|------|------|------------------|------|------|
| Col1 | Col2 | Col3 | Col1 | Col2 | Col3 |
| A | B | C | A | B | C |
| E | F | G | E | F | G |
| A | G | G | A | G | G |

```
+=====+=====+=====+
Cell 1	Cell 2	Cell 3
Multi-	Multi-	Multi-
line	line	line
+-----+-----+-----+

: Complex table with multiple lines {#tbl-complex}
```

Table 1: Demonstration of pipe table syntax

| Right | Left | Default | Center |
|-------|------|---------|--------|
| 12 | 12 | 12 | 12 |
| 123 | 123 | 123 | 123 |
| 1 | 1 | 1 | 1 |

Table 2: My Caption 1

| Col1 | Col2 | Col3 |
|------|------|------|
| A | B | C |
| E | F | G |
| A | G | G |

Referencing tables with `@tbl-KEY`: See Table 2.

See Table 3 for details, especially Table 3b.

Citations and References

Citation Styles

markdown

```
Narrative: @author2024 argues that...
Parenthetical: This is supported by evidence [@author2024].
```

Multiple: Several studies confirm this [@author2024; @other2023].
 Page specific: See discussion in [@author2024, pp. 45-67].
 Author only: As [-@author2024] demonstrates...

Bibliography Entry

bibtex

```
@article{author2024,
  title = {Article Title},
  author = {Author, First and Other, Second},
  date = {2024},
  journaltitle = {Journal Name},
  volume = {10},
  number = {2},
  pages = {45--67},
  doi = {10.1234/example},
  url = {https://example.com}
}
```

Cross-References

Section References

markdown

See @sec-introduction for background.
 As discussed in @sec-methodology...

Figure and Table References

markdown

@fig-pipeline shows the workflow.
 Results are summarized in @tbl-results.

Equation References

markdown

```
$$
E = mc^2
$$ {#eq-einstein}

Einstein's equation (@eq-einstein) shows...
```

Mathematics

Inline Math

markdown

```
The probability  $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$ 
```

Display Math

markdown

```
$$
\begin{align}
\nabla \times \vec{\mathbf{B}} &= \frac{1}{c} \frac{\partial \vec{\mathbf{E}}}{\partial t} \\
\nabla \cdot \vec{\mathbf{E}} &= 4 \pi \rho \\
\nabla \times \vec{\mathbf{E}} &= -\frac{1}{c} \frac{\partial \vec{\mathbf{B}}}{\partial t} \\
\nabla \cdot \vec{\mathbf{B}} &= 0
\end{align}
$$
```

inline math: $E = mc^2$

display math:

$$E = mc^2$$

If you want to define custom TeX macros, include them within $\mathbb{}$ delimiters enclosed in a .hidden block. For example:

For HTML math processed using MathJax (the default) you can use the `\def`, `\newcommand`, `\renewcommand`, `\newenvironment`, `\renewenvironment`, and `\let` commands to create your own macros and environments.

Footnotes

Footnotes are to be used as much as possible!

Simple Footnote

markdown

```
This needs clarification.[This is an inline footnote.]
```

Referenced Footnote

markdown

```
This is important.[^1]

[^1]: This is a longer footnote with multiple paragraphs.

    Second paragraph of the footnote.

    Even code blocks are possible:
    ```python
 print("In footnote")
    ```
```

Here is an inline note.¹

Here is a footnote reference,²

Another Text with a footnote³ but this time a “longnote”.

This paragraph won’t be part of the note, because it isn’t indented.

Appendices

Structure

markdown

```
# Appendices {.unnumbered}

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

### A.1 Implementation {.unnumbered}

## Appendix B: Additional Results {#sec-appendix-b .unnumbered}
```

Best Practices for Appendices

1. Include all supplementary material
2. Reference from main text

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.

³Here’s one with multiple blocks.

Subsequent paragraphs are indented to show that they belong to the previous footnote.

```
{ some.code }
```

The whole paragraph can be indented, or just the first line. In this way, multi-paragraph footnotes work like multi-paragraph list items.

3. Number consistently
4. Provide clear descriptions
5. Maintain same formatting standards

Glossary and Abbreviations

Glossary Format

markdown

```
# Glossary {.unnumbered}
Term
: Definition

AI
: Artificial Intelligence - Computer systems performing tasks requiring human intelligence

ML
: Machine Learning - Algorithms that improve through experience

DL
: Deep Learning - Neural networks with multiple layers
```

Interactive Elements

Jupyter Notebook Embedding

```
{{< embed notebook.ipynb#cell-label >}}
```

Mermaid Diagrams

```
```{mermaid}
graph TD
 A[Start] --> B{Decision}
 B -->|Yes| C[Action 1]
 B -->|No| D[Action 2]
 C --> E[End]
 D --> E
```

## Line Breaks and Spacing

### Spacing Rules

1. **Between sections:** 2 blank lines
2. **Between paragraphs:** 1 blank line
3. **Around code blocks:** 1 blank line before and after
4. **Around figures/tables:** 1 blank line before and after
5. **After headings:** 1 blank line
6. **Between list items:** No blank lines unless containing multiple paragraphs

### Page Breaks

## Comments and Metadata

### HTML Comments

```
<!-- This is a comment not shown in output -->
```

# Comprehensive Task Management System for Quarto Thesis

## Overview

This task management system uses HTML comments with specific formatting to create trackable, categorized tasks that integrate with VS Code's `ToDo-Tree` extension while remaining invisible or visible depending on the status in rendered output.

## Task Categories and Syntax

Write and track tasks with HTML comments in markdown blocks or with `verbatim` code ticks but ALWAYS add line breaks between tasks:

```
`<!-- [] TODO: Task description -->`

`<!-- [] FIND: @missing-citation: "Description" -->`

`<!-- [] VERIFY: @suggested-citation: "Source" -->`

`<!-- [] CREATE: {#fig-name}: "Figure description" -->`
```

### 1. General Tasks

In markdown blocks or with `verbatim` code ticks:

```
<!-- [] TODO: General task description -->

<!-- [] TODO: High-priority task with deadline (2024-12-31) -->

<!-- [] TODO: Task with subtasks
 - [] Subtask 1
 - [] Subtask 2
 - [] Subtask 3
```

```
-->
```

## 2. Citation Tasks

In markdown blocks or with `verbatim` code ticks:

```
`<!-- [] FIND: @missing-citation-key: "Description of needed source, keywords, search terms" -->`

`<!-- [] VERIFY: @suggested-citation: "Author (Year). Title. Journal." [Include BibTeX if a] -->`

`<!-- [] UPDATE: @outdated-citation: "Check for newer edition or updated data" -->`

`<!-- [] VERIFIED: @citation: "URL" -->`
```

## 3. Figure/Graphic Tasks

In markdown blocks or with `verbatim` code ticks:

```
<!-- [] CREATE: {#fig-diagram-name}: "Description of needed diagram, style, data to include" -->

<!-- [] FIND: {#fig-example-image}: "Stock photo of X, preferably showing Y" -->

<!-- [] UPDATE: {#fig-outdated-chart}: "Update with 2024 data" -->

<!-- [] IMPROVE: {#fig-low-quality}: "Higher resolution version needed" -->
```

## 4. Content Tasks

In markdown blocks or with `verbatim` code ticks:

```
<!-- [] WRITE: Section 3.2 - Methodology details -->

<!-- [] EXPAND: Background section needs 500 more words -->

<!-- [] REVISE: Introduction for clarity and flow -->

<!-- [] REVIEW: Chapter 4 for consistency -->
```

## 5. Technical Tasks

In markdown blocks or with `verbatim` code ticks:

```
<!-- [] FIX: Code block in section 2.3 has syntax error -->

<!-- [] TEST: Jupyter notebook embedding -->
```



```
<!-- [] OPTIMIZE: Large figure file sizes -->

<!-- [] IMPLEMENT: Cross-reference checking script -->
```

## Task States

### Open or In-Progress Tasks

In markdown blocks or with `verbatim` code ticks:

```
<!-- [] Task description -->
```

### Completed Tasks (Visible in ToDo-Tree)

Either markdown blocks or `verbatim` code ticks or without (to remain hidden in output):

```
<!-- [x] Task description (completed 2024-01-20) -->
```

### Verified/Archived Tasks (Hidden from ToDo-Tree)

markdown

```
<!-- [-] Task description (verified and archived) -->
```

## Advanced Task Formatting

### Multi-line Tasks with Details

markdown

```
<!-- [] Major task with extensive details
```

Background:

- Context for why this task exists
- Related issues or dependencies

Requirements:

1. Specific requirement one
2. Specific requirement two
3. Specific requirement three

Implementation Plan:

- [ ] Step 1: Initial research
- [ ] Step 2: Draft content
- [ ] Step 3: Review and revise

```
Resources:
```

- Reference document: path/to/doc
- Example: url-to-example

```
-->
```

## Linked Tasks

markdown

```
<!-- [] PRIMARY: Main task description
```

```
 Related tasks:
```

- See also: Task in Chapter 2
- Depends on: Task in Appendix A
- Blocks: Task in Chapter 5

```
-->
```

## Conditional Tasks

markdown

```
<!-- [] IF: Hypothesis confirmed in Chapter 3
```

```
 THEN: Add supporting evidence section
```

```
 ELSE: Revise theoretical framework -->
```

# Task Tracking Best Practices

## 1. Task Creation Guidelines

- > Create tasks immediately when identified
- > Be specific and actionable
- > Include context and success criteria
- > Link related tasks

## 2. Task Organization

- > Group related tasks together
- > Place tasks near or inside relevant content
- > Use consistent formatting
- > Maintain task hierarchy

## 3. Priority System

In markdown blocks or with `verbatim` code ticks:

```
<!-- [] URGENT: Task needing immediate attention -->
```

```
<!-- [] HIGH: Task important for next milestone -->
```

```
<!-- [] MEDIUM: Standard priority task -->
```

```
<!-- [] LOW: Nice-to-have improvement -->
```

#### #### Simple "One-line tasks"

Use Code ticks and html comment and task format for tasks distinctly visible across all forms

```
<!-- [] ToDos for things to do / tasks / reminders (allows "jump to with Taks Tree extension"
```

Use html comment and task format for open or uncertain tasks, visible in the .qmd file:

```
<!-- [] ToDos for things to do / tasks / reminders (allows "jump to with Taks Tree extension"
```

#### #### More Complex Tasks with Notes

More Information about task

Relevant notes

Step-by-step implementation Plan

Etc.

#### #### Completed Tasks

Retain completed tasks in ToDo-Tree by adding an x in the brackets: `[x]`

```
<!-- [x] Tasks which have been finished but should remain for later verification -->
```

```
<!-- [x] Tasks which have been finished but should remain for later verification -->
```

Mark and remove completed tasks from ToDo-Tree by adding a minus in the brackets: `[-]`

```
<!-- [-] Tasks which have been finished but should remain visible for later verification -->
```

```
<!-- [-] Tasks which have been finished but should remain for later verification (only in .qmd)
```

## Task Management Workflow

### 1. Task Creation

In markdown blocks or with `verbatim code` ticks:

```
<!-- [] TODO: Write introduction paragraph
 Context: Need to introduce the concept of X
 Requirements:
 - Define key terms
 - Provide historical context
 - Connect to thesis argument
 Deadline: 2024-02-15
-->
```

### 2. Task Execution

In markdown blocks or with `verbatim code` ticks:

```
<!-- [] TODO: Write introduction paragraph
 Progress:
 - [x] Defined key terms
 - [-] Not Working on historical context
 - [] Connection to thesis argument
-->
```

### 3. Task Completion

In markdown blocks or with `verbatim code` ticks:

```
<!-- [x] TODO: Write introduction paragraph (completed 2024-02-14)
 Final version includes all requirements
 Word count: 523
 Review status: Approved by advisor
-->
```

### 4. Task Archival

In markdown blocks or with `verbatim code` ticks:

```
<!-- [-] TODO: Write introduction paragraph (archived 2024-02-20)
 Moved to version control history
-->
```

## Best Practices Summary — ALWAYS CONSISTENTLY:

1. **Be Specific:** Tasks should be actionable and measurable

2. **Stay Organized:** Group related tasks and maintain hierarchy
3. **Archive Completed:** Keep task list manageable
4. **Use Categories:** Leverage task types for better organization
5. **Add Context:** Include enough detail for future reference
6. **Link Related:** Connect interdependent tasks
7. **Maintain Consistency:** Use standard formatting throughout
8. **Use correct formatting:** Deploy the correct formatting and fix any inconsistencies
9. **Always add extra line breaks:** Add additional line breaks between and around tasks



# Tagging and Highlighting System for Content Merging

## Overview

When merging content from multiple sources, it's crucial to identify and manage duplicate, redundant, or overlapping material. This system uses Quarto formatting features to clearly mark such content for review and consolidation.

## Tagging Categories

### A. Duplicate Content Marking

In markdown blocks or with `verbatim` code ticks:

```
::: {.duplicate-content data-source="Chapter2.qmd" data-section="2.3"}
This paragraph appears to be duplicated from Chapter 2, Section 2.3.
Consider consolidating or removing.
:::

`<!-- DUPLICATE: This content also appears in Section 2.3 -->`
```

### B. Redundant Content Highlighting

In markdown blocks or with `verbatim` code ticks:

```
::: {.redundant-content}
[This section covers similar ground to Section 3.2 but with less detail]{.mark style="background-color: #f0f0f0"}
:::

<!-- REDUNDANT: Similar content in Section 3.2 with more comprehensive coverage -->
```

### C. Better Version Available

In markdown blocks or with `verbatim` code ticks:

```

::: {.superseded-content data-better-version="Chapter4.qmd#sec-4-5"}
This explanation is superseded by a more comprehensive version in Chapter 4, Section 4.5
:::

<!-- SUPERSEDED: See Chapter 4.5 for improved version -->

```

## D. Merge Candidate

In markdown blocks or with `verbatim` code ticks:

```

::: {.merge-candidate data-merge-with="Section 5.2"}
MERGE CANDIDATE: This content could be combined with Section 5.2 for better flow.

Original content here...
:::

<!-- MERGE: Consider combining with Section 5.2 -->

```

## Visual Marking System

### Color-Coded Highlighting

In markdown blocks or with `verbatim` code ticks:

```

[Duplicate content - exact match]{style="background-color: #ff6b6b; color: white"}
[Redundant content - similar coverage]{style="background-color: #ffeb3b"}
[Better version exists elsewhere]{style="background-color: #4ecdc4"}
[Merge candidate]{style="background-color: #45b7d1"}
[Review needed]{style="background-color: #fa8231"}

```

### Border Marking

In markdown blocks or with `verbatim` code ticks:

```

::: {style="border-left: 5px solid #ff6b6b; padding-left: 10px"}
This entire section is duplicated elsewhere.
:::

```

### Inline Marking

In markdown blocks or with `verbatim` code ticks:

```

This paragraph contains [duplicate phrase]{.duplicate} that appears
in multiple locations.

```



## Metadata Tracking

### Comprehensive Metadata

In markdown blocks or with `verbatim` code ticks:

```

::: {.content-status
 data-status="duplicate"
 data-original-source="intro.qmd#para-3"
 data-other-locations="chapter2.qmd#para-15, chapter5.qmd#para-8"
 data-recommendation="keep-original"
 data-reviewed-by="VM"
 data-review-date="2024-02-15"}
This content appears in multiple locations.
The original in intro.qmd is most comprehensive.
:::

```

### Quick Reference Tags

In markdown blocks or with `verbatim` code ticks:

```

<!--
 STATUS: Duplicate
 ORIGINAL: intro.qmd#para-3
 ALSO IN: ch2#para-15, ch5#para-8
 ACTION: Remove this version
-->

```

## Workflow for Content Merging

### 1. Initial Marking Phase

In markdown blocks or with `verbatim` code ticks:

```

<!-- PHASE 1: Initial marking -->
<!-- [] TODO: Mark all duplicate content in Chapter 1 -->
<!-- [] TODO: Identify redundant sections in Chapter 2 -->
<!-- [] TODO: Tag better versions throughout document -->

```

### 2. Review and Comparison

In markdown blocks or with `verbatim` code ticks:

```

<!-- COMPARISON NEEDED -->
::: {.comparison-block}
Version A (Current):
Brief explanation of concept X.

```

```

Version B (Chapter 3.2):
More detailed explanation of concept X with examples.

Recommendation: Keep Version B, update cross-references.
:::

```

### 3. Consolidation Actions

In markdown blocks or with `verbatim` code ticks:

```

<!-- CONSOLIDATION PLAN -->
::: {.consolidation-plan}
1. Keep primary version in Section 2.3
2. Remove duplicate from Section 4.1
3. Add cross-reference from Section 4.1 to Section 2.3
4. Merge unique insights from Section 4.1 into Section 2.3
:::

```

## Automated Detection Helpers

### Search Patterns

markdown

```

<!-- Common duplicate indicators -->
- "As mentioned earlier"
- "As discussed in"
- "Similar to"
- "Like we saw in"
- "Returning to"

```

### Duplicate Detection Checklist

In markdown blocks or with `verbatim` code ticks:

```

<!-- [] Check for repeated definitions -->
<!-- [] Identify similar examples -->
<!-- [] Find redundant explanations -->
<!-- [] Locate repeated figures/tables -->
<!-- [] Search for similar section headings -->

```

## Best Practices for Merging

### 1. Pre-Merge Preparation

- > Mark all content systematically
- > Create comparison documents
- > Track all locations of similar content
- > Document rationale for decisions

### 2. During Merge Process

- > Keep best version based on:
  - Completeness
  - Clarity
  - Placement in document flow
  - Citation quality
  - Figure/table quality

### 3. Post-Merge Cleanup

- > Update all cross-references
- > Remove duplicate citations
- > Consolidate figures/tables
- > Harmonize terminology
- > Verify logical flow

## Templates for Common Scenarios

### Duplicate Definition

markdown

```

::: {.duplicate-definition data-term="Bayesian Network"}
DUPLICATE DEFINITION: "Bayesian Network" is defined in:
- Section 2.1 (basic definition)
- Section 3.3 (technical definition) ← **KEEP THIS**
- Glossary (summary definition)

Action: Keep technical definition in 3.3, reference from 2.1
:::

```

### Redundant Example

markdown

```
::: {.redundant-example}
REDUNDANT EXAMPLE: Rain-Sprinkler-Lawn appears in:
1. Introduction (brief mention)
2. Chapter 2 (detailed walkthrough) ← **PRIMARY**
3. Chapter 4 (reference only)

Action: Keep detailed version, add cross-references from others
:::
```

## Overlapping Sections

markdown

```
::: {.section-overlap}
SECTION OVERLAP:
- Section 3.2 "Methodology Overview"
- Section 4.1 "Methods Used"

Content comparison:
- 70% overlap in general approach
- 3.2 has better technical detail
- 4.1 has better practical examples

Recommendation: Merge into 3.2, incorporate examples from 4.1
:::
```

## Visual Summary Blocks

### Merge Status Dashboard

markdown

```
::: {.merge-status-dashboard}
Chapter 2 Merge Status
- Total sections: 15
- Duplicates found: 4
- Redundant content: 7
- Unique content: 4
- Merge complete: 2/11
- Pending review: 9
:::
```

### Decision Log

markdown

```

::: {.merge-decision-log}
Merge Decisions - 2024-02-15
1. **Section 2.3 vs 4.1**: Kept 2.3, removed 4.1
2. **Definition of AI**: Consolidated in Glossary
3. **Example set A vs B**: Merged best of both into new set
4. **Figure 2.1 vs 3.2**: Kept 3.2 (higher quality)
:::

```

## Quality Assurance

### Pre-Publication Checklist

In markdown blocks or with `verbatim` code ticks:

```

<!-- [] All duplicate markers removed -->
<!-- [] All merge decisions documented -->
<!-- [] Cross-references updated -->
<!-- [] No broken links from removed content -->
<!-- [] Terminology harmonized -->
<!-- [] Flow tested after merging -->

```

### Final Verification

In markdown blocks or with `verbatim` code ticks:

```

<!-- FINAL CHECK: Content Merging -->
- [] No duplicate content remains untagged
- [] All redundancies resolved
- [] Best versions retained
- [] Smooth transitions between merged sections
- [] Complete citation consolidation
- [] Figure/table deduplication

```



# Master Thesis Checklist for Quarto Projects

## Content Creation Checklist

### Document Structure

- ☐ All chapters following consistent structure
- ☐ Proper heading hierarchy (##, ###, ####)
- ☐ Section labels added (`{#sec-label}`)

### Text Quality

- ☐ American spelling throughout (run spell check)
- ☐ Consistent terminology (maintain glossary, add entries)
- ☐ Active voice preferred
- ☐ Sentences clear and concise
- ☐ Paragraphs focused on single ideas
- ☐ Transitions between sections smooth
- ☐ No widows or orphans in paragraphs

### Formatting Elements

- ☐ Lists properly formatted and consistent
- ☐ Code blocks with appropriate syntax highlighting
- ☐ Blockquotes used for citations
- ☐ Callout boxes for important information
- ☐ Mathematical equations properly formatted
- ☐ Footnotes used wherever possible
- ☐ Page breaks inserted where needed

### Figures and Tables

- ☐ All figures have unique identifiers (`#fig-name`)
- ☐ Comprehensive alt text for accessibility
- ☐ Short captions for list of figures

- ☐ Full captions explaining content
- ☐ Consistent sizing and alignment
- ☐ All figures referenced in text
- ☐ Source attribution included
- ☐ File formats optimized (PNG/SVG for web, PDF for print)
- ☐ Tables have proper headers
- ☐ Table captions descriptive
- ☐ All tables referenced in text

### Citations and References

- ☐ All claims supported by citations
- ☐ Citation style consistent throughout
- ☐ Page numbers included where appropriate
- ☐ Bibliography entries complete
- ☐ No missing citations (check FIND tasks)
- ☐ No duplicate citations
- ☐ Citations verified (check VERIFY tasks)
- ☐ DOIs/URLs included and working

### Cross-References

- ☐ All sections labeled for referencing
- ☐ Figure references working (`@fig-name`)
- ☐ Table references working (`@tbl-name`)
- ☐ Section references working (`@sec-name`)
- ☐ No broken cross-references

## Revision Phase

### Content Review

- ☐ Argument flow logical and clear
- ☐ Evidence supports all claims
- ☐ Counterarguments addressed
- ☐ Conclusions follow from evidence
- ☐ No redundant content (check merge tags)
- ☐ All promises in introduction fulfilled

### Task Completion

- ☐ All TODO items addressed or documented
- ☐ All FIND items researched
- ☐ All VERIFY items confirmed
- ☐ All CREATE items completed



- ☐ Task status updated ([], [x], [-])
- ☐ Progress summaries updated

## Prime Directives

1. **Quarto supremacy** – exploit *every* reliable feature Quarto offers.
2. **Four-level heading discipline** – never skip a level.
3. **Redundancy tagged, not deleted** – see § Tagging System.
4. **Checklists rule every commit** – see § Rigorous Checklist.
5. **Footnotes galore** – default to footnotes for nuance, citations, side quests.
6. **Glossary, TOC, LOF, LOT, appendices, cross-refs** – keep fully synched; update on *every* save.
7. **One thought one line-break** – err on the side of whitespace also when formatting syntax.

## Quarto Syntax Cheat-Sheet    Best-Practice

Feature	Minimal Syntax	Best-Practice Guidance
<b>Headings</b> (h1–h4)	<b>#, ##, ###, ####</b>	Use all four levels; propose deeper sub-heads via <code>&lt;!-- SUGGEST-H5: ... --&gt;</code> .
<b>Paragraph breaks</b>	blank line	Generally wrap at 80 chars    git-diff clarity.
<b>Bold / <i>Italic</i> /</b> <i>Both** /</i> <b>Strike**</b>	<b><code>**b**, *i*, ***bi***,</code></b> <b><code>~~del~~</code></b>	Reserve bold for <i>semantic</i> emphasis, italics for <i>titles &amp; meta</i> .
<b>Lists</b>	<b><code>-, *, 1.</code></b>	Rarely nest > 3 levels; indent 2 spaces per level.
<b>Callouts</b>	<b><code>::: {.callout-note}</code></b>	Use <code>.tip</code> , <code>.warning</code> , <code>.important</code> , <code>.duplicate</code> (custom) for tagging; close with <code>:::</code> .
<b>Blockquotes</b>	<b><code>&gt;</code></b>	Ideal for verbatim interview excerpts; cite speaker.
<b>Code blocks</b>	<b><code>```r</code></b>	Always declare language; add caption: ““ <code>{python} fig.cap="..."</code> ”.
<b>Figures &amp; Tables</b>	<b><code>![] (fig.png){#fig-id}</code></b>	Always add <code>{#fig-id fig-cap="..."}</code> ; etc. cross-ref with <code>@fig-id</code> .

Feature	Minimal Syntax	Best-Practice Guidance
<b>Cross-refs</b>	<code>@sec-intro</code> , <code>@tbl-results</code>	Prefix: <code>sec-</code> , <code>fig-</code> , <code>tbl-</code> , <code>eq-</code> .
<b>Citations</b>	<code>[@smith2024]</code>	Maintain <code>.bib</code> via Zotero; nightly <code>quarto check</code> .
<b>Footnotes</b>	<code>[^1]</code>	Overuse for tangents, mini-proofs, data caveats.
<b>Glossary</b>	<code>term: Definition</code>	Append <code>glossary: true</code> in task; link in-text <code>{@term}</code> .
<b>Comments</b>	<code>&lt;!-- [ ] TODO: ... --&gt;</code>	Use for tasks; parse with <i>Todo Tree</i> VS Code plugin.

## Tagging / Highlighting System

Use the custom tagging and highlighting system for all materials

## Workflow Rules

### During writing

- *Every new idea*: decide *body*, *footnote*, or *appendix* and place instantly.
- Add glossary entries as soon as a term of art appears.
- Insert provisional graphics with stub `{#fig-TBD}` and create a TODO comment.

## Rigorous Checklist

- ☐ Use the full hierarchy of headings
- ☐ All figures/tables have IDs, captions, and are referenced in text.
- ☐ Glossary updated; new `{@term}` links render without warnings.
- ☐ Citation list reflects *every* `[@]` callout.
- ☐ Footnotes compile and are sorted numerically.
- ☐ Appendices contain overflow material only; each referenced at least once.
- ☐ `duplicate` callouts reviewed; none accidentally removed.
- ☐ “Outstanding graphics” & “Outstanding citations” subsections updated.

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

```
title: "Automating the Modeling of Transformative Artificial Intelligence Risks"
subtitle: "An Epistemic Framework for Leveraging Frontier AI Systems to Upscale Conditional

- name: Valentin Jakob Meyer orcid: 0009-0006-0889-5269 corresponding: true email: [Valentin
 - Graduate Author affiliations:
 - University of Bayreuth
 - MCMP - LMU Munich
- name: Dr. Timo Speith orcid: 0000-0002-6675-154X corresponding: false roles:
 - Supervisor affiliations:
 - University of Bayreuth keywords:
- AMTAIR
- AI Governance
- Bayesian Networks
- Transformative AI
- Risk Assessment
- Argument Extraction
- Existential Risk
- Coordination Crisis
- Epistemic Security
- Policy Evaluation abstract: | This thesis addresses coordination failures in AI safety by

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

The thesis contributes both theoretical foundations and practical implementation, validated

- A novel two-stage extraction pipeline transforms argument structures into Bayesian network
- Interactive visualizations make complex probabilistic relationships accessible to diverse
- Formal representation enables systematic comparison across different worldviews and assumpt
```

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

# Frontmatter: Preface

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

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

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

## Acknowledgments

I thank my supervisor Dr. Timo Speith for guidance throughout this project, the MTAIR team for pioneering the manual approach that inspired automation, and the AI safety community for creating the rich literature that made this work possible. Special recognition goes to technical advisors who provided invaluable feedback and Coleman Snell for his partnership and research collaboration with the AMTAIR project. Any errors or limitations remain my own responsibility.



# List of Figures





# List of Tables



# List of Abbreviations

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



# 1. Introduction: The Coordination Crisis in AI Governance

## Chapter Overview

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

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

## 1.1 Opening Scenario: The Policymaker’s Dilemma

Todd [69]

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

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

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

## 1.2 The Coordination Crisis in AI Governance

Maslej [47]

Samborska [61]

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

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

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

Tegmark [67]

European [28]

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

### 1.2.1 Safety Gaps from Misaligned Efforts

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

When different communities operate with incompatible frameworks, critical risks fall through

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

Paul [53]

### **1.2.2 Resource Misallocation**

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

### **1.2.3 Negative-Sum Dynamics**

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

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

Armstrong, Bostrom, and Shulman [3]

Samuel [62], Hunt [36]

## **1.3 Historical Parallels and Temporal Urgency**

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

Schelling [63]

Rehman [59]

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

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

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

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

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

## 1.4 Research Question and Scope

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

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

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

To break this down into its components:

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

The investigation encompasses both theoretical development and practical implementation, focusing specifically on existential risks from misaligned AI systems rather than broader AI ethics concerns. This narrowed scope enables deep technical development while addressing the highest-stakes coordination challenges.



## 1.5 The Multiplicative Benefits Framework

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

### 1.5.1 Automated Worldview Extraction

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

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

Automation using frontier language models addresses this scaling challenge. By developing systematic methods to extract causal structures and probability judgments from natural language, we can:

- > Process orders of magnitude more content
- > Incorporate diverse perspectives rapidly
- > Maintain models that evolve with the discourse
- > Reduce barriers to entry for contributing worldviews

### 1.5.2 Live Data Integration

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

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

This integration serves multiple purposes:

- > Grounding abstract models in empirical forecasts

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

Tetlock and Gardner [68]

### 1.5.3 Formal Policy Evaluation

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

This enables genuinely evidence-based policy development:

- > Comparing interventions across multiple worldviews
- > Identifying robust strategies that work across scenarios
- > Understanding which uncertainties most affect policy effectiveness
- > Prioritizing research to reduce decision-relevant uncertainty

Pearl [55] and Pearl [54]

### 1.5.4 The Synergy

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

The multiplicative benefits emerge from the interactions between components:

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

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

## 1.6 Thesis Structure and Roadmap

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

**Chapter 2: Context and Theoretical Foundations** establishes the intellectual groundwork, examining the epistemic challenges unique to AI governance, Bayesian networks as formal tools for uncertainty representation, argument mapping as a bridge from natural language to formal models, the MTAIR project’s achievements and limitations, and requirements for effective coordination infrastructure.

**Chapter 3: AMTAIR Design and Implementation** presents the technical system including overall architecture and design principles, the two-stage extraction pipeline (ArgDown → BayesDown), validation methodology and results, case studies from simple examples to complex AI risk models, and integration with prediction markets and policy evaluation.

**Chapter 4: Discussion - Implications and Limitations** critically examines technical limitations and failure modes, conceptual concerns about formalization, integration with existing governance frameworks, scaling challenges and opportunities, and broader implications for epistemic security.

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

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

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



## 2. Context and Theoretical Foundations

### Chapter Overview

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

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

### 2.1 AI Existential Risk: The Carlsmith Model

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

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

Carlsmith [15], Carlsmith [13] and Carlsmith [14]

#### 2.1.1 Six-Premise Decomposition

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

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

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

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

**Premise 1: APS Systems by 2070** ( $P \approx 0.65$ )<sup>4</sup> “By 2070, there will be AI systems with Advanced capability, Agentic planning, and Strategic awareness”—the conjunction of capabilities that could enable systematic pursuit of objectives in the world.

**Premise 2: Alignment Difficulty** ( $P = 0.40$ ) “It will be harder to build aligned APS systems than misaligned systems that are still attractive to deploy”—capturing the challenge that safety may conflict with capability or efficiency.

**Premise 3: Deployment Despite Misalignment** ( $P = 0.70$ ) “Conditional on 1 and 2, we will deploy misaligned APS systems”—reflecting competitive pressures and limited coordination.

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

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

**Premise 6: Existential Catastrophe** ( $P = 0.95$ ) “Conditional on 1-5, this disempowerment constitutes existential catastrophe”—connecting power loss to permanent curtailment of human potential.

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

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

### 2.1.2 Why Carlsmith Exemplifies Formalizable Arguments

Carlsmith’s model represents “low-hanging fruit” for automated formalization because it already exhibits explicit probabilistic reasoning with clear conditional dependencies. Success with this structured argument validates the approach for less explicit arguments throughout AI safety literature.

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

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

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

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

---

<sup>4</sup>The probability estimates vary between outlines; using more conservative estimates from 12.2

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

Christiano [17]

## 2.2 The Epistemic Challenge of Policy Evaluation

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

### 2.2.1 Unique Characteristics of AI Governance

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

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

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

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

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

### 2.2.2 Limitations of Traditional Approaches

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

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

Standard policy evaluation tools prove inadequate for these challenges:

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

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

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

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

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

Hallegatte et al. [32]

### 2.2.3 Toward New Epistemic Tools

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

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

#### Key Insight

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

McCaslin et al. [48]

Gruetzemacher [31]



## 2.3 Bayesian Networks as Knowledge Representation

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

### 2.3.1 Mathematical Foundations

A Bayesian network consists of:

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

The joint probability distribution factors according to the graph structure:

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

This factorization enables efficient inference and embodies causal assumptions explicitly.

Pearl [56]

### 2.3.2 The Rain-Sprinkler-Grass Example

The canonical example illustrates key concepts:<sup>5</sup>

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

Network Structure:

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

python

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

Conditional probability specification
P_wet_given_causes = {
```

---

<sup>5</sup>This example, while simple, demonstrates all essential features of Bayesian networks and serves as the foundation for understanding more complex applications

```

(True, True): 0.99, # Rain=T, Sprinkler=T
(True, False): 0.80, # Rain=T, Sprinkler=F
(False, True): 0.90, # Rain=F, Sprinkler=T
(False, False): 0.01 # Rain=F, Sprinkler=F
}

```

This simple network demonstrates:

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

### Rain-Sprinkler-Grass Network Rendering

```

from IPython.display import IFrame

IFrame(src="https://singularitysmith.github.io/AMTAIR_Prototype/bayesian_network.html", width=1000, height=500)

<IPython.lib.display.IFrame at 0x10b0521d0>

```

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

### 2.3.3 Advantages for AI Risk Modeling

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

These features address key requirements for AI governance:

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

## 2.4 Argument Mapping and Formal Representations

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

### 2.4.1 From Natural Language to Structure

Natural language arguments contain rich information expressed through:

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

Argument mapping extracts this structure, identifying:

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

Anderson [2]

Benn and Macintosh [7]

Khartabil et al. [39]

Khartabil [40]

Ngajie et al. [52]

Prokudin, Lisanyuk, and Baymurov [58]

Scheuer et al. [64]

**thomas1962**

Walton [71]

### 2.4.2 ArgDown: Structured Argument Notation

Voigt [70]

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

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

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

### 2.4.3 BayesDown: The Bridge to Bayesian Networks

BayesDown extends ArgDown with probabilistic metadata:

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

This representation:

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

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

## 2.5 The MTAIR Framework: Achievements and Limitations

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

### 2.5.1 MTAIR's Approach

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

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

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

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

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

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

Clarke et al. [18]

### 2.5.2 Key Achievements

MTAIR demonstrated several important possibilities:

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

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

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

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

### 2.5.3 Fundamental Limitations

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

However, MTAIR’s manual approach faces severe constraints:

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

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

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

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

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

### 2.5.4 The Automation Opportunity

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

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

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

## 2.6 Literature Review: Content and Technical Levels

### 2.6.1 AI Risk Models Evolution

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

#### Note

#### ## Key Developments

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

Yudkowsky [74]

Bostrom [9]

Amodei et al. [1]

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

### 2.6.2 Governance Proposals Taxonomy

AI governance proposals can be categorized along several dimensions:

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

Dafoe [21] and Dafoe [20]

Miotti et al. [50]

### 2.6.3 Bayesian Network Theory and Applications

The theoretical foundations of Bayesian networks rest on probability theory and graph theory.

Key concepts include:

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

Koller and Friedman [41]

### 2.6.4 Software Tools Landscape

The implementation of AMTAIR builds on established software libraries:

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

### 2.6.5 Formalization Approaches

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

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

Pollock [57]

### 2.6.6 Correlation Accounting Methods

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

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

Nelson [51]

## 2.7 Methodology

### 2.7.1 Research Design Overview

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

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

The research process follows four integrated phases:

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

### 2.7.2 Formalizing World Models from AI Safety Literature

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

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

The extraction approach combines several elements:

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

Large language models facilitate this process through specialized techniques:

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

### 2.7.3 From Natural Language to Computational Models

#### The Two-Stage Extraction Process

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



**Stage 1: Structural Extraction (ArgDown Generation)**

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 Enhancement)**

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)
```

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

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

Key mathematical properties essential for AI risk modeling:

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

The causal interpretation follows Pearl's framework:<sup>6</sup>

- > Edges represent direct causal influence
- > Intervention analysis through do-calculus

---

<sup>6</sup>Pearl's causal framework revolutionized how we think about causation in complex systems

- > Counterfactual reasoning for “what if” scenarios
- > Evidence integration through Bayesian updating

### 2.7.5 Quantification of Probabilistic Judgments

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

Standard linguistic mappings (with significant individual variation) include:

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

Expert elicitation methodologies:

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

Calibration challenges:

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

### 2.7.6 Inference Techniques for Complex Networks

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

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

Implementation considerations:

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

### 2.7.7 Integration with Prediction Markets and Forecasting Platforms

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

Live data sources for dynamic model updating include:

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

The data processing 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 challenges:

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

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



# 3. AMTAIR: Design and Implementation

## i Chapter Overview

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

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

## 3.1 System Architecture Overview

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

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

### 3.1.1 Five-Stage Pipeline Architecture

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

### 3.1.2 Design Principles

#### 💡 Core Design Philosophy

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

## 3.2 The Two-Stage Extraction Process

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

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

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

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

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

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

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

Example ArgDown extraction:

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

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

The second stage adds quantitative information to the structural skeleton:

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

Examples needed:

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

**Probability Extraction:**

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

**Coherence Enforcement:**

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

### 3.2.3 Why Two Stages?

This separation provides several benefits:

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

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

**Visibility:**

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

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

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

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

## 3.3 Implementation Technologies

### 3.3.1 Technology Stack

The system leverages established libraries while adding novel extraction capabilities:

Table 5: Technology stack components

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

### 3.3.2 Key Algorithms

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

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

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

**Visual Encoding Strategy:**

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

### 3.3.3 (Expected) Performance Characteristics

→

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

## 3.4 Case Study: Rain-Sprinkler-Grass

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

### 3.4.1 Processing Steps

The system processes this input through five steps:

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

### 3.4.2 Example Conversion Steps



## ArgDown Example

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

## Example of Questions for BayesDown extraction

BayesDown Format Preview:

# BayesDown Representation with Placeholder Probabilities

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

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

/\* What is the probability for Grass\_Wet=grass\_wet\_TRUE? \*/

/\* What is the probability for Grass\_Wet=grass\_wet\_TRUE if Rain=rain\_TRUE, Sprinkler=sprinkler\_TRUE? \*/

/\* What is the probability for Grass\_Wet=grass\_wet\_TRUE if Rain=rain\_TRUE, Sprinkler=sprinkler\_FALSE? \*/

/\* What is the probability for Grass\_Wet=grass\_wet\_TRUE if Rain=rain\_FALSE, Sprinkler=sprinkler\_TRUE? \*/

/\* What is the probability for Grass\_Wet=grass\_wet\_TRUE if Rain=rain\_FALSE, Sprinkler=sprinkler\_FALSE? \*/

/\* What is the probability for Grass\_Wet=grass\_wet\_FALSE? \*/

/\* What is the probability for Grass\_Wet=grass\_wet\_FALSE if Rain=rain\_TRUE, Sprinkler=sprinkler\_TRUE? \*/

/\* What is the probability for Grass\_Wet=grass\_wet\_FALSE if Rain=rain\_TRUE, Sprinkler=sprinkler\_FALSE? \*/

/\* What is the probability for Grass\_Wet=grass\_wet\_FALSE if Rain=rain\_FALSE, Sprinkler=sprinkler\_TRUE? \*/

/\* What is the probability for Grass\_Wet=grass\_wet\_FALSE if Rain=rain\_FALSE, Sprinkler=sprinkler\_FALSE? \*/

[Grass\_Wet]: Concentrated moisture on, between and around the blades of grass. {"instantiations": ["grass\_wet\_TRUE", "grass\_wet\_FALSE"]}

/\* What is the probability for Rain=rain\_TRUE? \*/

```

/* What is the probability for Rain=rain_FALSE? */
+ [Rain]: Tears of angles crying high up in the skies hitting the ground. {"instantiations": ["rain_TRUE", "rain_FALSE"], "priors": {"p(rain_TRUE)": "0.2", "p(rain_FALSE)": "0.8"}}, {"instantiations": ["rain_TRUE"], "priors": {"p(rain_TRUE)": "0.2", "p(rain_FALSE)": "0.8"}}, {"instantiations": ["rain_FALSE"], "priors": {"p(rain_TRUE)": "0.2", "p(rain_FALSE)": "0.8"}}]
/* What is the probability for Sprinkler=sprinkler_TRUE? */
/* What is the probability for Sprinkler=sprinkler_TRUE if Rain=rain_TRUE? */
/* What is the probability for Sprinkler=sprinkler_TRUE if Rain=rain_FALSE? */
/* What is the probability for Sprinkler=sprinkler_FALSE? */
/* What is the probability for Sprinkler=sprinkler_FALSE if Rain=rain_TRUE? */
/* What is the probability for Sprinkler=sprinkler_FALSE if Rain=rain_FALSE? */
+ [Sprinkler]: Activation of a centrifugal force based CO2 droplet distribution system. {"instantiations": ["sprinkler_TRUE", "sprinkler_FALSE"], "priors": {"p(sprinkler_TRUE)": "0.9", "p(sprinkler_FALSE)": "0.1"}}, {"instantiations": ["sprinkler_TRUE"], "priors": {"p(sprinkler_TRUE)": "0.9", "p(sprinkler_FALSE)": "0.1"}}, {"instantiations": ["sprinkler_FALSE"], "priors": {"p(sprinkler_TRUE)": "0.9", "p(sprinkler_FALSE)": "0.1"}}]
/* What is the probability for Rain=rain_TRUE? */
/* What is the probability for Rain=rain_FALSE? */
+ [Rain]

```

8

## Complete BayesDown Example

The source BayesDown syntax representating the fully specified network:

```

[Grass_Wet]: Concentrated moisture on grass. {"instantiations": ["grass_wet_TRUE", "grass_wet_FALSE"],
"priors": {"p(grass_wet_TRUE)": "0.322", "p(grass_wet_FALSE)": "0.678"},
"posteriors": {
 "p(grass_wet_TRUE|sprinkler_TRUE,rain_TRUE)": "0.99",
 "p(grass_wet_TRUE|sprinkler_TRUE,rain_FALSE)": "0.9",
 "p(grass_wet_TRUE|sprinkler_FALSE,rain_TRUE)": "0.8",
 "p(grass_wet_TRUE|sprinkler_FALSE,rain_FALSE)": "0.0"
}}
+ [Rain]: Water falling from sky. {"instantiations": ["rain_TRUE", "rain_FALSE"],
 "priors": {"p(rain_TRUE)": "0.2", "p(rain_FALSE)": "0.8"}}
+ [Sprinkler]: Artificial watering system. {"instantiations": ["sprinkler_TRUE", "sprinkler_FALSE"],
 "priors": {"p(sprinkler_TRUE)": "0.9", "p(sprinkler_FALSE)": "0.1"}}

```

```
"priors": {"p(sprinkler_TRUE)": "0.448", "p(sprinkler_FALSE)": "0.552"},
"posteriors": {
 "p(sprinkler_TRUE|rain_TRUE)": "0.01",
 "p(sprinkler_TRUE|rain_FALSE)": "0.4"
}}
+ [Rain]
```

Resulting Rain-Sprinkler-Grass DataFrame

#| column: page

Title	Description	line_number	model_id	instantiation	Prior	Device	instantiation	priors	posteriors	No_Parent	No_Children	instantiations
0	Grass_Wet concentrated moisture on, between and around t...	3	[3]	0	[0]	[Rain, []	[Grass_wet_TRUE, grass_wet_FALSE]	[grass_wet_TRUE, grass_wet_FALSE]	'p(grass_wet_FA...			[rain_TRUE, rain_FALSE], [sprinkler_TRUE, spr...
1	Rain Tears of angles crying high up in the skies hi...	4	[4, 6]	2	[1, 2]	[Grass_Wet, Sprinkler]	[Grass_Wet_TRUE, rain_FALSE]	{'p(rain_TRUE)': 0.2, 'p(rain_FALSE)': 0.8}		True	False	[]
2	Sprinkler Activation of a centrifugal force based CO2 dr...	5	[5]	1	[1]	[Rain]	[Grass_Wet, sprinkler_TRUE, sprinkler_FALSE]	{'p(sprinkler_TRUE rain_TRUE)': 0.44838, 'p(sprinkler_TRUE rain_FALSE)': 0.01, 'p(spr...				[rain_TRUE, rain_FALSE]]

### 3.4.3 Results

#### Rain-Sprinkler-Grass Network Rendering

```
from IPython.display import IFrame

IFrame(src="https://singularitysmith.github.io/AMTAIR_Prototype/bayesian_network.html", width=1000, height=500)

<IPython.lib.display.IFrame at 0x10e65f490>
```

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

#### 💡 Validation Success

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

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

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

### 3.5.1 Model Complexity

The Carlsmith model contains:

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

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

### 3.5.2 Extraction Results

The automated extraction successfully identifies:

#### Core Risk Pathway:

Existential\_Catastrophe

- ← Human\_Disempowerment
- ← Scale\_Of\_Power\_Seeking
- ← Misaligned\_Power\_Seeking
- ← [APS\_Systems, Difficulty\_Of\_Alignment, Deployment\_Decisions]

#### Supporting Structure:

- > Competitive dynamics influencing deployment

- > Technical factors affecting alignment difficulty
- > Corrective mechanisms and their limitations

**Probability Preservation:**

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

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

#### 3.5.4 Insights from Formalization

Formal representation reveals several insights:

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

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

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

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

## 3.6 Validation Methodology

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

### 3.6.1 Ground Truth Construction

Plan the process:

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

—>

### 3.6.2 Evaluation Metrics

—>

### 3.6.3 Results Summary

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

### 3.6.4 Error Analysis

Common failure modes to avoid:

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

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

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

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

Understanding these limitations guides both current usage and future improvements.

## 3.7 Policy Evaluation Capabilities

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

### 3.7.1 Intervention Representation

→

### 3.7.2 Example: Deployment Governance

Consider a policy requiring safety certification before deployment:

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

**Results:**

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

This hypothetical quantitative analysis enables comparison across interventions.

### 3.7.3 Robustness Analysis

#### Cross-Worldview Robustness

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

## 3.8 Interactive Visualization Design

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

### 3.8.1 Visual Encoding Strategy

The system uses multiple visual channels:

**Color:** Probability magnitude (green=high, red=low)

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

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

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

### 3.8.2 Progressive Disclosure

Information appears at appropriate levels:

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

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

### 3.8.3 User Interface Elements

## 3.9 Integration with Prediction Markets

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

### 3.9.1 Design for Integration

#### Integration Architecture

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

Design documentation needed:

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

### 3.9.2 Challenges and Opportunities

Key integration challenges:

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

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

## 3.10 Computational Performance Analysis

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

### 3.10.1 Exact vs. Approximate Inference

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

**Monte Carlo Methods:** Sample from probability distributions to estimate queries

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

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

The system automatically selects appropriate methods based on network properties.

### 3.10.2 Scaling Strategies

For very large networks:

Document strategies with benchmarks:

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

## 3.11 Results and Achievements

### 3.11.1 Extraction Quality Assessment

→

### 3.11.2 Computational Performance

### 3.11.3 Policy Impact Evaluation

→

## 3.12 Summary of Technical Contributions

AMTAIR successfully demonstrates:



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

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

These results demonstrate both the feasibility and value of automated model extraction for AI governance. However, several important considerations and limitations merit discussion. The next chapter critically examines these issues, addresses potential objections, and explores the broader implications of this approach for enhancing epistemic security in AI governance.



## 4. Discussion: Implications and Limitations

### Chapter Overview

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

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

### 4.1 Technical Limitations and Responses

#### 4.1.1 Objection 1: Extraction Quality Boundaries

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

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

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

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

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

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

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

#### 4.1.2 Objection 2: False Precision in Uncertainty

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

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

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

Technical requirements:

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

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

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

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

@Lindley [45]

#### 4.1.3 Objection 3: Correlation Complexity

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

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

this:

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

**Copula methods:** For known correlation structures, copula functions can model dependencies while preserving marginal distributions. This extends standard Bayesian networks significantly.<sup>7</sup>

Nelson [51]

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

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

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

## 4.2 Conceptual and Methodological Concerns

### 4.2.1 Objection 4: Democratic Exclusion

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

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

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

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

**Comparative advantage** comes from making implicit technical content explicit, not adding complexity. When experts debate AI risk, they already employ sophisticated probabilistic

---

<sup>7</sup>Copulas provide a mathematically elegant way to separate marginal behavior from dependence structure

reasoning—formalization reveals rather than creates this complexity. Making hidden assumptions visible arguably enhances rather than reduces democratic participation.

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

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

#### 4.2.2 Objection 5: Oversimplification of Complex Systems

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

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

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

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

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

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

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

Box [10]

### 4.3 Red-Teaming Results

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

### 4.3.1 Adversarial Extraction Attempts

→

### 4.3.2 Robustness Findings

Key vulnerabilities of LLMs (and human experts) identified:

Specific metrics need validation:

- Anchoring bias: measured effect size with confidence intervals
  - Authority sensitivity: controlled experiment design
  - Complexity degradation: performance curve analysis
  - Context loss: dependency distance metrics
1. **Anchoring bias:** System tends to over-weight first probability mentioned<sup>8</sup>
  2. **Authority sensitivity:** Extracted probabilities influenced by cited expert prominence
  3. **Complexity degradation:** Performance drops sharply beyond 50 nodes
  4. **Context loss:** Long-range dependencies in text sometimes missed

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

### 4.3.3 Implications for Deployment

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

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

## 4.4 Enhancing Epistemic Security

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

### 4.4.1 Making Models Inspectable

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

> Likely means what probability?

---

<sup>8</sup>This reflects how LLMs inherit human cognitive biases from training data

- > Through what causal pathways?
- > Under what assumptions?
- > With what evidence?

This explicitation serves multiple functions:

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

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

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

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

#### 4.4.2 Revealing Convergence and Divergence

Implement comparison of 3+ models:

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

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

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

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

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

—>

#### 4.4.3 Improving Collective Reasoning

AMTAIR enhances group epistemics through:

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

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

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



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

## 4.5 Scaling Challenges and Opportunities

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

### 4.5.1 Technical Scaling

**Computational complexity** grows with network size, but several approaches help: - Hierarchical decomposition for very large models - Caching and approximation for common queries - Distributed processing for extraction tasks - Incremental updating rather than full recomputation

**Data quality** varies dramatically across sources: - Academic papers provide structured arguments - Blog posts offer rich ideas with less formal structure - Policy documents mix normative and empirical claims - Social media presents extreme extraction challenges

**Integration complexity** increases with ecosystem growth: - Multiple LLM providers with different capabilities - Diverse visualization needs across users - Various export formats for downstream tools - Version control for evolving models

### 4.5.2 Social and Institutional Scaling

**Adoption barriers** include: - Learning curve for formal methods - Institutional inertia in established processes - Concerns about replacing human judgment - Resource requirements for implementation

**Trust building** requires: - Transparent methodology documentation - Published validation studies - High-profile successful applications - Community ownership and development

**Sustainability** depends on: - Open source development model - Diverse funding sources - Academic and industry partnerships - Clear value demonstration

### 4.5.3 Opportunities for Impact

Despite challenges, several factors favor adoption:

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

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

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

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

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

## 4.6 Integration with Governance Frameworks

AMTAIR complements and integrates rather than replaces existing governance approaches.

### 4.6.1 Standards Development

Technical standards bodies could use AMTAIR to: - Model how proposed standards affect risk pathways - Compare different standard options systematically - Identify unintended consequences through pathway analysis - Build consensus through explicit model negotiation

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

### 4.6.2 Regulatory Design

Regulators could apply the framework to: - Assess regulatory impact across different scenarios - Identify enforcement challenges through explicit modeling - Compare international approaches systematically - Design adaptive regulations responsive to evidence

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

Cuomo, Mallin, and Zattoni [19], Demirag, Sudarsanam, and WRIGHT [23], De Villiers and Dimes [22], Di Vito and Trottier [24], Kaur [38], List and Pettit [46] and Solomon [65]

### 4.6.3 International Coordination

Multilateral bodies could leverage shared models for: - Establishing common risk assessments - Negotiating agreements with explicit assumptions - Monitoring compliance through parameter tracking - Adapting agreements as evidence emerges

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

### 4.6.4 Organizational Decision-Making

Individual organizations could use AMTAIR for: - Internal risk assessment and planning - Board-level communication about AI strategies - Research prioritization based on model sensitivity - Safety case development with explicit assumptions

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

## 4.7 Future Research Directions

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

### 4.7.1 Technical Enhancements

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

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

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

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

### 4.7.2 Methodological Extensions

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

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

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

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

Babakov et al. [5], Ban et al. [6], Bethard [8], Chen et al. [16], Duhem [27], Heinze-Deml, Maathuis, and Meinshausen [33], Meyer [49], Squires and Uhler [66], Squires and Uhler [66], Yang, Han, and Poon [73]

### 4.7.3 Application Domains

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

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

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

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

### 4.7.4 Ecosystem Development

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

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

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

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

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

## 4.8 Known Unknowns and Deep Uncertainties

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

### 4.8.1 Categories of Deep Uncertainty

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

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

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

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

### 4.8.2 Adaptation Strategies for Deep Uncertainty

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

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

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

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

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

### 4.8.3 Robust Decision-Making Principles

Given deep uncertainty, certain decision principles become paramount:

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

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

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

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

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

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



# 5. Conclusion: Toward Coordinated AI Governance

## Chapter Overview

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

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

## 5.1 Summary of Key Contributions

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

### 5.1.1 Theoretical Contributions

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

**The Multiplicative Benefits Framework:** The combination of automated extraction, prediction market integration, and formal policy evaluation creates value exceeding the sum of parts. Automation enables scale, markets provide empirical grounding, and policy analysis delivers actionable insights. Together, they address different facets of the coordination challenge while reinforcing each other's strengths.

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

### 5.1.2 Methodological Innovations

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

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

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

### 5.1.3 Technical Achievements

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

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

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

### 5.1.4 Empirical Findings

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

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

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



## 5.2 Limitations and Honest Assessment

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

### 5.2.1 Technical Constraints

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

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

**Computational Scaling:** Very large networks (»500 nodes) require approximations that may affect accuracy. As models grow to represent richer phenomena, computational constraints increasingly bind.

### 5.2.2 Conceptual Limitations

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

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

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

### 5.2.3 Practical Constraints

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

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

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

## 5.3 Implications for AI Governance

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

### 5.3.1 Near-Term Applications

**Research Coordination:** Research organizations can use formal models to: - Map the landscape of current arguments and identify gaps - Prioritize investigations targeting high-sensitivity parameters - Build cumulative knowledge through explicit model updating - Facilitate collaboration through shared representations

**Policy Development:** Governance bodies can apply the framework to: - Evaluate proposals across multiple expert worldviews - Identify robust interventions effective under uncertainty - Make assumptions explicit for democratic scrutiny - Track how evidence changes optimal policies over time

**Stakeholder Communication:** The visualization and analysis tools enable: - Clearer communication between technical and policy communities - Public engagement with complex risk assessments - Board-level strategic discussions grounded in formal analysis - International negotiations with explicit shared models

### 5.3.2 Medium-Term Transformation

As adoption spreads, we might see:

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

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

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

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

### 5.3.3 Long-Term Vision

Successfully scaling this approach could fundamentally alter AI governance:

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

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

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

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

## 5.4 Recommendations for Stakeholders

Different communities can take concrete steps to realize these benefits:

### 5.4.1 For Researchers

1. **Experiment with formalization:** Try extracting your own arguments into ArgDown/BayesDown format to discover implicit assumptions
2. **Contribute to validation:** Provide expert annotations for building benchmark datasets and improving extraction quality
3. **Develop extensions:** Build on the open-source foundation to add capabilities for your specific domain needs
4. **Publish formally:** Include formal model representations alongside traditional papers to enable cumulative building

#### Quick Start Guide

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

### 5.4.2 For Policymakers

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

### 5.4.3 For Technologists

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

4. **Scale infrastructure:** Address computational challenges for larger models and broader deployment

#### 5.4.4 For Funders

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

## 5.5 Future Research Agenda

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

### 5.5.1 Technical Priorities

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

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

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

### 5.5.2 Methodological Development

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

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

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

### 5.5.3 Application Expansion

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

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

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

## 5.6 Closing Reflections

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

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

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

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

### ! The Stakes

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

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

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

# References

- [1] Dario Amodei et al. *Concrete Problems in AI Safety*. July 25, 2016. DOI: 10.48550/arXiv.1606.06565. arXiv: 1606.06565 [cs]. URL: <http://arxiv.org/abs/1606.06565> (visited on 05/25/2025). Pre-published.
- [2] Terence J. Anderson. “Visualization Tools and Argument Schemes: A Question of Standpoint”. In: *Law, Prob. & Risk* 6 (2007), p. 97. URL: [https://heinonline.org/hol-cgi-bin/get\\_pdf.cgi?handle=hein.journals/lawprisk6&section=9](https://heinonline.org/hol-cgi-bin/get_pdf.cgi?handle=hein.journals/lawprisk6&section=9) (visited on 05/25/2025).
- [3] Stuart Armstrong, Nick Bostrom, and Carl Shulman. “Racing to the Precipice: A Model of Artificial Intelligence Development”. In: *AI & SOCIETY* 31.2 (May 1, 2016), pp. 201–206. ISSN: 1435-5655. DOI: 10.1007/s00146-015-0590-y. URL: <https://doi.org/10.1007/s00146-015-0590-y> (visited on 05/26/2025).
- [4] Amanda Askell et al. *A General Language Assistant as a Laboratory for Alignment*. Dec. 9, 2021. DOI: 10.48550/arXiv.2112.00861. arXiv: 2112.00861 [cs]. URL: <http://arxiv.org/abs/2112.00861> (visited on 05/25/2025). Pre-published.
- [5] Nikolay Babakov et al. “Reusability of Bayesian Networks Case Studies: A Survey”. In: *Applied Intelligence* 55.6 (Feb. 7, 2025), p. 417. ISSN: 1573-7497. DOI: 10.1007/s10489-025-06289-5. URL: <https://doi.org/10.1007/s10489-025-06289-5> (visited on 05/15/2025).
- [6] Taiyu Ban et al. *Causal Structure Learning Supervised by Large Language Model*. Nov. 20, 2023. DOI: 10.48550/arXiv.2311.11689. arXiv: 2311.11689 [cs]. URL: <http://arxiv.org/abs/2311.11689> (visited on 05/26/2025). Pre-published.
- [7] Neil Benn and Ann Macintosh. “Argument Visualization for eParticipation: Towards a Research Agenda and Prototype Tool”. In: *Electronic Participation*. Ed. by Efthimios Tambouris, Ann Macintosh, and Hans De Bruijn. Red. by David Hutchison et al. Vol. 6847. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 60–73. ISBN: 978-3-642-23332-6 978-3-642-23333-3. DOI: 10.1007/978-3-642-23333-3\_6. URL: [http://link.springer.com/10.1007/978-3-642-23333-3\\_6](http://link.springer.com/10.1007/978-3-642-23333-3_6) (visited on 05/25/2025).
- [8] Steven John Bethard. “Finding Event, Temporal and Causal Structure in Text: A Machine Learning Approach”. PhD thesis. University of Colorado at Boulder, 2007. URL: <https://search.proquest.com/openview/405fe32503123d9b5f4836dc3be4c011/1?pq-origsite=gscholar&cbl=18750> (visited on 05/26/2025).
- [9] Nick Bostrom. *Superintelligence: Paths, Strategies, Dangers*. Oxford: Oxford University Press, 2014. ISBN: 978-0-19-967811-2. URL: <https://scholar.dominican.edu/cynthia-stokes-brown-books-big-history/47>.

- [10] George E. P. Box. “Science and Statistics”. In: *Journal of the American Statistical Association* 71.356 (Dec. 1976), pp. 791–799. ISSN: 0162-1459, 1537-274X. DOI: 10.1080/01621459.1976.10480949. URL: <http://www.tandfonline.com/doi/abs/10.1080/01621459.1976.10480949> (visited on 05/26/2025).
- [11] Miles Brundage et al. *The Malicious Use of Artificial Intelligence: Forecasting, Prevention, and Mitigation*. Version 1. 2018. DOI: 10.48550/ARXIV.1802.07228. URL: <https://arxiv.org/abs/1802.07228> (visited on 11/13/2024). Pre-published.
- [12] Miles Brundage et al. “The Malicious Use of Artificial Intelligence: Forecasting, Prevention, and Mitigation”. 2018. arXiv: 1802.07228.
- [13] Joseph Carlsmith. *Is Power-Seeking AI an Existential Risk?* 2021. DOI: 10.48550/arXiv.2206.13353. URL: <https://arxiv.org/abs/2206.13353>. Pre-published.
- [14] Joseph Carlsmith. “Is Power-Seeking AI an Existential Risk?” 2022. arXiv: 2206.13353.
- [15] Joseph Carlsmith. *Is Power-Seeking AI an Existential Risk?* Aug. 13, 2024. DOI: 10.48550/arXiv.2206.13353. arXiv: 2206.13353 [cs]. URL: <http://arxiv.org/abs/2206.13353> (visited on 05/25/2025). Pre-published.
- [16] Lu Chen et al. “Inducing Causal Structure for Abstractive Text Summarization”. In: *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*. CIKM ’23: The 32nd ACM International Conference on Information and Knowledge Management. Birmingham United Kingdom: ACM, Oct. 21, 2023, pp. 213–223. ISBN: 979-8-4007-0124-5. DOI: 10.1145/3583780.3614934. URL: <https://dl.acm.org/doi/10.1145/3583780.3614934> (visited on 05/26/2025).
- [17] Paul F. Christiano. “What Failure Looks Like”. In: (Mar. 17, 2019). URL: <https://www.alignmentforum.org/posts/HBxe6wdjxK239zajf/what-failure-looks-like> (visited on 05/25/2025).
- [18] Sam Clarke et al. *Modeling Transformative AI Risks (MTAIR) Project – Summary Report*. Version 1. 2022. DOI: 10.48550/ARXIV.2206.09360. URL: <https://arxiv.org/abs/2206.09360> (visited on 11/13/2024). Pre-published.
- [19] Francesca Cuomo, Christine Mallin, and Alessandro Zattoni. “Corporate Governance Codes: A Review and Research Agenda”. In: *Corporate governance: an international review* 24.3 (2016), pp. 222–241. URL: <https://ueaeprints.uea.ac.uk/id/eprint/57664/> (visited on 05/26/2025).
- [20] Allan Dafoe. “AI Governance: A Research Agenda”. In: *Governance of AI Program, Future of Humanity Institute, University of Oxford: Oxford, UK* 1442 (2018), p. 1443. URL: <http://www.fhi.ox.ac.uk/wp-content/uploads/GovAI-Agenda.pdf> (visited on 05/25/2025).
- [21] Allan Dafoe. *AI Governance: A Research Agenda*. 2021. URL: <https://www.fhi.ox.ac.uk/wp-content/uploads/GovAI-Agenda.pdf> (visited on 05/25/2025). Pre-published.
- [22] Charl De Villiers and Ruth Dimes. “Determinants, Mechanisms and Consequences of Corporate Governance Reporting: A Research Framework”. In: *Journal of Management and Governance* 25.1 (Mar. 2021), pp. 7–26. ISSN: 1385-3457, 1572-963X. DOI: 10.1007/s10997-020-09530-0. URL: <https://link.springer.com/10.1007/s10997-020-09530-0> (visited on 05/26/2025).



- [23] Istemi Demirag, Sudi Sudarsanam, and MIKE WRIGHT. “Corporate Governance: Overview and Research Agenda”. In: *The British Accounting Review* 32.4 (2000), pp. 341–354. URL: <https://www.academia.edu/download/49469624/bare.2000.014620161009-3955-1dt4aq5.pdf> (visited on 05/26/2025).
- [24] Jackie Di Vito and Kim Trottier. “A Literature Review on Corporate Governance Mechanisms: Past, Present, and Future\*”. In: *Accounting Perspectives* 21.2 (June 2022), pp. 207–235. ISSN: 1911-382X, 1911-3838. DOI: 10.1111/1911-3838.12279. URL: <https://onlinelibrary.wiley.com/doi/10.1111/1911-3838.12279> (visited on 05/26/2025).
- [25] K. Eric Drexler. “Reframing Superintelligence: Comprehensive AI Services as General Intelligence”. In: (2019).
- [26] KE Drexler. *Reframing Superintelligence: Comprehensive AI Services as General Intelligence*. Technical Report. Future of Humanity Institute, 2019. URL: [https://owainevans.github.io/pdfs/Reframing\\_Superintelligence\\_FHI-TR-2019.pdf](https://owainevans.github.io/pdfs/Reframing_Superintelligence_FHI-TR-2019.pdf).
- [27] Pierre Maurice Marie Duhem. *The Aim and Structure of Physical Theory*. 1. Princeton University Press, 1954. 85–87.
- [28] Union European. *The Act Texts / EU Artificial Intelligence Act*. 2024. URL: <https://artificialintelligenceact.eu/the-act/> (visited on 05/25/2025).
- [29] Irving John Good. “Speculations Concerning the First Ultraintelligent Machine”. In: *Advances in Computers* (1966), p. 31. DOI: 10.1016/S0065-2458(08)60418-0. URL: <https://www.sciencedirect.com/science/article/abs/pii/S0065245808604180>.
- [30] Jakub Growiec. “Existential Risk from Transformative AI: An Economic Perspective”. In: *Technological and Economic Development of Economy* 30.6 (2024), pp. 1682–1708.
- [31] Ross Gruetzmacher. “Bayesian Networks vs. Conditional Trees for Creating Questions for Forecasting Tournaments”. In: (2022).
- [32] Stéphane Hallegatte et al. “Investment Decision-Making under Deep Uncertainty-Application to Climate Change”. In: *Policy research working paper* 6193 (2012). URL: <https://enpc.hal.science/hal-00802049/document> (visited on 05/25/2025).
- [33] Christina Heinze-Deml, Marloes H. Maathuis, and Nicolai Meinshausen. “Causal Structure Learning”. In: *Annual Review of Statistics and Its Application* 5.1 (Mar. 7, 2018), pp. 371–391. ISSN: 2326-8298, 2326-831X. DOI: 10.1146/annurev-statistics-031017-100630. URL: <https://www.annualreviews.org/doi/10.1146/annurev-statistics-031017-100630> (visited on 05/26/2025).
- [34] Dan Hendrycks et al. *Unsolved Problems in ML Safety*. Version 5. 2021. DOI: 10.48550/ARXIV.2109.13916. URL: <https://arxiv.org/abs/2109.13916> (visited on 11/13/2024). Pre-published.
- [35] Dan Hendrycks et al. “Unsolved Problems in ML Safety”. 2021. arXiv: 2109.13916.
- [36] Tam Hunt. *The Insane “Logic” of the AI Arms Race*. Medium. Mar. 3, 2025. URL: <https://tamhunt.medium.com/the-insane-logic-of-the-ai-arms-race-45a5f79f4c0e> (visited on 05/26/2025).
- [37] Edwin T Jaynes. *Probability Theory: The Logic of Science*. Cambridge university press, 2003.

- [38] Kawaljit Kaur. “Corporate Governance and Legal Accountability: A Critical Review of Global Practices”. In: *Journal of Law* 2.6 (2024), pp. 1–7. URL: <https://joi.shodhsagar.org/index.php/SSJOI/article/view/16> (visited on 05/26/2025).
- [39] D. Khartabil et al. “Design and Evaluation of Visualization Techniques to Facilitate Argument Exploration”. In: *Computer Graphics Forum* 40.6 (Sept. 2021), pp. 447–465. ISSN: 0167-7055, 1467-8659. DOI: 10.1111/cgf.14389. URL: <https://onlinelibrary.wiley.com/doi/10.1111/cgf.14389> (visited on 05/25/2025).
- [40] Dana Khartabil. “Visualisation Techniques to Facilitate Argument Exploration”. PhD thesis. 2020. URL: <https://napier-repository.worktribe.com/output/2694675> (visited on 05/25/2025).
- [41] Daphne Koller and Nir Friedman. *Probabilistic Graphical Models: Principles and Techniques*. MIT press, 2009. URL: [https://books.google.ca/books?hl=en&lr=&id=7dzpHCHzNQ4C&oi=fnd&pg=PR9&dq=Koller,+D.,+%26+Friedman,+N.+\(2009\).+Probabilistic+Graphical+Models&ots=py2HAh0VAL&sig=gpaID3x6-TY8x5SOopuXpZDXfzs](https://books.google.ca/books?hl=en&lr=&id=7dzpHCHzNQ4C&oi=fnd&pg=PR9&dq=Koller,+D.,+%26+Friedman,+N.+(2009).+Probabilistic+Graphical+Models&ots=py2HAh0VAL&sig=gpaID3x6-TY8x5SOopuXpZDXfzs) (visited on 05/25/2025).
- [42] Ram Shankar Siva Kumar et al. *Failure Modes in Machine Learning Systems*. Version 1. 2019. DOI: 10.48550/ARXIV.1911.11034. URL: <https://arxiv.org/abs/1911.11034> (visited on 11/13/2024). Pre-published.
- [43] Ram Shankar Siva Kumar et al. “Failure Modes in Machine Learning Systems”. 2019. arXiv: 1911.11034.
- [44] Robert J Lempert, Steven W Popper, and Steven C Bankes. *Shaping the next One Hundred Years: New Methods for Quantitative, Long-Term Policy Analysis*. RAND Corporation, 2003.
- [45] Dennis V. Lindley. *Understanding Uncertainty*. John Wiley & Sons, 2013. URL: [https://books.google.ca/books?hl=en&lr=&id=Tfk8AgAAQBAJ&oi=fnd&pg=PR11&dq=Lindley,+D.+\(2006\).+Understanding+Uncertainty&ots=55HS6lTOVP&sig=0gKCDvRu5rUKhuPyJqhOzW23upU](https://books.google.ca/books?hl=en&lr=&id=Tfk8AgAAQBAJ&oi=fnd&pg=PR11&dq=Lindley,+D.+(2006).+Understanding+Uncertainty&ots=55HS6lTOVP&sig=0gKCDvRu5rUKhuPyJqhOzW23upU) (visited on 05/26/2025).
- [46] Christian List and Philip Pettit. *Group Agency: The Possibility, Design, and Status of Corporate Agents*. Oxford University Press, 2011.
- [47] Nestor Maslej. “Artificial Intelligence Index Report 2025”. In: *Artificial Intelligence* (2025).
- [48] Tegan McCaslin et al. “Conditional Trees: A Method for Generating Informative Questions about Complex Topics”. In: *Forecasting Research Institute* (2024). URL: <https://static1.squarespace.com/static/635693acf15a3e2a14a56a4a/t/66ba37a144f1d6095de467df/1723479995772/AIConditionalTrees.pdf>.
- [49] Valentin Jakob Meyer. “A Structure of Knowledge & the Process of Science”. In: *Philosophy of the Social Sciences* First Course Paper (University Bayreuth 2022). DOI: <https://www.vjmeyer.com/papers/essays>.
- [50] Andrea Miotti et al. *A Narrow Path*. 2024. URL: <https://www.narrowpath.co/> (visited on 05/19/2025).
- [51] Roger B. Nelson. *An Introduction to Copulas*. Springer Series in Statistics. New York, NY: Springer New York, 2006. ISBN: 978-0-387-28659-4. DOI: 10.1007/0-387-28678-0. URL: <http://link.springer.com/10.1007/0-387-28678-0> (visited on 05/25/2025).

- [52] Berty Nsolly Ngajie et al. “Investigating the Effects of a Systematic and Model-Based Design of Computer-Supported Argument Visualization on Critical Thinking”. In: *Thinking Skills and Creativity* 38 (2020), p. 100742. URL: <https://www.sciencedirect.com/science/article/pii/S1871187120302169> (visited on 05/25/2025).
- [53] Paul. *The elephAInt – Are We All like the Six Blind Men When It Comes to AI? / PRISMAGuard LLC*. 2023. URL: <https://www.prismaguard.com/the-elephaint-are-we-all-like-the-six-blind-men-when-it-comes-to-ai/> (visited on 05/25/2025).
- [54] Judea Pearl. *Causality: Models, Reasoning and Inference*. 2nd ed. Cambridge University Press, 2009.
- [55] Judea Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge, U.K. ; New York: Cambridge University Press, 2000. 384 pp. ISBN: 978-0-521-89560-6 978-0-521-77362-1.
- [56] Judea Pearl. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Elsevier, 2014. URL: [https://books.google.ca/books?hl=en&lr=&id=mn2jBQAAQBAJ&oi=fnd&pg=PP1&dq=Pearl,+J.+\(1988\).+Probabilistic+Reasoning+in+Intelligent+Systems&ots=4tEX2A4Ha8&sig=lgUs\\_RCoeXEEuGwM5xMEoyJy4HI](https://books.google.ca/books?hl=en&lr=&id=mn2jBQAAQBAJ&oi=fnd&pg=PP1&dq=Pearl,+J.+(1988).+Probabilistic+Reasoning+in+Intelligent+Systems&ots=4tEX2A4Ha8&sig=lgUs_RCoeXEEuGwM5xMEoyJy4HI) (visited on 05/25/2025).
- [57] John L. Pollock. *Cognitive Carpentry: A Blueprint for How to Build a Person*. Mit Press, 1995. URL: [https://books.google.ca/books?hl=en&lr=&id=JAfHrHTqswAC&oi=fnd&pg=PA1&dq=Pollock,+J.+\(1995\).+Cognitive+Carpentry&ots=rq-qSCBcxV&sig=aAfHGSGUosxl\\_1-JuxIEA7C2QO4](https://books.google.ca/books?hl=en&lr=&id=JAfHrHTqswAC&oi=fnd&pg=PA1&dq=Pollock,+J.+(1995).+Cognitive+Carpentry&ots=rq-qSCBcxV&sig=aAfHGSGUosxl_1-JuxIEA7C2QO4) (visited on 05/25/2025).
- [58] D. E. Prokudin, E. N. Lisanyuk, and I. R. Baymuratov. “Visualization Functions in Argumentation Representation Software”. In: *Scientific Visualization* 16.3 (2024). URL: <https://sv-journal.org/2024-3/11/en.pdf> (visited on 05/25/2025).
- [59] Iskander Rehman. *The Battle for Brilliant Minds: From the Nuclear Age to AI*. War on the Rocks. Jan. 13, 2025. URL: <https://warontherocks.com/2025/01/the-battle-for-brilliant-minds-from-the-nuclear-age-to-ai/> (visited on 05/25/2025).
- [60] Stuart Russell et al. “Research Priorities for Robust and Beneficial Artificial Intelligence: An Open Letter”. In: *AI Magazine* 36.4 (2015), pp. 3–4. DOI: 10.1609/aimag.v36i4.2621. URL: <https://ojs.aaai.org/aimagazine/index.php/aimagazine/article/view/2621>.
- [61] Veronika Samborska. “Scaling up: How Increasing Inputs Has Made Artificial Intelligence More Capable”. In: *Our World in Data* (Jan. 20, 2025). URL: <https://ourworldindata.org/scaling-up-ai> (visited on 05/25/2025).
- [62] Sigal Samuel. *AI Is a “Tragedy of the Commons.” We’ve Got Solutions for That*. Vox. July 7, 2023. URL: <https://www.vox.com/future-perfect/2023/7/7/23787011/ai-arms-race-tragedy-commons-risk-safety> (visited on 05/26/2025).
- [63] Thomas C. Schelling. “1960. The Strategy of Conflict”. In: *Cambridge, Mass* (1960).
- [64] Oliver Scheuer et al. “Computer-Supported Argumentation: A Review of the State of the Art”. In: *International Journal of Computer-Supported Collaborative Learning* 5.1 (Mar. 2010), pp. 43–102. ISSN: 1556-1607, 1556-1615. DOI: 10.1007/s11412-009-9080-x. URL: <http://link.springer.com/10.1007/s11412-009-9080-x> (visited on 05/25/2025).
- [65] Jill Solomon. *Corporate Governance and Accountability*. John Wiley & Sons, 2020. URL: <https://books.google.ca/books?hl=en&lr=&id=JAX9DwAAQBAJ&oi=fnd&pg=PR1&d>

- q=review+of+the+effects+of+liability+frameworks+on+corporate+governance+&ots=ny23\_vd-U0&sig=3LuNNhvSWXriEeg-ipAdDIQGAgo (visited on 05/26/2025).
- [66] Chandler Squires and Caroline Uhler. “Causal Structure Learning: A Combinatorial Perspective”. In: *Foundations of Computational Mathematics* 23.5 (Oct. 2023), pp. 1781–1815. ISSN: 1615-3375, 1615-3383. DOI: 10.1007/s10208-022-09581-9. URL: <https://link.springer.com/10.1007/s10208-022-09581-9> (visited on 05/26/2025).
- [67] Max Tegmark. *Asilomar AI Principles*. Future of Life Institute. 2024. URL: <https://futureoflife.org/open-letter/ai-principles/> (visited on 05/25/2025).
- [68] Philip E. Tetlock and Dan Gardner. *Superforecasting: The Art and Science of Prediction*. First paperback edition. New York: Broadway Books, 2015. 340 pp. ISBN: 978-0-8041-3671-6.
- [69] Benjamin Todd. *It Looks like There Are Some Good Funding Opportunities in AI Safety Right Now*. Benjamin Todd. Dec. 21, 2024. URL: <https://benjamintodd.substack.com/p/it-looks-like-there-are-some-good-funding> (visited on 05/25/2025).
- [70] Christian Voigt. *Christianvoigt/Argdown*. May 23, 2025. URL: <https://github.com/christianvoigt/argdown> (visited on 05/25/2025).
- [71] Douglas Walton. “Argument Visualization Tools for Corroborative Evidence”. In: *Proc. of the 2nd International Conference on Evidence Law and Forensic Science*. 2009, pp. 32–49. URL: <https://www.academia.edu/download/37718171/09ArguVis.pdf> (visited on 05/25/2025).
- [72] Nick Wilson et al. “The Need for Long-Term Thinking—Especially for Preventing Catastrophic Risks”. In: *Public Health Expert Briefing* (2023).
- [73] Jie Yang, Soyeon Caren Han, and Josiah Poon. “A Survey on Extraction of Causal Relations from Natural Language Text”. In: *Knowledge and Information Systems* 64.5 (May 2022), pp. 1161–1186. ISSN: 0219-1377, 0219-3116. DOI: 10.1007/s10115-022-01665-w. URL: <https://link.springer.com/10.1007/s10115-022-01665-w> (visited on 05/26/2025).
- [74] Eliezer Yudkowsky. “Artificial Intelligence as a Positive and Negative Factor in Global Risk”. In: *Global Catastrophic Risks*. Oxford University Press, July 3, 2008. ISBN: 978-0-19-857050-9 978-0-19-191810-0. DOI: 10.1093/oso/9780198570509.003.0021. URL: <https://academic.oup.com/book/40615/chapter/348239228> (visited on 11/15/2024).

{ embed I.Appendices.qmd }

# References (.md)

## Error Watch

### Catch ALL Potential Hallucinations

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

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

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

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

## Figure Inventory and Tracking

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

```
<!-- FIGURE INVENTORY -->
```

```
<!-- Last updated: 2024-02-15 -->
```

```
Implemented Figures
```

```
Section to keep track of all Figures
```

```

`<!-- [] ALWAYS include the "inclusions" of all figures/graphics below -->`
`<!-- [] ALWAYS keep the #fig-KEYS up-to-date -->`

```markdown
{{
[![Example Caption/Title 4] (/images/cover.png){
  #fig-Unique_identifier_for_crossreferencing
  fig-scrap="Short caption 4 list of figures as seen in LoF"
  fig-alt="Detailed alt text that describes the image content, type, purpose, and meaning.
    [CHART TYPE]: [Short description].
      DATA: [What data is shown, x/y axes].
      PURPOSE: [Why it's included, what to look for].
      DETAILS: [Longer description of patterns, anomalies, or key insights].
      SOURCE: Data from [source name/year and url/link]
    "
  fig-align="left"
  width="30%"
  }](https://github.com/VJMeyer/submission)
}}

```

Chapter 1

- ☒ {#fig-overview}: System overview diagram
 - File: images/system-overview.png
 - Source: Created by author using Draw.io

```

{{
[![Example Caption/Title 4] (/images/cover.png){
  #fig-Unique_identifier_for_crossreferencing
  fig-scrap="Short caption 4 list of figures as seen in LoF"
  fig-alt="Detailed alt text that describes the image content, type, purpose, and meaning.
    [CHART TYPE]: [Short description].
      DATA: [What data is shown, x/y axes].
      PURPOSE: [Why it's included, what to look for].
      DETAILS: [Longer description of patterns, anomalies, or key insights].
      SOURCE: Data from [source name/year and url/link]
    "
  fig-align="left"
  width="30%"
  }](https://github.com/VJMeyer/submission)
}}

```

Chapter 2

- ☒ {#fig-methodology}: Research methodology flowchart
 - File: images/methodology-flow.svg
 - Source: Author original

```

{{
[![Example Caption/Title 4] (/images/cover.png){
  #fig-Unique_identifier_for_crossreferencing
  fig-scap="Short caption 4 list of figures as seen in LoF"
  fig-alt="Detailed alt text that describes the image content, type, purpose, and meaning.
    [CHART TYPE]: [Short description].
      DATA: [What data is shown, x/y axes].
      PURPOSE: [Why it's included, what to look for].
      DETAILS: [Longer description of patterns, anomalies, or key insights].
      SOURCE: Data from [source name/year and url/link]

  "
  fig-align="left"
  width="30%"
  }](https://github.com/VJMeyer/submission)
}}

```

Pending Figures

```

### High Priority
- [ ] {#fig-results-chart}: Main results visualization
  - Status: Data ready, needs visualization

{{
[![Example Caption/Title 4] (/images/cover.png){
  #fig-Unique_identifier_for_crossreferencing
  fig-scap="Short caption 4 list of figures as seen in LoF"
  fig-alt="Detailed alt text that describes the image content, type, purpose, and meaning.
    [CHART TYPE]: [Short description].
      DATA: [What data is shown, x/y axes].
      PURPOSE: [Why it's included, what to look for].
      DETAILS: [Longer description of patterns, anomalies, or key insights].
      SOURCE: Data from [source name/year and url/link]

  "
  fig-align="left"
  width="30%"
  }](https://github.com/VJMeyer/submission)
}}

```

```

### Medium Priority
- [ ] {#fig-architecture}: System architecture diagram
  - Status: Sketch complete, needs professional rendering

{{
[![Example Caption/Title 4] (/images/cover.png){
  #fig-Unique_identifier_for_crossreferencing
  fig-scap="Short caption 4 list of figures as seen in LoF"
  fig-alt="Detailed alt text that describes the image content, type, purpose, and meaning.
    [CHART TYPE]: [Short description].
      DATA: [What data is shown, x/y axes].
      PURPOSE: [Why it's included, what to look for].
      DETAILS: [Longer description of patterns, anomalies, or key insights].
      SOURCE: Data from [source name/year and url/link]
    "
  fig-align="left"
  width="30%"
  }](https://github.com/VJMeyer/submission)
}}

```

Master Citation Registry

```

## BibTeX of Main Citations Included

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

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

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

## Update in ref/MAref.bib

## Core Citations (Must Have)

### Foundational Works
- [x] @carlsmith2021 - Power-seeking AI framework
  - Chapter usage: 1, 2, 4
  - Key concepts: Six premises, existential risk

```



```

- Notes: Central to thesis argument

- [x] @bostrom2014 - Superintelligence paths
  - Chapter usage: 1, 2, 3, 5
  - Key concepts: Orthogonality, convergence
  - Notes: Historical foundation

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

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

@article{bostrom2016,
  title = {The {{Unilateralist}}'s {{Curse}} and the {{Case}} for a {{Principle}} of {{Confidence}}},
  author = {Bostrom, Nick and Douglas, Thomas and Sandberg, Anders},
  date = {2016},
  journaltitle = {Social Epistemology},
  volume = {30},
  number = {4},
  pages = {350--371},
  publisher = {Routledge, part of the Taylor & Francis Group},

```

```
doi = {10.1080/02691728.2015.1108373},
url = {https://www.tandfonline.com/doi/full/10.1080/02691728.2015.1108373}
}
```

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

Pending Citations

Need to Find

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

Need to Verify

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

Citation Health Check

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

Bibliography

AMTAIR Thesis Relevant Literature & Citations

Items from MAref.bib

@carlsmith2021: Carlsmith [13]

Carlsmith, Joseph (2021)

Is Power-Seeking AI an Existential Risk?

DOI: 10.48550/arXiv.2206.13353

arXiv ID: 2206.13353

Better alternative: None - this is the primary case study

Relevant thesis section(s):

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

Potential claims supported (with certainty %):

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

@bostrom2014: Bostrom [9]

Bostrom, Nick (2014)

Superintelligence: Paths, Dangers, Strategies

ISBN: 978-0-19-967811-2

Better alternative: None - foundational text

Relevant thesis section(s):

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

Potential claims supported (with certainty %):

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

@clarke2022: Clarke et al. [18]

Clarke, Sam et al. (2022)

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

DOI: 10.48550/ARXIV.2206.09360

arXiv ID: 2206.09360

Better alternative: None - this is what AMTAIR builds upon

Relevant thesis section(s):

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

Potential claims supported (with certainty %):

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

@pearl2009 and @pearl2000: Pearl [55] and Pearl [54]

Pearl, Judea (2009)

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

ISBN: 978-0-521-89560-6

DOI: 10.1017/CB09780511803161

Better alternative: None - theoretical foundation

Relevant thesis section(s):

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

Potential claims supported (with certainty %):

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

@jaynes2003: Jaynes [37]

Jaynes, Edwin T. (2003)

Probability Theory: The Logic of Science

ISBN: 978-0-521-59271-0

DOI: 10.1017/CB09780511790423

Better alternative: None for foundational probability theory

Relevant thesis section(s):

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

Potential claims supported (with certainty %):

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

@tetlock2015: Tetlock and Gardner [68]

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

Superforecasting: The Art and Science of Prediction

ISBN: 978-0-8041-3671-6

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

Relevant thesis section(s):

- Section 1.5.2: Live Data Integration
- Section 3.9: Integration with Prediction Markets
- Forecasting methodology context

Potential claims supported (with certainty %):

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

@lempert2003: Lempert, Popper, and Bankes [44]

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

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

ISBN: 978-0-8330-3275-8

Better alternative: None for deep uncertainty methods

Relevant thesis section(s):

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

Potential claims supported (with certainty %):

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

@good1966: Good [29]

Good, Irving John (1966)

Speculations Concerning the First Ultraintelligent Machine

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

Relevant thesis section(s):

- Historical context in Introduction
- Background for intelligence explosion concept

Potential claims supported (with certainty %):

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

@yudkowsky2008: Yudkowsky [74]

Yudkowsky, Eliezer (2008)

Artificial Intelligence as a Positive and Negative Factor in Global Risk

DOI: 10.1093/oso/9780198570509.003.0021

Better alternative: @yudkowsky2022 for more recent formulation

Relevant thesis section(s):

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

Potential claims supported (with certainty %):

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

@russell2015: Russell et al. [60]

Russell, Stuart et al. (2015)

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

DOI: 10.1609/aimag.v36i4.2621

Better alternative: None - important consensus document

Relevant thesis section(s):

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

Potential claims supported (with certainty %):

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

New Suggested Citations

New Items to Consider:

@amodei2016: Amodei et al. [1]

Amodei, Dario et al. (2016)

Concrete Problems in AI Safety

arXiv ID: 1606.06565

Relevant thesis section(s):

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

Potential claims supported (with certainty %):

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

@christiano2019: Christiano [17]

Christiano, Paul (2019)
What Failure Looks Like

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

Relevant thesis section(s):

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

Potential claims supported (with certainty %):

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

@critch2019: critch2019

Critch, Andrew (2019)
ARCHES: AI Research Considerations for Human Existential Safety

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

Relevant thesis section(s):

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

Potential claims supported (with certainty %):

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

@dafoe2018 and updated @dafoe2021: Dafoe [21] and Dafoe [20]

Dafoe, Allan (2021)
AI Governance: A Research Agenda

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

Relevant thesis section(s):

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

Potential claims supported (with certainty %):

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

@askell12021: Askill et al. [4]

Askill, Amanda et al. (2021)

A General Language Assistant as a Laboratory for Alignment

arXiv ID: 2112.00861

Relevant thesis section(s):

- LLM capabilities for extraction tasks
- Alignment considerations for AMTAIR

Potential claims supported (with certainty %):

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

Further Citations to Integrate:

Growiec [30]

Clarke et al. [18]

Drexler [26] and Drexler [25]

Brundage et al. [11] and Brundage et al. [12]

Kumar et al. [42] and Kumar et al. [43]

Carlsmith [13] and Carlsmith [14] and Carlsmith [15]

Hendrycks et al. [34] and Hendrycks et al. [35]

Wilson et al. [72]

Appendices

Appendices

Appendix A: Technical Implementation Details

Contents:

- Full API specifications
- Architectural diagrams with component details
- Code structure and organization
- Deployment instructions
- Performance optimization guides

A.1 Core Data Structures

The AMTAIR system employs several custom data structures optimized for representing hierarchical arguments with probabilistic metadata:

A.2 Extraction Algorithm Details

A.3 API Specifications

Appendix B: Validation Datasets and Procedures

B.1 Expert Annotation Protocol

B.2 Benchmark Dataset Construction

B.3 Validation Results

Appendix C: Extended Case Studies

Include:

- Christiano's "What failure looks like"
- Critch's ARCHES model
- Additional policy evaluation scenarios
- Comparative analysis across models

C.1 Christiano’s “What Failure Looks Like” Extraction

C.2 Critch’s ARCHES Model

C.3 Policy Evaluation: A Narrow Path

Appendix D: BayesDown Syntax Specification

Contents:

- Full syntax definition
- Validation rules
- Example transformations
- Implementation notes
- Extension possibilities

Appendix E: Prompt Engineering Details

Include:

- Full extraction prompts with annotations
- Iterative refinement history
- Ablation study results
- Best practices guide
- Common failure patterns

Appendix F: User Guide

Sections:

- Getting started with AMTAIR
- Creating your first extraction
- Interpreting visualizations
- Policy evaluation walkthrough
- Troubleshooting common issues

Appendix G: Jupyter Notebook Implementation

The complete implementation is available as an interactive Jupyter notebook demonstrating:

- > Environment setup and configuration
- > Step-by-step extraction pipeline
- > Visualization generation
- > Policy evaluation examples
- > Performance benchmarking

Appendix H: Ethical Considerations and Governance

H.1 Potential Misuse Scenarios

H.2 Democratic Participation Frameworks

H.3 Responsibility Assignment

Appendix I: Full Extraction Examples

Appendix J: Software Installation and Usage Guide

AMTAIR Prototype Demonstration (Public Colab Notebook)

AMTAIR Prototype: Automating Transformative AI Risk Modeling

Executive Summary

This notebook implements a prototype of the AMTAIR (Automating Transformative AI Risk Modeling) project, which addresses the critical coordination failure in AI governance by developing computational tools that automate the extraction of probabilistic world models from AI safety literature.

The prototype demonstrates the transformation pipeline from structured argument representations (ArgDown) to probabilistic Bayesian networks (BayesDown), enabling the visualization and analysis of causal relationships and probability distributions that underlie AI risk assessments and policy evaluations.

Purpose Within the Master's Thesis

This notebook serves as the technical implementation component of the Master's thesis "Automating Transformative AI Risk Modeling: A Computational Approach to Policy Impact Evaluation." It demonstrates the feasibility of automating the extraction and formalization of world models, focusing on the core extraction pipeline and visualization capabilities that form the foundation for more sophisticated analysis.

Relevance to AI Governance

The coordination crisis in AI governance stems from different stakeholders working with incompatible assumptions, terminologies, and priorities. By making implicit models explicit through automated extraction and formalization, this work helps bridge communication gaps between technical researchers, policy specialists, and other stakeholders, contributing to more effective coordination in addressing existential risks from advanced AI.

Notebook Structure and Workflow

This notebook implements a multi-stage pipeline for transforming argument structures into interactive Bayesian network visualizations:

0. **Environment Setup** (Sections 0.1-0.3): Establishes the technical environment with necessary libraries and data connections
1. **Argument Extraction** (Sections 1.0-1.8): Processes source documents into structured ArgDown representations
2. **Probability Integration** (Sections 2.0-2.8): Enhances ArgDown with probability information to create BayesDown
3. **Data Transformation** (Section 3.0): Converts BayesDown into structured DataFrame format
4. **Visualization and Analysis** (Section 4.0): Creates interactive Bayesian network visualizations
5. **Archiving and Export** (Sections 5.0-6.0): Provides utilities for saving and sharing results

Throughout this notebook, we use the classic rain-sprinkler-lawn example as a canonical test case, demonstrating how a simple causal scenario (rain and sprinkler use affecting wet grass) can be represented, processed, and visualized using our automated pipeline.

Instructions — How to use this notebook:

1. **Import Libraries & Install Packages:** Run Section 0.1 to set up the necessary dependencies for data processing and visualization.
2. **Connect to GitHub Repository & Load Data files:** Run Section 0.2 to establish connections to the data repository and load example datasets. This step retrieves sample ArgDown files and extracted data for demonstration.
3. **Process Source Documents to ArgDown:** Sections 1.0-1.8 demonstrate the extraction of argument structures from source documents (such as PDFs) into ArgDown format, a markdown-like notation for structured arguments.
4. **Convert ArgDown to BayesDown:** Sections 2.0-2.3 handle the transformation of ArgDown files into BayesDown format, which incorporates probabilistic information into the argument structure.
5. **Extract Data into Structured Format:** Section 3.0 processes BayesDown format into structured database entries (CSV) that can be used for analysis.
6. **Create and Analyze Bayesian Networks:** Section 4.0 demonstrates how to build Bayesian networks from the extracted data and provides tools for analyzing risk pathways.
7. **Save and Export Results:** Sections 5.0-6.0 provide methods for archiving results and exporting visualizations.

AMTAIR Prototype Demonstration (Public Colab Notebook)

AMTAIR Prototype: Automating Transformative AI Risk Modeling

Executive Summary

Purpose Within the Master's Thesis

Relevance to AI Governance

Notebook Structure and Workflow

Instructions — How to use this notebook:

Key Concepts:

Example Workflow:

Troubleshooting:

0 Environment Setup and Data Access

0.1 Prepare Colab/Python Environment — Import Libraries & Packages

0.2 Connect to GitHub Repository

0.3 File Import

1 Sources (PDF's of Papers) to ArgDown (.md file)

1.0 Sources to ArgDown: Structured Argument Extraction

Process Overview

What is ArgDown?

1.1 Specify Source Document (e.g. PDF)

1.2 Generate ArgDown Extraction Prompt

1.3 Prepare LLM API Call

1.4 Make ArgDown Extraction LLM API Call

1.5 Save ArgDown Extraction Response

1.6 Review and Check ArgDown.md File

1.6.0 Check the Graph Structure with the ArgDown Sandbox Online

1.7 Extract ArgDown Graph Information as DataFrame

1.8 Store ArgDown Information as 'ArgDown.csv' file

2 Probability Extractions: ArgDown (.csv) to BayesDown (.md + plugin JSON syntax)

2.0 ArgDown to BayesDown: Adding Probability Information

Process Overview

What is BayesDown?

- 2.1 Probability Extraction Questions — ‘ArgDown.csv’ to ‘ArgDown_WithQuestions.csv’
- 2.2 ‘ArgDown_WithQuestions.csv’ to ‘BayesDownQuestions.md’
- 2.3 Generate BayesDown Probability Extraction Prompt
 - 2.3.0 BayesDown Format Specification
 - Core Structure
 - 2.3.1 Rain-Sprinkler-Lawn Example
- 2.4 Prepare 2nd API call
- 2.5 Make BayesDown Probability Extraction API Call
- 2.6 Save BayesDown with Probability Estimates (.csv)
- 2.7 Review & Verify BayesDown Probability Estimates
 - 2.7.2 Check the Graph Structure with the ArgDown Sandbox Online
- 2.8 Extract BayesDown with Probability Estimates as Dataframe
- 3 Data Extraction: BayesDown (.md) to Database (.csv)
 - 3.0 BayesDown to Structured Data: Network Construction
- Extraction Pipeline Overview
 - Theoretical Foundation
 - Role in Thesis Research
 - 3.0.0 ExtractBayesDown-Data_v1
 - 3.0.1 Test BayesDown Extraction
 - 3.0.2 Check the Graph Structure with the ArgDown Sandbox Online
- 3.1 Extraction
- 3.2 Data-Post-Processing
- 3.4 Download and save finished data frame as .csv file
- 4 Analysis & Inference: Bayesian Network Visualization
 - 4.0 Bayesian Network Visualization Approach
 - Visualization Philosophy
 - Connection to AMTAIR Goals
 - Implementation Structure
 - 4.1 Phase 1: Dependencies/Functions

- 4.2 Phase 2: Node Classification and Styling Module
- 4.3 Phase 3: HTML Content Generation Module
- 4.4 Phase 4: Main Visualization Function
- 5 Quick check HTML Outputs
- Conclusion: From Prototype to Production
 - Summary of Achievements
 - Limitations and Future Work
 - Connection to AMTAIR Project
- 6 Save Outputs
 - 6.0 Saving and Exporting Results
 - 6.1 Convert .ipynb Notebook to Markdown
 - 6.2 Convert Notebook to Markdown Documentation
 - 6.3 Create PDF and Latex

Key Concepts:

- > **ArgDown:** A structured format for representing arguments, with hierarchical relationships between statements.
- > **BayesDown:** An extension of ArgDown that incorporates probabilistic information, allowing for Bayesian network construction.
- > **Extraction Pipeline:** The process of converting unstructured text to structured argument representations.
- > **Bayesian Networks:** Probabilistic graphical models that represent variables and their conditional dependencies.

Example Workflow:

1. Load a sample ArgDown file from the repository
2. Extract the hierarchical structure and relationships
3. Add probabilistic information to create a BayesDown representation
4. Generate a Bayesian network visualization
5. Analyze conditional probabilities and risk pathways

Troubleshooting:

- > If connectivity issues occur, ensure you have access to the GitHub repository
- > For visualization errors, check that all required libraries are properly installed
- > When processing custom files, ensure they follow the expected format conventions

0 Environment Setup and Data Access

This section establishes the technical foundation for the AMTAIR prototype by: 1. Installing and importing necessary libraries 2. Setting up access to the GitHub repository 3. Loading example data files

The environment setup is designed to be run once per session, with flags to prevent redundant installations and imports. This section forms the basis for the subsequent extraction and analysis steps in the pipeline.

The key goal is to create a reproducible environment where the Bayesian network extraction and visualization can be performed consistently, with appropriate error handling and resource management.

0.1 Prepare Colab/Python Environment — Import Libraries & Packages

```
# @title 0.1.0 --- Install & Import Libraries & Packages (One-Time Setup) --- [install_import]

"""
BLOCK PURPOSE:
Establishes the core technical environment for the AMTAIR prototype.
Sets up all required libraries for Bayesian network processing, visualization,
and data manipulation.
Uses a flag-based approach to ensure setup only runs once per session,
enhancing efficiency.

The setup follows a three-stage process:
1. Install required packages not available in Colab by default
2. Import all necessary libraries with error handling
3. Set a global flag to prevent redundant execution

DEPENDENCIES: Requires internet connection for package installation
```

```

OUTPUTS: Global variable _setup_imports_done and loaded Python libraries
"""

# Check if setup has already been completed in this session using environment flag
try:
    # If this variable exists, setup was already done successfully
    _setup_imports_done
    print(" Libraries already installed and imported in this session. Skipping setup.")

except NameError:
    print(" Performing one-time library installation and imports...")

    # --- STAGE 1: Install required packages ---
    # Install visualization and network analysis libraries
    !pip install -q pyvis # Network visualization library
    !apt-get install pandoc -y # Document conversion utility

    # Install Google API and data processing packages
    # Data manipulation and Google integration
    !pip install -q --upgrade gspread pandas google-auth google-colab

    # Install Bayesian network and probabilistic modeling tools
    !pip install -q pgmpy # Probabilistic graphical models library

    # Install notebook conversion tools
    !pip install -q nbconvert # Often pre-installed, but ensures availability

    print(" --> Installations complete.")

    # --- STAGE 2: Import libraries with error handling ---
    try:
        # Network and HTTP libraries
        import requests # For making HTTP requests to APIs and GitHub
        import io # For handling in-memory file-like objects

        # Data processing libraries
        import pandas as pd # For structured data manipulation
        import numpy as np # For numerical operations
        import json # For JSON parsing and serialization
        import re # For regular expression pattern matching

        # Visualization libraries

```

```

import matplotlib.pyplot as plt # For creating plots and charts
from IPython.display import HTML, display, Markdown # For rich output in notebook

# --- Specialized libraries requiring installation ---
# Network analysis library
import networkx as nx # For graph representation and analysis

# Probabilistic modeling libraries
from pgmpy.models import BayesianNetwork # For Bayesian network structure
from pgmpy.factors.discrete import TabularCPD # For conditional probability tables
from pgmpy.inference import VariableElimination # For probabilistic inference

# Interactive network visualization
from pyvis.network import Network # For interactive network visualization

# Output version information for key libraries
print(f"    pandas version: {pd.__version__}")
print(f"    networkx version: {nx.__version__}")
# Add others if specific versions are critical

print("    --> Imports complete.")

# --- STAGE 3: Set flag to indicate successful setup ---
_setup_imports_done = True
print("    One-time setup finished successfully.")

except ImportError as e:
    # Handle specific import failures
    print(f"    ERROR during import: {e}")
    print("    --> Setup did not complete successfully. Please check installations.")
except Exception as e:
    # Handle unexpected errors
    print(f"    UNEXPECTED ERROR during setup: {e}")
    print("    --> Setup did not complete successfully.")

# Environment is now ready for AMTAIR processing

```

Libraries already installed and imported in this session. Skipping setup.

0.2 Connect to GitHub Repository

The Public GitHub Repo Url in use:

https://raw.githubusercontent.com/SingularitySmith/AMTAIR_Prototype/main/

Note: When encountering errors, accessing the data, try using “RAW” Urls.

```
# @title 0.2.0 --- Connect to GitHub Repository --- Load Files [connect_to_github_repository]

"""
BLOCK PURPOSE: Establishes connection to the AMTAIR GitHub repository and provides
functions to load example data files for processing.

This block creates a reusable function for accessing files from the project's
GitHub repository, enabling access to example files like the rain-sprinkler-lawn
Bayesian network that serves as our canonical test case.

DEPENDENCIES: requests library, io library
OUTPUTS: load_file_from_repo function and test file loads
"""

from requests.exceptions import HTTPError

# Specify the base repository URL for the AMTAIR project
repo_url = "https://raw.githubusercontent.com/SingularitySmith/AMTAIR_Prototype/main/data/ex
print(f"Connecting to repository: {repo_url}")

def load_file_from_repo(relative_path):
    """
    Loads a file from the specified GitHub repository using a relative path.

    Args:
        relative_path (str): Path to the file relative to the repo_url

    Returns:
        For CSV/JSON: pandas DataFrame
        For MD: string containing file contents

    Raises:
        HTTPError: If file not found or other HTTP error occurs
        ValueError: If unsupported file type is requested
    """
    file_url = repo_url + relative_path
    print(f"Attempting to load: {file_url}")

    # Fetch the file content from GitHub
    response = requests.get(file_url)
```

```

# Check for bad status codes with enhanced error messages
if response.status_code == 404:
    raise HTTPError(f"File not found at URL: {file_url}. Check the file path/name and en
else:
    response.raise_for_status() # Raise for other error codes

# Convert response to file-like object
file_object = io.StringIO(response.text)

# Process different file types appropriately
if relative_path.endswith(".csv"):
    return pd.read_csv(file_object) # Return DataFrame for CSV
elif relative_path.endswith(".json"):
    return pd.read_json(file_object) # Return DataFrame for JSON
elif relative_path.endswith(".md"):
    return file_object.read() # Return raw content for MD files
else:
    raise ValueError(f"Unsupported file type: {relative_path.split('.')[-1]}. Add support

# Load example files to test connection
try:
    # Load the extracted data CSV file
    # df = load_file_from_repo("extracted_data.csv")

    # Load the ArgDown test text
    md_content = load_file_from_repo("ArgDown.md")

    print(" Successfully connected to repository and loaded test files.")
except Exception as e:
    print(f" Error loading files: {str(e)}")
    print("Please check your internet connection and the repository URL.")

# Display preview of loaded content (commented out to avoid cluttering output)
print(md_content)

```

Connecting to repository: https://raw.githubusercontent.com/SingularitySmith/AMTAIR_Prototype/main/main.py

Attempting to load: https://raw.githubusercontent.com/SingularitySmith/AMTAIR_Prototype/main/main.py

Successfully connected to repository and loaded test files.

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

- [Human_Disempowerment]: Permanent and collective disempowerment of humanity relative to AI

- [Scale_Of_Power_Seeking]: Power-seeking by AI systems scaling to the point of permanent

- [Misaligned_Power_Seeking]: Deployed AI systems seeking power in unintended and hi

- [APS_Systems]: AI systems with advanced capabilities, agentic planning, and st
 - [Advanced_AI_Capability]: AI systems that outperform humans on tasks that
 - [Agentic_Planning]: AI systems making and executing plans based on world m
 - [Strategic_Awareness]: AI systems with models accurately representing power
- [Difficulty_Of_Alignment]: It is harder to build aligned systems than misalign
 - [Instrumental_Convergence]: AI systems with misaligned objectives tend to
 - [Problems_With_Proxies]: Optimizing for proxy objectives breaks correlatio
 - [Problems_With_Search]: Search processes can yield systems pursuing differ
- [Deployment_Decisions]: Decisions to deploy potentially misaligned AI systems.
 - [Incentives_To_Build_APS]: Strong incentives to build and deploy APS syste
 - [Usefulness_Of_APS]: APS systems are very useful for many valuable tas
 - [Competitive_Dynamics]: Competitive pressures between AI developers. {
 - [Deception_By_AI]: AI systems deceiving humans about their true objectives
- [Corrective_Feedback]: Human society implementing corrections after observing prob
 - [Warning_Shots]: Observable failures in weaker systems before catastrophic ris
 - [Rapid_Capability_Escalation]: AI capabilities escalating very rapidly, allowi

[Barriers_To_Understanding]: Difficulty in understanding the internal workings of advanced A

- [Misaligned_Power_Seeking]: Deployed AI systems seeking power in unintended and high-impac

[Adversarial_Dynamics]: Potentially adversarial relationships between humans and power-seeking

- [Misaligned_Power_Seeking]: Deployed AI systems seeking power in unintended and high-impac

[Stakes_Of_Error]: The escalating impact of mistakes with power-seeking AI systems. {"instan

- [Misaligned_Power_Seeking]: Deployed AI systems seeking power in unintended and high-impac

```
# Specify the relative path to the HTML file
html_file_path = "bayesian_network.html"

try:
    # Load the HTML file content using the existing function
    # The function returns raw content (string) for .md files, and we'll treat .html similar
    # We'll modify the function's behavior slightly if needed, or handle the string directly
    html_content = load_file_from_repo(html_file_path)

    print(f" Successfully loaded {html_file_path}.")

    # Render the HTML content directly in the notebook
    display(HTML(html_content))

except ValueError as e:
    # Handle the case where the function might raise ValueError for unsupported types
    print(f" Error loading HTML file: {e}")
    print("Make sure the load_file_from_repo function supports .html or handle the string co")
except Exception as e:
    # Catch any other potential errors during loading or display
```

```
print(f" Error loading or displaying {html_file_path}: {str(e)}")  
print("Please check the file path and your internet connection.")
```

Attempting to load: https://raw.githubusercontent.com/SingularitySmith/AMTAIR_Prototype/main/bayesian_network.html

Error loading or displaying bayesian_network.html: File not found at URL: https://raw.githubusercontent.com/SingularitySmith/AMTAIR_Prototype/main/bayesian_network.html

Please check the file path and your internet connection.

0.3 File Import

```
# @title  
md_content
```

```
'[Existential_Catastrophe]: The destruction of humanity\'s long-term potential due to AI sys
```


1 Sources (PDF's of Papers) to ArgDown (.md file)

1.0 Sources to ArgDown: Structured Argument Extraction

Process Overview

This section implements the first major stage of the AMTAIR pipeline: transforming source documents (such as research papers, blog posts, or expert analyses) into structured argument representations using the ArgDown format.

ArgDown is a markdown-like notation for representing arguments in a hierarchical structure. In the context of AMTAIR, it serves as the first step toward creating formal Bayesian networks by:

1. Identifying key variables/statements in the text
2. Capturing their hierarchical relationships
3. Preserving their descriptive content
4. Defining their possible states (instantiations)

The extraction process uses Large Language Models (LLMs) to identify the structure and relationships in the text, though in this notebook we focus on processing pre-formatted examples rather than performing the full extraction from raw text.

What is ArgDown?

ArgDown uses a simple syntax where:

- Statements are represented as `[Statement]: Description`
- Relationships are indicated with `+` symbols and indentation
- Metadata is added in JSON format, including possible states of each variable

For example:

```
[MainClaim]: Description of the main claim. {"instantiations": ["claim_TRUE", "claim_FALSE"]}

+ [SupportingEvidence]: Description of evidence. {"instantiations": ["evidence_TRUE", "evidence_FALSE"]}
```

This structure will later be enhanced with probability information to create BayesDown, which can be transformed into a Bayesian network for analysis and visualization.

1.1 Specify Source Document (e.g. PDF)

Review the source document, ensure it is suitable for API call and upload to / store it in the correct location.

```
# @title 1.1.0 --- MTAIR Online Model (Analytica) --- [online_model]

from IPython.display import IFrame

IFrame(src="https://acp.analytica.com/view0?invite=4560&code=3000289064591444815", width="100%", height="400px")

<IPython.lib.display.IFrame at 0x7b9cc9929f50>

MTAIR Online Model (Analytica)
```

1.2 Generate ArgDown Extraction Prompt

Generate Extraction Prompt

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

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

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

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

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

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

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

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

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

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

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

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

## Epistemic Foundation & Purpose

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

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

## ArgDown Format Specification

### Core Syntax

ArgDown represents causal relationships using a hierarchical structure:

1. Variables appear in square brackets with descriptive text:
```

```
`[Variable_Name]: Description of the variable.`

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

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

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

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

### JSON Metadata Format

The JSON metadata must follow this exact structure:

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

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

For binary variables (most common case):
{"instantiations": ["variable_TRUE", "variable_FALSE"]}

For multi-state variables (when clearly specified in the text):
{"instantiations": ["variable_HIGH", "variable_MEDIUM", "variable_LOW"]}

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

Complex Structural Patterns

Variables Influencing Multiple Effects

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

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

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

```
+ [Cause_A]
+ [Cause_B]: Description of cause B. {"instantiations": ["cause_b_TRUE", "cause_b_FALSE"]}

Multiple Causes of the Same Effect
Multiple causes can influence the same effect by being listed at the same indentation level.
[Effect]: Description of effect. {"instantiations": ["effect_TRUE", "effect_FALSE"]}
+ [Cause1]: Description of first cause. {"instantiations": ["cause1_TRUE", "cause1_FALSE"]}
+ [Cause2]: Description of second cause. {"instantiations": ["cause2_TRUE", "cause2_FALSE"]}
+ [Deeper_Cause]: A cause that influences Cause2. {"instantiations": ["deeper_cause_TRUE", "deeper_cause_FALSE"]}

Causal Chains
Causal chains are represented through multiple levels of indentation:
[Ultimate_Effect]: The final outcome. {"instantiations": ["ultimate_effect_TRUE", "ultimate_effect_FALSE"]}
+ [Intermediate_Effect]: A mediating variable. {"instantiations": ["intermediate_effect_TRUE", "intermediate_effect_FALSE"]}
+ [Root_Cause]: The initial cause. {"instantiations": ["root_cause_TRUE", "root_cause_FALSE"]}
+ [2nd_Intermediate_Effect]: A mediating variable. {"instantiations": ["intermediate_effect2_TRUE", "intermediate_effect2_FALSE"]}

Common Cause of Multiple Variables
A common cause affecting multiple variables is represented by referencing the same variable name.
[Effect1]: First effect description. {"instantiations": ["effect1_TRUE", "effect1_FALSE"]}
+ [Common_Cause]: Description of common cause. {"instantiations": ["common_cause_TRUE", "common_cause_FALSE"]}

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

Detailed Extraction Workflow
Please follow this step-by-step process, documenting your reasoning in XML tags:
<analysis>
First, conduct a holistic analysis of the document:
1. Identify the main subject matter or domain
2. Note key concepts, variables, and factors discussed
3. Pay attention to language indicating causal relationships (causes, affects, influences, etc.)
4. Look for the ultimate outcomes or effects that are the focus of the document
5. Record your general understanding of the document's implicit causal structure
</analysis>
<variable_identification>
Next, identify and list the key variables in the causal model:
* Focus on factors that are discussed as having an influence or being influenced
* For each variable:
 * Create a descriptive name in [square_brackets]
 * Write a concise description based directly on the text
```

```
* Determine possible states (usually binary TRUE/FALSE unless clearly specified)
* Distinguish between:
 * Outcome variables (effects the author is concerned with)
 * Intermediate variables (both causes and effects in chains)
 * Root cause variables (exogenous factors in the model)
* List all identified variables with their descriptions and possible states
</variable_identification>

<causal_structure>
Then, determine the causal relationships between variables:
* For each variable, identify what factors influence it
* Note the direction of causality (what causes what)
* Look for mediating variables in causal chains
* Identify common causes of multiple effects
* Capture feedback loops if present (though they must be represented as DAGs)
* Map out the hierarchical structure of the causal model
</causal_structure>

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

<validation>

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

</validation>

Source Document Analysis Guidance
```

When analyzing the source document:

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

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

""")

```
BayesDown probability extraction prompt - enhances ArgDown with probability information
BAYESDOWN_EXTRACTION = PromptTemplate("""
```

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

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

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

Here is the format to follow:

```
[Node]: Description. { "instantiations": ["node_TRUE", "node_FALSE"], "priors": { "p(node_TRUE)": 0.5, "p(node_FALSE)": 0.5 } }
[Parent]: Parent description. {...}
```

Here are the specific probability questions to answer:

\$questions

ArgDown structure to enhance:

\$argdown

Provide the complete BayesDown representation with probabilities:

""")

```
@classmethod
def get_template(cls, template_name: str) -> PromptTemplate:
 """Get a prompt template by name"""
 if hasattr(cls, template_name):
```

```

 return getattr(cls, template_name)
 else:
 raise ValueError(f"Template not found: {template_name}")

```

## 1.3 Prepare LLM API Call

Combine Systemprompt + API Specifications + ArgDown Instructions + Prompt + Source PDF for API Call

```

@title 1.3.0 --- Provider-Agnostic LLM API Interface --- [provider_agnostic-interface]

"""
BLOCK PURPOSE: Provides a unified interface for interacting with different LLM providers.

This block implements a flexible, provider-agnostic system for making LLM API calls:
1. Base abstract class (LLMProvider) defining the common interface
2. Implementation classes for specific providers (OpenAI and Anthropic)
3. Factory class for creating appropriate provider instances

This abstraction allows the extraction pipeline to work with different LLM providers
without changing the core code, supporting both current and future LLM backends.

DEPENDENCIES: requests for API calls, os for environment variables, abstract base classes
OUTPUTS: LLMProvider abstract class and concrete implementations for OpenAI and Anthropic
"""

import os
import json
import time
import requests
from abc import ABC, abstractmethod
from typing import Dict, List, Optional, Union, Any
from dataclasses import dataclass

@dataclass
class LLMResponse:
 """Standard response object for LLM completions"""
 content: str # The generated text response
 model: str # The model used for generation
 usage: Dict[str, int] # Token usage statistics
 raw_response: Dict[str, Any] # Complete provider-specific response
 created_at: float = time.time() # Timestamp of response creation

```



```

class LLMPProvider(ABC):
 """Abstract base class for LLM providers"""

 @abstractmethod
 def complete(self,
 prompt: str,
 system_prompt: Optional[str] = None,
 temperature: float = 0.7,
 max_tokens: int = 4000) -> LLMResponse:
 """Generate a completion from the LLM"""
 pass

 @abstractmethod
 def get_available_models(self) -> List[str]:
 """Return a list of available models from this provider"""
 pass

class OpenAIProvider(LLMPProvider):
 """OpenAI API implementation"""

 def __init__(self, api_key: Optional[str] = None, organization: Optional[str] = None):
 """Initialize with API key from args or environment"""
 self.api_key = api_key or os.environ.get("OPENAI_API_KEY")
 if not self.api_key:
 raise ValueError("OpenAI API key is required. Provide as argument "
 + "or set OPENAI_API_KEY environment variable.")

 self.organization = organization or os.environ.get("OPENAI_ORGANIZATION")
 self.api_base = "https://api.openai.com/v1"

 def complete(self,
 prompt: str,
 system_prompt: Optional[str] = None,
 model: str = "gpt-4-turbo",
 temperature: float = 0.7,
 max_tokens: int = 4000) -> LLMResponse:
 """Generate a completion using OpenAI's API"""

 # Prepare request headers
 headers = {
 "Content-Type": "application/json",

```

```

 "Authorization": f"Bearer {self.api_key}"
 }

 if self.organization:
 headers["OpenAI-Organization"] = self.organization

 # Create message structure
 messages = []
 if system_prompt:
 messages.append({"role": "system", "content": system_prompt})

 messages.append({"role": "user", "content": prompt})

 # Prepare request data
 data = {
 "model": model,
 "messages": messages,
 "temperature": temperature,
 "max_tokens": max_tokens
 }

 # Make API call
 response = requests.post(
 f"{self.api_base}/chat/completions",
 headers=headers,
 json=data
)

 response.raise_for_status()
 result = response.json()

 # Transform into standardized response format
 return LLMResponse(
 content=result["choices"][0]["message"]["content"],
 model=result["model"],
 usage=result["usage"],
 raw_response=result
)

def get_available_models(self) -> List[str]:
 """Return a list of available OpenAI models"""
 headers = {

```

```

 "Authorization": f"Bearer {self.api_key}"
 }

 if self.organization:
 headers["OpenAI-Organization"] = self.organization

 response = requests.get(
 f"{self.api_base}/models",
 headers=headers
)

 response.raise_for_status()
 models = response.json()["data"]
 return [model["id"] for model in models]

class AnthropicProvider(LLMProvider):
 """Anthropic Claude API implementation"""

 def __init__(self, api_key: Optional[str] = None):
 """Initialize with API key from args or environment"""
 self.api_key = api_key or os.environ.get("ANTHROPIC_API_KEY")
 if not self.api_key:
 raise ValueError("Anthropic API key is required. Provide as"
 + " argument or set ANTHROPIC_API_KEY environment variable.")

 self.api_base = "https://api.anthropic.com/v1"

 def complete(self,
 prompt: str,
 system_prompt: Optional[str] = None,
 model: str = "claude-3-opus-20240229",
 temperature: float = 0.7,
 max_tokens: int = 4000) -> LLMResponse:
 """Generate a completion using Anthropic's API"""

 # Prepare request headers
 headers = {
 "Content-Type": "application/json",
 "X-API-Key": self.api_key,
 "anthropic-version": "2023-06-01"
 }

```

```

Prepare request data in Anthropic-specific format
data = {
 "model": model,
 "messages": [{"role": "user", "content": prompt}],
 "temperature": temperature,
 "max_tokens": max_tokens
}

Add system prompt if provided (Anthropic uses a different format)
if system_prompt:
 data["system"] = system_prompt

Make API call
response = requests.post(
 f"{self.api_base}/messages",
 headers=headers,
 json=data
)

response.raise_for_status()
result = response.json()

Transform into standardized response format
return LLMResponse(
 content=result["content"][0]["text"],
 model=result["model"],
 usage={"prompt_tokens": result.get("usage", {}).get("input_tokens", 0),
 "completion_tokens": result.get("usage", {}).get("output_tokens", 0)},
 raw_response=result
)

def get_available_models(self) -> List[str]:
 """Return a list of available Anthropic models"""
 # Anthropic doesn't have a models endpoint, so we return a static list
 return [
 "claude-3-opus-20240229",
 "claude-3-sonnet-20240229",
 "claude-3-haiku-20240307"
]

class LLMFactory:
 """Factory for creating LLM providers"""

```

```

 @staticmethod
 def create_provider(provider_name: str, **kwargs) -> LLMProvider:
 """Create and return an LLM provider instance"""
 if provider_name.lower() == "openai":
 return OpenAIProvider(**kwargs)
 elif provider_name.lower() == "anthropic":
 return AnthropicProvider(**kwargs)
 else:
 raise ValueError(f"Unsupported provider: {provider_name}")

@title 1.3.1 --- API Call Function Definitions --- [api_call_function_definitions]

"""
BLOCK PURPOSE: Provides core functions for extracting ArgDown representations from text using LLMs.

This block implements the main extraction functionality:
1. extract_argdown_from_text: Sends text to LLM to extract structured ArgDown representation
2. validate_argdown: Verifies the extracted ArgDown for correctness and completeness
3. process_source_document: Handles source files (PDF, TXT, MD) and manages extraction
4. save_argdown_extraction: Saves extraction results with metadata for further processing

These functions form the first stage of the AMTAIR pipeline, transforming
unstructured text into structured argument representations.

DEPENDENCIES: LLMFactory from previous cell, re for pattern matching
OUTPUTS: Functions for ArgDown extraction, validation, and storage
"""

def extract_argdown_from_text(text: str, provider_name: str = "openai", model: str = None) -> ArgDown:
 """
 Extract ArgDown representation from text using LLM

 Args:
 text: The source text to extract arguments from
 provider_name: The LLM provider to use (openai or anthropic)
 model: Specific model to use, or None for default

 Returns:
 Extracted ArgDown representation
 """
 # Create LLM provider
 provider = LLMFactory.create_provider(provider_name)

```

```

Get extraction prompt
prompt_template = PromptLibrary.get_template("ARGDOWN_EXTRACTION")
prompt = prompt_template.format(text=text)

Set model-specific parameters
if provider_name.lower() == "openai":
 model = model or "gpt-4-turbo"
 temperature = 0.3 # Lower temperature for more deterministic extraction
 max_tokens = 4000
elif provider_name.lower() == "anthropic":
 model = model or "claude-3-opus-20240229"
 temperature = 0.2
 max_tokens = 4000

Call the LLM
system_prompt = "You are an expert in argument mapping and causal reasoning."
response = provider.complete(
 prompt=prompt,
 system_prompt=system_prompt,
 model=model,
 temperature=temperature,
 max_tokens=max_tokens
)

Extract the ArgDown content (remove any markdown code blocks if present)
argdown_content = response.content
if "```" in argdown_content:
 # Extract content between code blocks if present
 import re
 matches = re.findall(r"```(?:argdown)?\n([\s\S]*?)\n```", argdown_content)
 if matches:
 argdown_content = matches[0]

return argdown_content

def validate_argdown(argdown_text: str) -> Dict[str, Any]:
 """
 Validate ArgDown representation to ensure it's well-formed

 Args:
 argdown_text: ArgDown representation to validate

```

```

Returns:
 Dictionary with validation results
"""
Initialize validation results
results = {
 "is_valid": True,
 "errors": [],
 "warnings": [],
 "stats": {
 "node_count": 0,
 "relationship_count": 0,
 "max_depth": 0
 }
}

Basic syntax checks
lines = argdown_text.split("\n")
node_pattern = r'\[(.*?)\]:'
instantiation_pattern = r'{"instantiations":'

Track nodes and relationships
nodes = set()
relationships = []
current_depth = 0
max_depth = 0

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

 # Calculate indentation depth
 indent = 0
 if '+' in line:
 indent = line.find('+') // 2

 current_depth = indent
 max_depth = max(max_depth, current_depth)

 # Check for node definitions
 import re

```

```

node_matches = re.findall(node_pattern, line)
if node_matches:
 node = node_matches[0]
 nodes.add(node)
 results["stats"]["node_count"] += 1

 # Check for instantiations
 if instantiation_pattern not in line:
 results["warnings"].append(f"Line {i+1}: Node '{node}' is missing instantiation")

Check parent-child relationships
if indent > 0 and '+' in line and node_matches:
 # This is a child node; find its parent
 parent_indent = indent - 1
 j = i - 1
 while j >= 0:
 if '+' in lines[j] and lines[j].find('+') // 2 == parent_indent:
 parent_matches = re.findall(node_pattern, lines[j])
 if parent_matches:
 parent = parent_matches[0]
 relationships.append((parent, node))
 results["stats"]["relationship_count"] += 1
 break
 j -= 1

results["stats"]["max_depth"] = max_depth

If we didn't find any nodes, that's a problem
if results["stats"]["node_count"] == 0:
 results["is_valid"] = False
 results["errors"].append("No valid nodes found in ArgDown representation")

return results

def process_source_document(file_path: str, provider_name: str = "openai") -> Dict[str, Any]:
 """
 Process a source document to extract ArgDown representation

 Args:
 file_path: Path to the source document
 provider_name: The LLM provider to use

```



```

Returns:
 Dictionary with extraction results
"""
Load the source document
text = ""
if file_path.endswith(".pdf"):
 # PDF handling requires additional libraries
 try:
 import PyPDF2
 with open(file_path, 'rb') as file:
 reader = PyPDF2.PdfReader(file)
 text = ""
 for page in reader.pages:
 text += page.extract_text() + "\n"
 except ImportError:
 raise ImportError("PyPDF2 is required for PDF processing. "
 + "Install it with: pip install PyPDF2")
elif file_path.endswith(".txt"):
 with open(file_path, 'r') as file:
 text = file.read()
elif file_path.endswith(".md"):
 with open(file_path, 'r') as file:
 text = file.read()
else:
 raise ValueError(f"Unsupported file format: {file_path}")

Extract ArgDown
argdown_content = extract_argdown_from_text(text, provider_name)

Validate the extraction
validation_results = validate_argdown(argdown_content)

Prepare results
results = {
 "source_path": file_path,
 "extraction_timestamp": time.time(),
 "argdown_content": argdown_content,
 "validation": validation_results,
 "provider": provider_name
}

return results

```

```
def save_argdown_extraction(results: Dict[str, Any], output_path: str) -> None:
 """
 Save ArgDown extraction results

 Args:
 results: Extraction results dictionary
 output_path: Path to save the results
 """
 # Save the ArgDown content
 with open(output_path, 'w') as file:
 file.write(results["argdown_content"])

 # Save metadata alongside
 metadata_path = output_path.replace('.md', '_metadata.json')
 metadata = {
 "source_path": results["source_path"],
 "extraction_timestamp": results["extraction_timestamp"],
 "validation": results["validation"],
 "provider": results["provider"]
 }

 with open(metadata_path, 'w') as file:
 json.dump(metadata, file, indent=2)

@title 1.3.2 --- Prepare LLM API Call --- [prepare_api_call]

"""
BLOCK PURPOSE: Prepares parameters for LLM API calls used in ArgDown extraction.

This function handles the configuration for LLM API calls, including:
1. Source document path validation
2. LLM provider selection and validation
3. Model selection with appropriate defaults

The function returns a configuration dictionary that can be passed to the
extraction function in the next step of the pipeline.

DEPENDENCIES: None (uses standard Python functionality)
OUTPUTS: Dictionary with extraction configuration parameters
"""

def prepare_extraction_call(source_path, provider_name="openai", model=None):
```

```

"""
Prepare the LLM API call for ArgDown extraction

Args:
 source_path (str): Path to the source document to extract from
 provider_name (str): LLM provider to use ('openai' or 'anthropic')
 model (str, optional): Specific model to use. Defaults to None (uses provider's default)

Returns:
 dict: Configuration parameters for extraction

Raises:
 ValueError: If an unsupported provider is specified
"""

Load the source document
print(f"Processing source document: {source_path}")

Determine provider and model
provider = provider_name.lower()
if provider not in ["openai", "anthropic"]:
 raise ValueError(f"Unsupported provider: {provider}. Use 'openai' or 'anthropic'.")

Set default model if none provided
if model is None:
 if provider == "openai":
 model = "gpt-4-turbo"
 elif provider == "anthropic":
 model = "claude-3-opus-20240229"

Print configuration
print(f"Using provider: {provider}")
print(f"Selected model: {model}")

return {
 "source_path": source_path,
 "provider": provider,
 "model": model
}

Usage example:
source_path = "example_document.pdf" # Replace with actual document path
extraction_config = prepare_extraction_call(source_path, provider_name="openai")

```

## 1.4 Make ArgDown Extraction LLM API CallSources (PDF's of Papers) to ArgDown (.md file)

Processing source document: example\_document.pdf

Using provider: openai

Selected model: gpt-4-turbo

### 1.4 Make ArgDown Extraction LLM API Call

```
@title 1.4.0 --- Make ArgDown Extraction LLM API Call --- [extraction_api_call]

"""
BLOCK PURPOSE: Executes the ArgDown extraction process using the LLM API.

This function performs the actual extraction of ArgDown representations from
source documents:

1. Takes the configuration parameters prepared in the previous step
2. Processes the document using the LLM API
3. Validates the extraction results
4. Provides timing and statistics about the extraction

The extraction process transforms unstructured text into a structured argument
representation following the ArgDown syntax defined in the AMTAIR project.

DEPENDENCIES: process_source_document function from previous cells
OUTPUTS: Dictionary with extraction results including ArgDown content and validation info
"""

def execute_extraction(extraction_config):
 """
 Execute the ArgDown extraction using the LLM API

 Args:
 extraction_config (dict): Configuration parameters for extraction

 Returns:
 dict: Extraction results including ArgDown content and validation info

 Raises:
 Exception: For any errors during extraction
 """
 print(f"Starting extraction from {extraction_config['source_path']}")
 start_time = time.time()

 try:
```

```
Process the document
results = process_source_document(
 extraction_config["source_path"],
 provider_name=extraction_config["provider"]
)

Print success message
elapsed_time = time.time() - start_time
print(f"Extraction completed in {elapsed_time:.2f} seconds")
print(f"Extracted {results['validation']['stats']['node_count']} nodes with "
 f"{results['validation']['stats']['relationship_count']} relationships")

Print any warnings
if results['validation']['warnings']:
 print("\nWarnings:")
 for warning in results['validation']['warnings']:
 print(f"- {warning}")

return results

except Exception as e:
 print(f"Error during extraction: {str(e)}")
 raise

Usage example:
extraction_results = execute_extraction(extraction_config)
```

Starting extraction from example\_document.pdf

Error during extraction: PyPDF2 is required for PDF processing. Install it with: pip install

ImportError: PyPDF2 is required for PDF processing. Install it with: pip install PyPDF2

```

ModuleNotFoundError Traceback (most recent call last)
<ipython-input-19-fd592eb962ab> in process_source_document(file_path, provider_name)
 166 try:
--> 167 import PyPDF2
 168 with open(file_path, 'rb') as file:
ModuleNotFoundError: No module named 'PyPDF2'
During handling of the above exception, another exception occurred:
ImportError Traceback (most recent call last)
<ipython-input-21-27555067c1d2> in <cell line: 0>()
 59
 60 # Usage example:
```

```

---> 61 extraction_results = execute_extraction(extraction_config)

<ipython-input-21-27555067c1d2> in execute_extraction(extraction_config)
 35 try:
 36 # Process the document
---> 37 results = process_source_document(
 38 extraction_config["source_path"],
 39 provider_name=extraction_config["provider"]
<ipython-input-19-fd592eb962ab> in process_source_document(file_path, provider_name)
 172 text += page.extract_text() + "\n"
 173 except ImportError:
--> 174 raise ImportError("PyPDF2 is required for PDF processing. Install it with: pip install PyPDF2")
 175 elif file_path.endswith(".txt"):
 176 with open(file_path, 'r') as file:
ImportError: PyPDF2 is required for PDF processing. Install it with: pip install PyPDF2

NOTE: If your import is failing due to a missing package, you can
manually install dependencies using either !pip or !apt.
To view examples of installing some common dependencies, click the
"Open Examples" button below.

```

## 1.5 Save ArgDown Extraction Response

1. Save and log API return
2. Save ArgDown.md file for further Processing

```

@title 1.5.0 --- Save ArgDown Extraction Response --- [save_extraction_response]

"""
BLOCK PURPOSE: Saves the extracted ArgDown content to files for further processing.

This function handles saving the extraction results:
1. Creates an output directory if it doesn't exist
2. Saves the extracted ArgDown content with a timestamp in the filename
3. Saves accompanying metadata in a JSON file
4. Saves a copy at a standard location for the next steps in the pipeline
5. Provides a preview of the extracted content

The saved files serve as inputs for the next stage of the pipeline where
probability information will be added to create BayesDown.

DEPENDENCIES: os module for directory operations

```

OUTPUTS: Saved ArgDown files and preview of extracted content

"""

```
def save_extraction_results(results, output_directory="./outputs"):
 """
 Save the extraction results to file

 Args:
 results (dict): Extraction results from execute_extraction
 output_directory (str): Directory to save results

 Returns:
 str: Path to the saved ArgDown file
 """
 # Ensure output directory exists
 import os
 os.makedirs(output_directory, exist_ok=True)

 # Create base filename from source
 import os.path
 base_name = os.path.basename(results["source_path"]).split('.')[0]
 timestamp = time.strftime("%Y%m%d-%H%M%S")
 output_filename = f"{base_name}_argdown_{timestamp}.md"
 output_path = os.path.join(output_directory, output_filename)

 # Save the results
 save_argdown_extraction(results, output_path)

 print(f"Saved ArgDown extraction to: {output_path}")
 print(f"Metadata saved to: {output_path.replace('.md', '_metadata.json')}")

 # Also save to standard location for further processing
 standard_path = os.path.join(output_directory, "ArgDown.md")
 with open(standard_path, 'w') as f:
 f.write(results["argdown_content"])
 print(f"Also saved to standard location: {standard_path}")

 return output_path

Usage example:
output_path = save_extraction_results(extraction_results)
```

```
Preview the extracted ArgDown
from IPython.display import Markdown, display

Display the first 500 characters of the extracted ArgDown
preview = extraction_results["argdown_content"][:500] + "..." if len(extraction_results["argdown_content"]) > 500 else extraction_results["argdown_content"]
display(Markdown(f"## Extracted ArgDown Preview\n\n```\n{preview}\n```"))
```

```
NameError: name 'extraction_results' is not defined
```

```

NameError Traceback (most recent call last)
<ipython-input-57-84ee4ea64739> in <cell line: 0>()
 55
 56 # Usage example:
--> 57 output_path = save_extraction_results(extraction_results)
 58
 59 # Preview the extracted ArgDown
NameError: name 'extraction_results' is not defined
```

## 1.6 Review and Check ArgDown.md File

```
display(Markdown(md_content))
```

[Existential\_Catastrophe]: The destruction of humanity's long-term potential due to AI systems we've lost control over. {"instantiations": ["existential\_catastrophe\_TRUE", "existential\_catastrophe\_FALSE"]} - [Human\_Disempowerment]: Permanent and collective disempowerment of humanity relative to AI systems. {"instantiations": ["human\_disempowerment\_TRUE", "human\_disempowerment\_FALSE"]} - [Scale\_Of\_Power\_Seeking]: Power-seeking by AI systems scaling to the point of permanently disempowering all of humanity. {"instantiations": ["scale\_of\_power\_seeking\_TRUE", "scale\_of\_power\_seeking\_FALSE"]} - [Misaligned\_Power\_Seeking]: Deployed AI systems seeking power in unintended and high-impact ways due to problems with their objectives. {"instantiations": ["misaligned\_power\_seeking\_TRUE", "misaligned\_power\_seeking\_FALSE"]} - [APS\_Systems]: AI systems with advanced capabilities, agentic planning, and strategic awareness. {"instantiations": ["aps\_systems\_TRUE", "aps\_systems\_FALSE"]} - [Advanced\_AI\_Capability]: AI systems that outperform humans on tasks that grant significant power in the world. {"instantiations": ["advanced\_ai\_capability\_TRUE", "advanced\_ai\_capability\_FALSE"]} - [Agentic\_Planning]: AI systems making and executing plans based on world models to achieve objectives. {"instantiations": ["agentic\_planning\_TRUE", "agentic\_planning\_FALSE"]} - [Strategic\_Awareness]: AI systems with models accurately representing power dynamics with humans. {"instantiations": ["strategic\_awareness\_TRUE", "strategic\_awareness\_FALSE"]} - [Difficulty\_Of\_Alignment]: It is harder to build aligned systems than misaligned systems that are attractive to deploy. {"instantiations": ["difficulty\_of\_alignment\_TRUE", "difficulty\_of\_alignment\_FALSE"]} - [Instrumental\_Convergence]: AI systems with



misaligned objectives tend to seek power as an instrumental goal. {"instantiations": ["instrumental\_convergence\_TRUE", "instrumental\_convergence\_FALSE"]} - [Problems\_With\_Proxies]: Optimizing for proxy objectives breaks correlations with intended goals. {"instantiations": ["problems\_with\_proxies\_TRUE", "problems\_with\_proxies\_FALSE"]} - [Problems\_With\_Search]: Search processes can yield systems pursuing different objectives than intended. {"instantiations": ["problems\_with\_search\_TRUE", "problems\_with\_search\_FALSE"]} - [Deployment\_Decisions]: Decisions to deploy potentially misaligned AI systems. {"instantiations": ["deployment\_decisions\_DEPLOY", "deployment\_decisions\_WITHHOLD"]} - [Incentives\_To\_Build\_APS]: Strong incentives to build and deploy APS systems. {"instantiations": ["incentives\_to\_build\_aps\_STRONG", "incentives\_to\_build\_aps\_WEAK"]} - [Usefulness\_Of\_APS]: APS systems are very useful for many valuable tasks. {"instantiations": ["usefulness\_of\_aps\_HIGH", "usefulness\_of\_aps\_LOW"]} - [Competitive\_Dynamics]: Competitive pressures between AI developers. {"instantiations": ["competitive\_dynamics\_STRONG", "competitive\_dynamics\_WEAK"]} - [Deception\_By\_AI]: AI systems deceiving humans about their true objectives. {"instantiations": ["deception\_by\_ai\_TRUE", "deception\_by\_ai\_FALSE"]} - [Corrective\_Feedback]: Human society implementing corrections after observing problems. {"instantiations": ["corrective\_feedback\_EFFECTIVE", "corrective\_feedback\_INEFFECTIVE"]} - [Warning\_Shots]: Observable failures in weaker systems before catastrophic risks. {"instantiations": ["warning\_shots\_OBSERVED", "warning\_shots\_UNOBSERVED"]} - [Rapid\_Capability\_Escalation]: AI capabilities escalating very rapidly, allowing little time for correction. {"instantiations": ["rapid\_capability\_escalation\_TRUE", "rapid\_capability\_escalation\_FALSE"]} - [Barriers\_To\_Understanding]: Difficulty in understanding the internal workings of advanced AI systems. {"instantiations": ["barriers\_to\_understanding\_HIGH", "barriers\_to\_understanding\_LOW"]} - [Misaligned\_Power\_Seeking]: Deployed AI systems seeking power in unintended and high-impact ways due to problems with their objectives. {"instantiations": ["misaligned\_power\_seeking\_TRUE", "misaligned\_power\_seeking\_FALSE"]} - [Adversarial\_Dynamics]: Potentially adversarial relationships between humans and power-seeking AI. {"instantiations": ["adversarial\_dynamics\_TRUE", "adversarial\_dynamics\_FALSE"]} - [Misaligned\_Power\_Seeking]: Deployed AI systems seeking power in unintended and high-impact ways due to problems with their objectives. {"instantiations": ["misaligned\_power\_seeking\_TRUE", "misaligned\_power\_seeking\_FALSE"]} - [Stakes\_Of\_Error]: The escalating impact of mistakes with power-seeking AI systems. {"instantiations": ["stakes\_of\_error\_HIGH", "stakes\_of\_error\_LOW"]} - [Misaligned\_Power\_Seeking]: Deployed AI systems seeking power in unintended and high-impact ways due to problems with their objectives. {"instantiations": ["misaligned\_power\_seeking\_TRUE", "misaligned\_power\_seeking\_FALSE"]}

### 1.6.0 Check the Graph Structure with the ArgDown Sandbox Online

Copy and paste the BayesDown formatted ... in the ArgDown Sandbox below to quickly verify that the network renders correctly.

```
@title 1.6.1 --- ArgDown Online Sandbox --- [argdown_online_sandbox]

from IPython.display import IFrame

IFrame(src="https://argdown.org/sandbox/map/", width="100%", height="600px")
```

<IPython.lib.display.IFrame at 0x7b9ccf0ea210>

ArgDown Online Sandbox

## 1.7 Extract ArgDown Graph Information as DataFrame

Extract:

- > Nodes (Variable\_Title)
- > Edges (Parents)
- > Instantiations
- > Description

Implementation nodes: - One function for ArgDown and BayesDown extraction, but: - IF YOU ONLY WANT ARGDOWN EXTRACTION: USE ARGUMENT IN FUNCTION CALL “parse\_markdown\_hierarchy(markdown\_text, ArgDown = True)” - so if you set ArgDown = True, it gives you only instantiations, no probabilities.

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

```
"""
```

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

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

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

Key steps in the parsing process:

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

DEPENDENCIES: pandas, re, json libraries

```

INPUTS: Markdown text in ArgDown/BayesDown format
OUTPUTS: Structured DataFrame with node information, relationships, and properties
"""

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

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

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

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

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

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

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

 return df

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

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

 Returns:
 """

```

```

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

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

 Args:
 text (str): Cleaned markdown text

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

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

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

 title = title_match.group(1)

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

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

 # Split description and metadata at the first "{"
 if "{" in full_text:
 split_index = full_text.find("{")
 description = full_text[:split_index].strip()
 metadata = full_text[split_index:].strip()

```

```

else:
 # Keep the entire description and no metadata
 description = full_text
 metadata = '' # Initialize as empty string
else:
 description = ''
 metadata = '' # Ensure metadata is initialized

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

If neither symbol exists, indentation remains 0

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

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

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

 # Do NOT update metadata here - keep the original metadata
else:
 # First time seeing this title, create a new entry

```

```

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

 return titles_info

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

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

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

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

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

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

```

```

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

title = title_match.group(1)

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

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

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

line_number += 1

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

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

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

```

```

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

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

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

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

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

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

 j -= 1

return titles_info

```



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

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

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

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

```

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

```

```

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

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

return df

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

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

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

 for _, row in dataframe.iterrows():

```

```

 no_parent.append(not row['Parents']) # True if Parents list is empty
 no_children.append(not row['Children']) # True if Children list is empty

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

 return dataframe

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

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

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

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

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

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

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

```

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

return dataframe
```

```
example use case:
ex_csv = parse_markdown_hierarchy_fixed(md_content, ArgDown = True)
ex_csv
```

	Title	Description	line	line_number
0	Existential_Catastrophe	The destruction of humanity's long-term potent...	0	[0]
1	Human_Disempowerment	Permanent and collective disempowerment of hum...	1	[1]
2	Scale_Of_Power_Seeking	Power-seeking by AI systems scaling to the poi...	2	[2]
3	Misaligned_Power_Seeking	Deployed AI systems seeking power in unintende...	3	[3, 21, 23, 24]
4	APS_Systems	AI systems with advanced capabilities, agentic...	4	[4]
5	Advanced_AI_Capability	AI systems that outperform humans on tasks tha...	5	[5]
6	Agentic_Planning	AI systems making and executing plans based on...	6	[6]
7	Strategic_Awareness	AI systems with models accurately representing...	7	[7]
8	Difficulty_Of_Alignment	It is harder to build aligned systems than mis...	8	[8]
9	Instrumental_Convergence	AI systems with misaligned objectives tend to ...	9	[9]
10	Problems_With_Proxies	Optimizing for proxy objectives breaks correla...	10	[10]
11	Problems_With_Search	Search processes can yield systems pursuing di...	11	[11]
12	Deployment_Decisions	Decisions to deploy potentially misaligned AI ...	12	[12]
13	Incentives_To_Build_APS	Strong incentives to build and deploy APS syst...	13	[13]
14	Usefulness_Of_APS	APS systems are very useful for many valuable ...	14	[14]
15	Competitive_Dynamics	Competitive pressures between AI developers.	15	[15]
16	Deception_By_AI	AI systems deceiving humans about their true o...	16	[16]
17	Corrective_Feedback	Human society implementing corrections after o...	17	[17]
18	Warning_Shots	Observable failures in weaker systems before c...	18	[18]
19	Rapid_Capability_Escalation	AI capabilities escalating very rapidly, allow...	19	[19]
20	Barriers_To_Understanding	Difficulty in understanding the internal worki...	20	[20]
21	Adversarial_Dynamics	Potentially adversarial relationships between ...	22	[22]
22	Stakes_Of_Error	The escalating impact of mistakes with power-s...	24	[24]

example use case

## 1.8 Store ArgDown Information as 'ArgDown.csv' file

```
Assuming 'md_content' holds the markdown text
Store the results of running the function parse_markdown_hierarchy(md_content, ArgDown = True)
result_df = parse_markdown_hierarchy_fixed(md_content, ArgDown = True)
```

## 1.8 Store ArgDown Information as 'ArgDownSource' (PDF's of Papers) to ArgDown (.md file)

```
Save to CSV
result_df.to_csv('ArgDown.csv', index=False)

Test if 'ArgDown.csv' has been saved correctly with the correct information
Load the data from the CSV file
argdown_df = pd.read_csv('ArgDown.csv')

Display the DataFrame
print(argdown_df)
```

	Title \
0	Existential_Catastrophe
1	Human_Disempowerment
2	Scale_Of_Power_Seeking
3	Misaligned_Power_Seeking
4	APS_Systems
5	Advanced_AI_Capability
6	Agentic_Planning
7	Strategic_Awareness
8	Difficulty_Of_Alignment
9	Instrumental_Convergence
10	Problems_With_Proxies
11	Problems_With_Search
12	Deployment_Decisions
13	Incentives_To_Build_APS
14	Usefulness_Of_APS
15	Competitive_Dynamics
16	Deception_By_AI
17	Corrective_Feedback
18	Warning_Shots
19	Rapid_Capability_Escalation
20	Barriers_To_Understanding
21	Adversarial_Dynamics
22	Stakes_Of_Error

	Description	line	line_numbers \
0	The destruction of humanity's long-term potent...	0	[0]
1	Permanent and collective disempowerment of hum...	1	[1]
2	Power-seeking by AI systems scaling to the poi...	2	[2]
3	Deployed AI systems seeking power in unintende...	3	[3, 21, 23, 25]
4	AI systems with advanced capabilities, agentic...	4	[4]
5	AI systems that outperform humans on tasks tha...	5	[5]
6	AI systems making and executing plans based on...	6	[6]

7	AI systems with models accurately representing...	7	[7]
8	It is harder to build aligned systems than mis...	8	[8]
9	AI systems with misaligned objectives tend to ...	9	[9]
10	Optimizing for proxy objectives breaks correla...	10	[10]
11	Search processes can yield systems pursuing di...	11	[11]
12	Decisions to deploy potentially misaligned AI ...	12	[12]
13	Strong incentives to build and deploy APS syst...	13	[13]
14	APS systems are very useful for many valuable ...	14	[14]
15	Competitive pressures between AI developers.	15	[15]
16	AI systems deceiving humans about their true o...	16	[16]
17	Human society implementing corrections after o...	17	[17]
18	Observable failures in weaker systems before c...	18	[18]
19	AI capabilities escalating very rapidly, allow...	19	[19]
20	Difficulty in understanding the internal worki...	20	[20]
21	Potentially adversarial relationships between ...	22	[22]
22	The escalating impact of mistakes with power-s...	24	[24]

	indentation	indentation_levels \
0	0	[0]
1	0	[0]
2	4	[4]
3	8	[8, 0, 0, 0]
4	12	[12]
5	16	[16]
6	16	[16]
7	16	[16]
8	12	[12]
9	16	[16]
10	16	[16]
11	16	[16]
12	12	[12]
13	16	[16]
14	20	[20]
15	20	[20]
16	16	[16]
17	8	[8]
18	12	[12]
19	12	[12]
20	0	[0]
21	0	[0]
22	0	[0]

```

 Parents \
0 []
1 ['Scale_Of_Power_Seeking']
2 ['Misaligned_Power_Seeking', 'Corrective_Feedb...
3 ['APS_Systems', 'Difficulty_Of_Alignment', 'De...
4 ['Advanced_AI_Capability', 'Agentic_Planning',...
5 []
6 []
7 []
8 ['Instrumental_Convergence', 'Problems_With_Pr...
9 []
10 []
11 []
12 ['Incentives_To_Build_APS', 'Deception_By_AI']
13 ['Usefulness_Of_APS', 'Competitive_Dynamics']
14 []
15 []
16 []
17 ['Warning_Shots', 'Rapid_Capability_Escalation']
18 []
19 []
20 []
21 []
22 []

```

```

 Children \
0 []
1 []
2 ['Human_Disempowerment']
3 ['Scale_Of_Power_Seeking']
4 ['Misaligned_Power_Seeking']
5 ['APS_Systems']
6 ['APS_Systems']
7 ['APS_Systems']
8 ['Misaligned_Power_Seeking']
9 ['Difficulty_Of_Alignment']
10 ['Difficulty_Of_Alignment']
11 ['Difficulty_Of_Alignment']
12 ['Misaligned_Power_Seeking']
13 ['Deployment_Decisions']
14 ['Incentives_To_Build_APS']
15 ['Incentives_To_Build_APS']

```



```

16 ['Deployment_Decisions']
17 ['Scale_Of_Power_Seeking']
18 ['Corrective_Feedback']
19 ['Corrective_Feedback']
20 []
21 []
22 []

```

	instantiations	No_Parent	No_Children	\
0	['existential_catastrophe_TRUE', 'existential_...	True	True	
1	['human_disempowerment_TRUE', 'human_disempowe...	False	True	
2	['scale_of_power_seeking_TRUE', 'scale_of_powe...	False	False	
3	['misaligned_power_seeking_TRUE', 'misaligned_...	False	False	
4	['aps_systems_TRUE', 'aps_systems_FALSE']	False	False	
5	['advanced_ai_capability_TRUE', 'advanced_ai_c...	True	False	
6	['agentic_planning_TRUE', 'agentic_planning_FA...	True	False	
7	['strategic_awareness_TRUE', 'strategic_aware...	True	False	
8	['difficulty_of_alignment_TRUE', 'difficulty_o...	False	False	
9	['instrumental_convergence_TRUE', 'instrumenta...	True	False	
10	['problems_with_proxies_TRUE', 'problems_with_...	True	False	
11	['problems_with_search_TRUE', 'problems_with_s...	True	False	
12	['deployment_decisions_DEPLOY', 'deployment_de...	False	False	
13	['incentives_to_build_aps_STRONG', 'incentives...	False	False	
14	['usefulness_of_aps_HIGH', 'usefulness_of_aps_...	True	False	
15	['competitive_dynamics_STRONG', 'competitive_d...	True	False	
16	['deception_by_ai_TRUE', 'deception_by_ai_FALSE']	True	False	
17	['corrective_feedback_EFFECTIVE', 'corrective_...	False	False	
18	['warning_shots_OBSERVED', 'warning_shots_UNOB...	True	False	
19	['rapid_capability_escalation_TRUE', 'rapid_ca...	True	False	
20	['barriers_to_understanding_HIGH', 'barriers_t...	True	True	
21	['adversarial_dynamics_TRUE', 'adversarial_dyn...	True	True	
22	['stakes_of_error_HIGH', 'stakes_of_error_LOW']	True	True	

```

 parent_instantiations
0 []
1 [['scale_of_power_seeking_TRUE', 'scale_of_pow...
2 [['misaligned_power_seeking_TRUE', 'misaligned...
3 [['aps_systems_TRUE', 'aps_systems_FALSE'], ['...
4 [['advanced_ai_capability_TRUE', 'advanced_ai_...
5 []
6 []
7 []

```

```
8 [['instrumental_convergence_TRUE', 'instrument...
9 []
10 []
11 []
12 [['incentives_to_build_aps_STRONG', 'incentive...
13 [['usefulness_of_aps_HIGH', 'usefulness_of_aps...
14 []
15 []
16 []
17 [['warning_shots_OBSERVED', 'warning_shots_UNO...
18 []
19 []
20 []
21 []
22 []
```

## 2 Probability Extractions: ArgDown (.csv) to BayesDown (.md + plugin JSON syntax)

### 2.0 ArgDown to BayesDown: Adding Probability Information

#### Process Overview

This section implements the second major stage of the AMTAIR pipeline: enhancing the structured argument representation (ArgDown) with probability information to create BayesDown.

BayesDown extends ArgDown by adding: 1. Prior probabilities for each variable (unconditional beliefs) 2. Conditional probabilities representing the relationships between variables 3. The full parameter specification needed for a Bayesian network

The process follows these steps: 1. Generate probability questions for each node and its relationships 2. Create a BayesDown template with placeholders for these probabilities 3. Answer the probability questions (manually or via LLM) 4. Substitute the answers into the BayesDown representation

This enhanced representation contains all the information needed to construct a formal Bayesian network, enabling probabilistic reasoning and policy evaluation.

#### What is BayesDown?

BayesDown maintains the ArgDown structure but adds probability metadata:

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

The result is a hybrid representation that preserves the narrative structure of arguments while adding the mathematical precision of Bayesian networks.

## 2.1 Probability Extraction Questions — ‘ArgDown.csv’ to ‘ArgDown\_WithQuestions.csv’

```
@title 2.1.0 --- Probability Extraction Questions Generation --- [probability_extraction_c

"""
BLOCK PURPOSE: Generates probability questions for ArgDown nodes to prepare for BayesDown co

This block implements a key step in the pipeline where structure (from ArgDown)
is prepared for probability integration (to create BayesDown). It:

1. Processes a CSV file containing ArgDown structure
2. For each node, generates appropriate probability questions:
 - Prior probability questions for all nodes
 - Conditional probability questions for nodes with parents
3. Creates a new CSV file with these questions ready for the next stage

The generated questions serve as placeholders that will be answered in the
probability extraction phase to complete the Bayesian network.

DEPENDENCIES: pandas, json, itertools libraries
INPUTS: ArgDown CSV file
OUTPUTS: Enhanced CSV with probability questions for each node
"""

import pandas as pd
import re
import json
import itertools
from IPython.display import Markdown, display

def parse_instantiations(instantiations_str):
 """
 Parse instantiations from string or list format.
 Handles various input formats flexibly.

 Args:
 instantiations_str: Instantiations in string or list format

 Returns:
 list: Parsed instantiations as a list
 """
```

```

"""
if pd.isna(instantiations_str) or instantiations_str == '':
 return []

if isinstance(instantiations_str, list):
 return instantiations_str

try:
 # Try to parse as JSON
 return json.loads(instantiations_str)
except:
 # Try to parse as string list
 if isinstance(instantiations_str, str):
 # Remove brackets and split by comma
 clean_str = instantiations_str.strip('[]"\')
 if not clean_str:
 return []
 return [s.strip(' "') for s in clean_str.split(',') if s.strip()]

return []

def parse_parents(parents_str):
 """
 Parse parents from string or list format.
 Handles various input formats flexibly.

 Args:
 parents_str: Parents in string or list format

 Returns:
 list: Parsed parents as a list
 """
 if pd.isna(parents_str) or parents_str == '':
 return []

 if isinstance(parents_str, list):
 return parents_str

 try:
 # Try to parse as JSON
 return json.loads(parents_str)
 except:

```

```

 # Try to parse as string list
 if isinstance(parents_str, str):
 # Remove brackets and split by comma
 clean_str = parents_str.strip('[]"\'')
 if not clean_str:
 return []
 return [s.strip(' "') for s in clean_str.split(',') if s.strip()]

 return []

def get_parent_instantiations(parent, df):
 """
 Get the instantiations for a parent node from the DataFrame.
 Returns default instantiations if not found.

 Args:
 parent (str): Parent node name
 df (DataFrame): DataFrame containing node information

 Returns:
 list: Instantiations for the parent node
 """
 parent_row = df[df['Title'] == parent]
 if parent_row.empty:
 return [f"{parent}_TRUE", f"{parent}_FALSE"]

 instantiations = parse_instantiations(parent_row.iloc[0]['instantiations'])
 if not instantiations:
 return [f"{parent}_TRUE", f"{parent}_FALSE"]

 return instantiations

def generate_instantiation_questions(title, instantiation, parents, df):
 """
 Generate questions for a specific instantiation of a node.

 Args:
 title (str): The title of the node
 instantiation (str): The specific instantiation (e.g., "title_TRUE")
 parents (list): List of parent nodes
 df (DataFrame): The full DataFrame for looking up parent instantiations

```

```

Returns:
 dict: Dictionary mapping questions to estimate keys
"""
questions = {}

Always generate a prior probability question, regardless of parents
prior_question = f"What is the probability for {title}={instantiation}?"
questions[prior_question] = 'prior' # Question is the key, 'prior' is the value

If no parents, return only the prior question
if not parents:
 return questions

For nodes with parents, generate conditional probability questions
Get all combinations of parent instantiations
parent_instantiations = []
for parent in parents:
 parent_insts = get_parent_instantiations(parent, df)
 parent_instantiations.append([(parent, inst) for inst in parent_insts])

Generate all combinations
all_combinations = list(itertools.product(*parent_instantiations))

Create conditional probability questions for each combination
and use questions as keys, estimate_i as values
for i, combination in enumerate(all_combinations):
 condition_str = ", ".join([f"{parent}={inst}" for parent, inst in combination])
 question = f"What is the probability for {title}={instantiation} if {condition_str}?"
 questions[question] = f'estimate_{i + 1}' # Question is the key,
 # estimate_i is the value

return questions

def generate_argdown_with_questions(argdown_csv_path, output_csv_path):
 """
 Generate probability questions based on the ArgDown CSV file and save
 to a new CSV file.

 Args:
 argdown_csv_path (str): Path to the input ArgDown CSV file
 output_csv_path (str): Path to save the output CSV file with questions

```

```

Returns:
 DataFrame: Enhanced DataFrame with probability questions

Raises:
 Exception: If CSV loading fails or required columns are missing
"""
print(f"Loading ArgDown CSV from {argdown_csv_path}...")

Load the ArgDown CSV file
try:
 df = pd.read_csv(argdown_csv_path)
 print(f"Successfully loaded CSV with {len(df)} rows.")
except Exception as e:
 raise Exception(f"Error loading ArgDown CSV: {e}")

Validate required columns
required_columns = ['Title', 'Parents', 'instantiations']
missing_columns = [col for col in required_columns if col not in df.columns]
if missing_columns:
 raise Exception(f"Missing required columns: {'', '.join(missing_columns)}")

Initialize columns for questions
df['Generate_Positive_Instantiation_Questions'] = None
df['Generate_Negative_Instantiation_Questions'] = None

print("Generating probability questions for each node...")

Process each row to generate questions
for idx, row in df.iterrows():
 title = row['Title']
 instantiations = parse_instantiations(row['instantiations'])
 parents = parse_parents(row['Parents'])

 if len(instantiations) < 2:
 # Default instantiations if not provided
 instantiations = [f"{title}_TRUE", f"{title}_FALSE"]

 # Generate positive instantiation questions
 positive_questions = generate_instantiation_questions(title, instantiations[0], parents)

 # Generate negative instantiation questions

```



```

negative_questions = generate_instantiation_questions(title, instantiations[1], parent_id)

Update the DataFrame
df.at[idx, 'Generate_Positive_Instantiation_Questions'] = json.dumps(positive_questions)
df.at[idx, 'Generate_Negative_Instantiation_Questions'] = json.dumps(negative_questions)

Save the enhanced DataFrame
df.to_csv(output_csv_path, index=False)
print(f"Generated questions saved to {output_csv_path}")

return df

Example usage:
df_with_questions = generate_argdown_with_questions("ArgDown.csv", "ArgDown_WithQuestions.csv")

```

Loading ArgDown CSV from ArgDown.csv...

Successfully loaded CSV with 23 rows.

Generating probability questions for each node...

Generated questions saved to ArgDown\_WithQuestions.csv

```

Load the data from the ArgDown_WithQuestions CSV file
argdown_with_questions_df = pd.read_csv('ArgDown_WithQuestions.csv')

Display the DataFrame
print(argdown_with_questions_df)
argdown_with_questions_df

```

	Title \
0	Existential_Catastrophe
1	Human_Disempowerment
2	Scale_Of_Power_Seeking
3	Misaligned_Power_Seeking
4	APS_Systems
5	Advanced_AI_Capability
6	Agentic_Planning
7	Strategic_Awareness
8	Difficulty_Of_Alignment
9	Instrumental_Convergence
10	Problems_With_Proxies
11	Problems_With_Search
12	Deployment_Decisions
13	Incentives_To_Build_APS
14	Usefulness_Of_APS

15 Competitive\_Dynamics  
 16 Deception\_By\_AI  
 17 Corrective\_Feedback  
 18 Warning\_Shots  
 19 Rapid\_Capability\_Escalation  
 20 Barriers\_To\_Understanding  
 21 Adversarial\_Dynamics  
 22 Stakes\_Of\_Error

	Description	line	line_numbers \
0	The destruction of humanity's long-term potent...	0	[0]
1	Permanent and collective disempowerment of hum...	1	[1]
2	Power-seeking by AI systems scaling to the poi...	2	[2]
3	Deployed AI systems seeking power in unintende...	3	[3, 21, 23, 25]
4	AI systems with advanced capabilities, agentic...	4	[4]
5	AI systems that outperform humans on tasks tha...	5	[5]
6	AI systems making and executing plans based on...	6	[6]
7	AI systems with models accurately representing...	7	[7]
8	It is harder to build aligned systems than mis...	8	[8]
9	AI systems with misaligned objectives tend to ...	9	[9]
10	Optimizing for proxy objectives breaks correla...	10	[10]
11	Search processes can yield systems pursuing di...	11	[11]
12	Decisions to deploy potentially misaligned AI ...	12	[12]
13	Strong incentives to build and deploy APS syst...	13	[13]
14	APS systems are very useful for many valuable ...	14	[14]
15	Competitive pressures between AI developers.	15	[15]
16	AI systems deceiving humans about their true o...	16	[16]
17	Human society implementing corrections after o...	17	[17]
18	Observable failures in weaker systems before c...	18	[18]
19	AI capabilities escalating very rapidly, allow...	19	[19]
20	Difficulty in understanding the internal worki...	20	[20]
21	Potentially adversarial relationships between ...	22	[22]
22	The escalating impact of mistakes with power-s...	24	[24]

	indentation	indentation_levels \
0	0	[0]
1	0	[0]
2	4	[4]
3	8	[8, 0, 0, 0]
4	12	[12]
5	16	[16]
6	16	[16]

7	16	[16]
8	12	[12]
9	16	[16]
10	16	[16]
11	16	[16]
12	12	[12]
13	16	[16]
14	20	[20]
15	20	[20]
16	16	[16]
17	8	[8]
18	12	[12]
19	12	[12]
20	0	[0]
21	0	[0]
22	0	[0]

	Parents \
0	[]
1	['Scale_Of_Power_Seeking']
2	['Misaligned_Power_Seeking', 'Corrective_Feedb...]
3	['APS_Systems', 'Difficulty_Of_Alignment', 'De...]
4	['Advanced_AI_Capability', 'Agentic_Planning',...]
5	[]
6	[]
7	[]
8	['Instrumental_Convergence', 'Problems_With_Pr...]
9	[]
10	[]
11	[]
12	['Incentives_To_Build_APS', 'Deception_By_AI']
13	['Usefulness_Of_APS', 'Competitive_Dynamics']
14	[]
15	[]
16	[]
17	['Warning_Shots', 'Rapid_Capability_Escalation']
18	[]
19	[]
20	[]
21	[]
22	[]

```

 Children \
0 []
1 []
2 ['Human_Disempowerment']
3 ['Scale_Of_Power_Seeking']
4 ['Misaligned_Power_Seeking']
5 ['APS_Systems']
6 ['APS_Systems']
7 ['APS_Systems']
8 ['Misaligned_Power_Seeking']
9 ['Difficulty_Of_Alignment']
10 ['Difficulty_Of_Alignment']
11 ['Difficulty_Of_Alignment']
12 ['Misaligned_Power_Seeking']
13 ['Deployment_Decisions']
14 ['Incentives_To_Build_APS']
15 ['Incentives_To_Build_APS']
16 ['Deployment_Decisions']
17 ['Scale_Of_Power_Seeking']
18 ['Corrective_Feedback']
19 ['Corrective_Feedback']
20 []
21 []
22 []

```

	instantiations	No_Parent	No_Children \
0	['existential_catastrophe_TRUE', 'existential_...	True	True
1	['human_disempowerment_TRUE', 'human_disempowe...	False	True
2	['scale_of_power_seeking_TRUE', 'scale_of_powe...	False	False
3	['misaligned_power_seeking_TRUE', 'misaligned_...	False	False
4	['aps_systems_TRUE', 'aps_systems_FALSE']	False	False
5	['advanced_ai_capability_TRUE', 'advanced_ai_c...	True	False
6	['agentic_planning_TRUE', 'agentic_planning_FA...	True	False
7	['strategic_awareness_TRUE', 'strategic_aware...	True	False
8	['difficulty_of_alignment_TRUE', 'difficulty_o...	False	False
9	['instrumental_convergence_TRUE', 'instrumenta...	True	False
10	['problems_with_proxies_TRUE', 'problems_with_...	True	False
11	['problems_with_search_TRUE', 'problems_with_s...	True	False
12	['deployment_decisions_DEPLOY', 'deployment_de...	False	False
13	['incentives_to_build_aps_STRONG', 'incentives...	False	False
14	['usefulness_of_aps_HIGH', 'usefulness_of_aps_...	True	False
15	['competitive_dynamics_STRONG', 'competitive_d...	True	False

16	['deception_by_ai_TRUE', 'deception_by_ai_FALSE']	True	False
17	['corrective_feedback_EFFECTIVE', 'corrective_...]	False	False
18	['warning_shots_OBSERVED', 'warning_shots_UNOB...]	True	False
19	['rapid_capability_escalation_TRUE', 'rapid_ca...]	True	False
20	['barriers_to_understanding_HIGH', 'barriers_t...]	True	True
21	['adversarial_dynamics_TRUE', 'adversarial_dyn...]	True	True
22	['stakes_of_error_HIGH', 'stakes_of_error_LOW']	True	True

```

 parent_instantiations \
0 []
1 [['scale_of_power_seeking_TRUE', 'scale_of_pow...
2 [['misaligned_power_seeking_TRUE', 'misaligned...
3 [['aps_systems_TRUE', 'aps_systems_FALSE'], ['...
4 [['advanced_ai_capability_TRUE', 'advanced_ai_...
5 []
6 []
7 []
8 [['instrumental_convergence_TRUE', 'instrument...
9 []
10 []
11 []
12 [['incentives_to_build_aps_STRONG', 'incentive...
13 [['usefulness_of_aps_HIGH', 'usefulness_of_aps...
14 []
15 []
16 []
17 [['warning_shots_OBSERVED', 'warning_shots_UNO...
18 []
19 []
20 []
21 []
22 []

```

```

Generate_Positive_Instantiation_Questions \
0 {"What is the probability for Existential_Cata...
1 {"What is the probability for Human_Disempower...
2 {"What is the probability for Scale_Of_Power_S...
3 {"What is the probability for Misaligned_Power...
4 {"What is the probability for APS_Systems=aps_...
5 {"What is the probability for Advanced_AI_Capa...
6 {"What is the probability for Agentic_Planning...
7 {"What is the probability for Strategic_Awaren...

```

```

8 {"What is the probability for Difficulty_Of_Al...
9 {"What is the probability for Instrumental_Con...
10 {"What is the probability for Problems_With_Pr...
11 {"What is the probability for Problems_With_Se...
12 {"What is the probability for Deployment_Decis...
13 {"What is the probability for Incentives_To_Bu...
14 {"What is the probability for Usefulness_Of_AP...
15 {"What is the probability for Competitive_Dyna...
16 {"What is the probability for Deception_By_AI=...
17 {"What is the probability for Corrective_Feedb...
18 {"What is the probability for Warning_Shots=wa...
19 {"What is the probability for Rapid_Capability...
20 {"What is the probability for Barriers_To_Unde...
21 {"What is the probability for Adversarial_Dyna...
22 {"What is the probability for Stakes_Of_Error=...

```

#### Generate\_Negative\_Instantiation\_Questions

```

0 {"What is the probability for Existential_Cata...
1 {"What is the probability for Human_Disempower...
2 {"What is the probability for Scale_Of_Power_S...
3 {"What is the probability for Misaligned_Power...
4 {"What is the probability for APS_Systems=aps_...
5 {"What is the probability for Advanced_AI_Capa...
6 {"What is the probability for Agentic_Planning...
7 {"What is the probability for Strategic_Awaren...
8 {"What is the probability for Difficulty_Of_Al...
9 {"What is the probability for Instrumental_Con...
10 {"What is the probability for Problems_With_Pr...
11 {"What is the probability for Problems_With_Se...
12 {"What is the probability for Deployment_Decis...
13 {"What is the probability for Incentives_To_Bu...
14 {"What is the probability for Usefulness_Of_AP...
15 {"What is the probability for Competitive_Dyna...
16 {"What is the probability for Deception_By_AI=...
17 {"What is the probability for Corrective_Feedb...
18 {"What is the probability for Warning_Shots=wa...
19 {"What is the probability for Rapid_Capability...
20 {"What is the probability for Barriers_To_Unde...
21 {"What is the probability for Adversarial_Dyna...
22 {"What is the probability for Stakes_Of_Error=...

```

	Title	Description	line	line_number
0	Existential_Catastrophe	The destruction of humanity's long-term potent...	0	[0]
1	Human_Disempowerment	Permanent and collective disempowerment of hum...	1	[1]
2	Scale_Of_Power_Seeking	Power-seeking by AI systems scaling to the poi...	2	[2]
3	Misaligned_Power_Seeking	Deployed AI systems seeking power in unintende...	3	[3, 21, 23, 24]
4	APS_Systems	AI systems with advanced capabilities, agentic...	4	[4]
5	Advanced_AI_Capability	AI systems that outperform humans on tasks tha...	5	[5]
6	Agentic_Planning	AI systems making and executing plans based on...	6	[6]
7	Strategic_Awareness	AI systems with models accurately representing...	7	[7]
8	Difficulty_Of_Alignment	It is harder to build aligned systems than mis...	8	[8]
9	Instrumental_Convergence	AI systems with misaligned objectives tend to ...	9	[9]
10	Problems_With_Proxies	Optimizing for proxy objectives breaks correla...	10	[10]
11	Problems_With_Search	Search processes can yield systems pursuing di...	11	[11]
12	Deployment_Decisions	Decisions to deploy potentially misaligned AI ...	12	[12]
13	Incentives_To_Build_APS	Strong incentives to build and deploy APS syst...	13	[13]
14	Usefulness_Of_APS	APS systems are very useful for many valuable ...	14	[14]
15	Competitive_Dynamics	Competitive pressures between AI developers.	15	[15]
16	Deception_By_AI	AI systems deceiving humans about their true o...	16	[16]
17	Corrective_Feedback	Human society implementing corrections after o...	17	[17]
18	Warning_Shots	Observable failures in weaker systems before c...	18	[18]
19	Rapid_Capability_Escalation	AI capabilities escalating very rapidly, allow...	19	[19]
20	Barriers_To_Understanding	Difficulty in understanding the internal worki...	20	[20]
21	Adversarial_Dynamics	Potentially adversarial relationships between ...	22	[22]
22	Stakes_Of_Error	The escalating impact of mistakes with power-s...	24	[24]

## 2.2 ‘ArgDown\_WithQuestions.csv’ to ‘BayesDownQuestions.md’

### 2.2 Save BayesDown Extraction Questions as ‘BayesDownQuestions.md’

```
@title 2.2.0 --- BayesDown Questions Generation --- [bayesdown_questions_generation]

"""
BLOCK PURPOSE: Transforms the ArgDown with questions into a BayesDown template with placeholders

This function creates a BayesDown representation with probability placeholders
based on the questions generated in the previous step. It:

1. Loads the CSV file with probability questions
2. Constructs a directed graph to represent the causal structure
3. Processes each node to create BayesDown syntax with probability placeholders
4. Optionally includes comments with the specific questions to be answered
5. Saves the result as a markdown file for the next stage of the pipeline
```

The output is a BayesDown template that can be used in the probability extraction phase, where the placeholders will be replaced with actual probability values.

DEPENDENCIES: networkx, pandas, json libraries

INPUTS: CSV file with ArgDown structure and probability questions

OUTPUTS: BayesDown markdown file with probability placeholders

"""

```
def extract_bayesdown_questions_fixed(argdown_with_questions_path, output_md_path, include_c
```

"""

Generate BayesDown syntax from the ArgDown\_WithQuestions CSV file with correct parent-child relationships.

Args:

argdown\_with\_questions\_path (str): Path to the CSV file with probability questions

output\_md\_path (str): Path to save the output BayesDown file

include\_questions\_as\_comments (bool, optional): Whether to include the original questions as comments. Defaults to True.

Returns:

str: The generated BayesDown content

Raises:

Exception: If CSV loading fails or required columns are missing

"""

```
print(f>Loading CSV from {argdown_with_questions_path}...")
```

```
Load the CSV file
```

```
try:
```

```
 df = pd.read_csv(argdown_with_questions_path)
```

```
 print(f>Successfully loaded CSV with {len(df)} rows.")
```

```
except Exception as e:
```

```
 raise Exception(f>Error loading CSV: {e}")
```

```
Validate required columns
```

```
required_columns = ['Title', 'Description', 'Parents', 'Children', 'instantiations']
```

```
missing_columns = [col for col in required_columns if col not in df.columns]
```

```
if missing_columns:
```

```
 raise Exception(f>Missing required columns: {'', '.join(missing_columns)}")
```



```

print("Generating BayesDown syntax with placeholder probabilities...")

Build a directed graph of nodes
G = nx.DiGraph()

Add nodes to the graph
for idx, row in df.iterrows():
 G.add_node(row['Title'], data=row)

Add edges to the graph based on parent-child relationships - CORRECTLY
for idx, row in df.iterrows():
 child = row['Title']

 # Parse parents and add edges
 parents = row['Parents']
 if isinstance(parents, str):
 # Handle string representation of list
 if parents.startswith('[') and parents.endswith(']'):
 parents = parents.strip('[]')
 if parents: # Check if not empty
 parent_list = [p.strip().strip('\\"') for p in parents.split(',')]
 for parent in parent_list:
 if parent in G.nodes():
 # In BayesDown: Parent (cause) -> Child (effect)
 G.add_edge(parent, child)
 elif isinstance(parents, list):
 # Handle actual list
 for parent in parents:
 if parent in G.nodes():
 G.add_edge(parent, child)

Function to safely parse JSON strings
def safe_parse_json(json_str):
 if pd.isna(json_str):
 return {}

 if isinstance(json_str, dict):
 return json_str

 try:
 return json.loads(json_str)
 except:

```

```

 return {}

Start building the BayesDown content
bayesdown_content = "" # Initialize as empty

if include_questions_as_comments:
 bayesdown_content = "# BayesDown Representation with Placeholder Probabilities\n\n"
 bayesdown_content += "/* This file contains BayesDown syntax with placeholder probabilities\n\n"
 bayesdown_content += " Replace the placeholders with actual probability values based on\n\n"
 bayesdown_content += " questions in the comments. */\n\n"

Get leaf nodes (nodes with no outgoing edges) - these are effects without children
leaf_nodes = [n for n in G.nodes() if G.out_degree(n) == 0]

Helper function to process a node and its parents recursively
def process_node(node, indent_level=0, processed_nodes=None):
 if processed_nodes is None:
 processed_nodes = set()

 # Create the indentation string
 indent = ' ' * (indent_level * 2)
 prefix = f"{indent}+ " if indent_level > 0 else ""

 # Get node data
 node_data = G.nodes[node]['data']
 title = node_data['Title']
 description = node_data['Description'] if not pd.isna(node_data['Description']) else ''

 # Parse instantiations from the row data
 instantiations = parse_instantiations_safely(node_data['instantiations'])

 # Build the node string
 node_output = ""

 # Add comments with questions if requested
 if include_questions_as_comments:
 # Add positive questions as comments
 if 'Generate_Positive_Instantiation_Questions' in node_data:
 positive_questions = safe_parse_json(node_data['Generate_Positive_Instantiation_Questions'])
 for question in positive_questions.keys():
 node_output += f"{indent}/* {question} */\n"

```

```

Add negative questions as comments
if 'Generate_Negative_Instantiation_Questions' in node_data:
 negative_questions = safe_parse_json(node_data['Generate_Negative_Instantiation_Questions'])
 for question in negative_questions.keys():
 node_output += f"{indent}/* {question} */\n"

Check if this node was already fully defined elsewhere
if node in processed_nodes:
 # Just add a reference to the node
 node_output += f"{prefix}[{title}]\n"
 return node_output

Mark this node as processed
processed_nodes.add(node)

Prepare the metadata JSON
metadata = {
 "instantiations": instantiations
}

Add priors with full questions as keys
priors = {}
if 'Generate_Positive_Instantiation_Questions' in node_data:
 positive_questions = safe_parse_json(node_data['Generate_Positive_Instantiation_Questions'])
 for question, estimate_key in positive_questions.items():
 if estimate_key == 'prior':
 priors[question] = "%?" # Default placeholder

if 'Generate_Negative_Instantiation_Questions' in node_data:
 negative_questions = safe_parse_json(node_data['Generate_Negative_Instantiation_Questions'])
 for question, estimate_key in negative_questions.items():
 if estimate_key == 'prior':
 priors[question] = "%?" # Default placeholder

metadata["priors"] = priors

Add posteriors with full questions as keys
parents = list(G.predecessors(node))
if parents:
 posteriors = {}
 if 'Generate_Positive_Instantiation_Questions' in node_data:
 positive_questions = safe_parse_json(node_data['Generate_Positive_Instantiation_Questions'])

```

```

 for question, estimate_key in positive_questions.items():
 if estimate_key.startswith('estimate_'):
 posteriors[question] = "?" # Default placeholder

 if 'Generate_Negative_Instantiation_Questions' in node_data:
 negative_questions = safe_parse_json(node_data['Generate_Negative_Instantiation_Questions'])
 for question, estimate_key in negative_questions.items():
 if estimate_key.startswith('estimate_'):
 posteriors[question] = "?" # Default placeholder

 metadata["posteriors"] = posteriors

Format the node with metadata
node_output += f"{prefix}[{title}]: {description} {json.dumps(metadata)}\n"

Process parent nodes
for parent in parents:
 if parent != node: # Avoid self-references
 parent_output = process_node(parent, indent_level + 1, processed_nodes)
 node_output += parent_output

return node_output

Helper function to parse instantiations safely
def parse_instantiations_safely(instantiations_data):
 if isinstance(instantiations_data, list):
 return instantiations_data if instantiations_data else [f"TRUE", f"FALSE"]

 if isinstance(instantiations_data, str):
 try:
 parsed = json.loads(instantiations_data)
 if isinstance(parsed, list):
 return parsed if parsed else [f"TRUE", f"FALSE"]
 except:
 if instantiations_data.startswith('[') and instantiations_data.endswith(']'):
 items = instantiations_data.strip('[]').split(',')
 result = [item.strip(' "') for item in items if item.strip()]
 return result if result else [f"TRUE", f"FALSE"]

 return [f"TRUE", f"FALSE"] # Default

Process each leaf node and its ancestors

```

```

for leaf in leaf_nodes:
 bayesdown_content += process_node(leaf)

Save the BayesDown content
with open(output_md_path, 'w') as f:
 f.write(bayesdown_content)

print(f"BayesDown Questions saved to {output_md_path}")
return bayesdown_content

Explicitly set the value of include_questions_as_comments
include_questions_as_comments=False # or False, depending on your needs

Get the markdown content
bayesdown_questions = extract_bayesdown_questions_fixed(
 "ArgDown_WithQuestions.csv",
 "BayesDownQuestions.md", include_questions_as_comments=include_questions_as_comments
)

Determine the output file path based on include_questions_as_comments
if include_questions_as_comments: # Assuming include_questions_as_comments is defined somewhere
 output_file_path = "FULL_BayesDownQuestions.md"
else:
 output_file_path = "BayesDownQuestions.md"

Save the markdown content to the appropriate file
with open(output_file_path, 'w') as f:
 f.write(md_content)

print(f"Markdown content saved to {output_file_path}")

```

Loading CSV from ArgDown\_WithQuestions.csv...

Successfully loaded CSV with 23 rows.

Generating BayesDown syntax with placeholder probabilities...

BayesDown Questions saved to BayesDownQuestions.md

Markdown content saved to BayesDownQuestions.md

```

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

```

```
Display a preview of the format
print("\nBayesDown Format Preview:")
print(bayesdown_questions[:50000] + "...\\n")
```

Loading CSV from ArgDown\_WithQuestions.csv...

Successfully loaded CSV with 23 rows.

Generating BayesDown syntax with placeholder probabilities...

BayesDown Questions saved to FULL\_BayesDownQuestions.md

BayesDown Format Preview:

# BayesDown Representation with Placeholder Probabilities

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

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

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

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

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

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

/\* What is the probability for Human\_Disempowerment=human\_disempowerment\_TRUE if Scale\_Of\_Power\_Seeking=scale\_of\_power\_seeking\_TRUE? \*/

/\* What is the probability for Human\_Disempowerment=human\_disempowerment\_TRUE if Scale\_Of\_Power\_Seeking=scale\_of\_power\_seeking\_FALSE? \*/

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

/\* What is the probability for Human\_Disempowerment=human\_disempowerment\_FALSE if Scale\_Of\_Power\_Seeking=scale\_of\_power\_seeking\_TRUE? \*/

/\* What is the probability for Human\_Disempowerment=human\_disempowerment\_FALSE if Scale\_Of\_Power\_Seeking=scale\_of\_power\_seeking\_FALSE? \*/

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

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

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

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

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

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

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

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

/\* What is the probability for Scale\_Of\_Power\_Seeking=scale\_of\_power\_seeking\_FALSE if Misaligned\_Power\_Seeking=misaligned\_power\_seeking\_FALSE? \*/

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

/\* What is the probability for Scale\_Of\_Power\_Seeking=scale\_of\_power\_seeking\_FALSE if Misaligned\_Power\_Seeking=misaligned\_power\_seeking\_FALSE? \*/

+ [Scale\_Of\_Power\_Seeking]: Power-seeking by AI systems scaling to the point of permanent loss of humanity

/\* What is the probability for Misaligned\_Power\_Seeking=misaligned\_power\_seeking\_TRUE? \*/

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

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

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

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

```

/* What is the probability for Misaligned_Power_Seeking=misaligned_power_seeking_TRUE if
/* What is the probability for Misaligned_Power_Seeking=misaligned_power_seeking_TRUE if
/* What is the probability for Misaligned_Power_Seeking=misaligned_power_seeking_TRUE if
/* What is the probability for Misaligned_Power_Seeking=misaligned_power_seeking_TRUE if
/* What is the probability for Misaligned_Power_Seeking=misaligned_power_seeking_FALSE?
/* What is the probability for Misaligned_Power_Seeking=misaligned_power_seeking_FALSE i
/* What is the probability for Misaligned_Power_Seeking=misaligned_power_seeking_FALSE i
/* What is the probability for Misaligned_Power_Seeking=misaligned_power_seeking_FALSE i
/* What is the probability for Misaligned_Power_Seeking=misaligned_power_seeking_FALSE i
/* What is the probability for Misaligned_Power_Seeking=misaligned_power_seeking_FALSE i
/* What is the probability for Misaligned_Power_Seeking=misaligned_power_seeking_FALSE i
/* What is the probability for Misaligned_Power_Seeking=misaligned_power_seeking_FALSE i
+ [Misaligned_Power_Seeking]: Deployed AI systems seeking power in unintended and high-i
/* What is the probability for APS_Systems=aps_systems_TRUE? */
/* What is the probability for APS_Systems=aps_systems_TRUE if Advanced_AI_Capability=
/* What is the probability for APS_Systems=aps_systems_TRUE if Advanced_AI_Capability=
/* What is the probability for APS_Systems=aps_systems_TRUE if Advanced_AI_Capability=
/* What is the probability for APS_Systems=aps_systems_TRUE if Advanced_AI_Capability=
/* What is the probability for APS_Systems=aps_systems_TRUE if Advanced_AI_Capability=
/* What is the probability for APS_Systems=aps_systems_TRUE if Advanced_AI_Capability=
/* What is the probability for APS_Systems=aps_systems_TRUE if Advanced_AI_Capability=
/* What is the probability for APS_Systems=aps_systems_TRUE if Advanced_AI_Capability=
/* What is the probability for APS_Systems=aps_systems_FALSE? */
/* What is the probability for APS_Systems=aps_systems_FALSE if Advanced_AI_Capability
/* What is the probability for APS_Systems=aps_systems_FALSE if Advanced_AI_Capability
/* What is the probability for APS_Systems=aps_systems_FALSE if Advanced_AI_Capability
/* What is the probability for APS_Systems=aps_systems_FALSE if Advanced_AI_Capability
/* What is the probability for APS_Systems=aps_systems_FALSE if Advanced_AI_Capability
/* What is the probability for APS_Systems=aps_systems_FALSE if Advanced_AI_Capability
/* What is the probability for APS_Systems=aps_systems_FALSE if Advanced_AI_Capability
/* What is the probability for APS_Systems=aps_systems_FALSE if Advanced_AI_Capability
+ [APS_Systems]: AI systems with advanced capabilities, agentic planning, and strategi
/* What is the probability for Advanced_AI_Capability=advanced_ai_capability_TRUE? */
/* What is the probability for Advanced_AI_Capability=advanced_ai_capability_FALSE?
+ [Advanced_AI_Capability]: AI systems that outperform humans on tasks that grant si
/* What is the probability for Agentic_Planning=agentic_planning_TRUE? */
/* What is the probability for Agentic_Planning=agentic_planning_FALSE? */
+ [Agentic_Planning]: AI systems making and executing plans based on world models to
/* What is the probability for Strategic_Awareness=strategic_awareness_TRUE? */
/* What is the probability for Strategic_Awareness=strategic_awareness_FALSE? */
+ [Strategic_Awareness]: AI systems with models accurately representing power dynami

```

```

/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE? */
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if Inc
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if Inc
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if Inc
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if Inc
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if Inc
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if Inc
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if Inc
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if Inc
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if Inc
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE?
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE if Inc
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE if Inc
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE if Inc
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE if Inc
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE if Inc
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE if Inc
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE if Inc
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE if Inc
+ [Difficulty_Of_Alignment]: It is harder to build aligned systems than misaligned sys
/* What is the probability for Instrumental_Convergence=instrumental_convergence_TRU
/* What is the probability for Instrumental_Convergence=instrumental_convergence_FAL
+ [Instrumental_Convergence]: AI systems with misaligned objectives tend to seek pow
/* What is the probability for Problems_With_Proxies=problems_with_proxies_TRUE? */
/* What is the probability for Problems_With_Proxies=problems_with_proxies_FALSE? */
+ [Problems_With_Proxies]: Optimizing for proxy objectives breaks correlations with
/* What is the probability for Problems_With_Search=problems_with_search_TRUE? */
/* What is the probability for Problems_With_Search=problems_with_search_FALSE? */
+ [Problems_With_Search]: Search processes can yield systems pursuing different obje
/* What is the probability for Deployment_Decisions=deployment_decisions_DEPLOY? */
/* What is the probability for Deployment_Decisions=deployment_decisions_DEPLOY if Inc
/* What is the probability for Deployment_Decisions=deployment_decisions_DEPLOY if Inc
/* What is the probability for Deployment_Decisions=deployment_decisions_DEPLOY if Inc
/* What is the probability for Deployment_Decisions=deployment_decisions_DEPLOY if Inc
/* What is the probability for Deployment_Decisions=deployment_decisions_WITHHOLD? */
/* What is the probability for Deployment_Decisions=deployment_decisions_WITHHOLD if I
/* What is the probability for Deployment_Decisions=deployment_decisions_WITHHOLD if I
/* What is the probability for Deployment_Decisions=deployment_decisions_WITHHOLD if I
/* What is the probability for Deployment_Decisions=deployment_decisions_WITHHOLD if I
+ [Deployment_Decisions]: Decisions to deploy potentially misaligned AI systems. {"ins
/* What is the probability for Incentives_To_Build_APS=incentives_to_build_aps_STRON
/* What is the probability for Incentives_To_Build_APS=incentives_to_build_aps_STRON
/* What is the probability for Incentives_To_Build_APS=incentives_to_build_aps_STRON

```







```

/* What is the probability for Misaligned_Power_Seeking=misaligned_power_seeking_TRUE if
/* What is the probability for Misaligned_Power_Seeking=misaligned_power_seeking_FALSE?
/* What is the probability for Misaligned_Power_Seeking=misaligned_power_seeking_FALSE i
/* What is the probability for Misaligned_Power_Seeking=misaligned_power_seeking_FALSE i
/* What is the probability for Misaligned_Power_Seeking=misaligned_power_seeking_FALSE i
/* What is the probability for Misaligned_Power_Seeking=misaligned_power_seeking_FALSE i
/* What is the probability for Misaligned_Power_Seeking=misaligned_power_seeking_FALSE i
/* What is the probability for Misaligned_Power_Seeking=misaligned_power_seeking_FALSE i
/* What is the probability for Misaligned_Power_Seeking=misaligned_power_seeking_FALSE i
+ [Misaligned_Power_Seeking]: Deployed AI systems seeking power in unintended and high-i
/* What is the probability for APS_Systems=aps_systems_TRUE? */
/* What is the probability for APS_Systems=aps_systems_TRUE if Advanced_AI_Capability=
/* What is the probability for APS_Systems=aps_systems_TRUE if Advanced_AI_Capability=
/* What is the probability for APS_Systems=aps_systems_TRUE if Advanced_AI_Capability=
/* What is the probability for APS_Systems=aps_systems_TRUE if Advanced_AI_Capability=
/* What is the probability for APS_Systems=aps_systems_TRUE if Advanced_AI_Capability=
/* What is the probability for APS_Systems=aps_systems_TRUE if Advanced_AI_Capability=
/* What is the probability for APS_Systems=aps_systems_TRUE if Advanced_AI_Capability=
/* What is the probability for APS_Systems=aps_systems_TRUE if Advanced_AI_Capability=
/* What is the probability for APS_Systems=aps_systems_FALSE? */
/* What is the probability for APS_Systems=aps_systems_FALSE if Advanced_AI_Capability
/* What is the probability for APS_Systems=aps_systems_FALSE if Advanced_AI_Capability
/* What is the probability for APS_Systems=aps_systems_FALSE if Advanced_AI_Capability
/* What is the probability for APS_Systems=aps_systems_FALSE if Advanced_AI_Capability
/* What is the probability for APS_Systems=aps_systems_FALSE if Advanced_AI_Capability
/* What is the probability for APS_Systems=aps_systems_FALSE if Advanced_AI_Capability
/* What is the probability for APS_Systems=aps_systems_FALSE if Advanced_AI_Capability
+ [APS_Systems]: AI systems with advanced capabilities, agentic planning, and strategi
/* What is the probability for Advanced_AI_Capability=advanced_ai_capability_TRUE? *
/* What is the probability for Advanced_AI_Capability=advanced_ai_capability_FALSE?
+ [Advanced_AI_Capability]: AI systems that outperform humans on tasks that grant si
/* What is the probability for Agentic_Planning=agentic_planning_TRUE? */
/* What is the probability for Agentic_Planning=agentic_planning_FALSE? */
+ [Agentic_Planning]: AI systems making and executing plans based on world models to
/* What is the probability for Strategic_Awareness=strategic_awareness_TRUE? */
/* What is the probability for Strategic_Awareness=strategic_awareness_FALSE? */
+ [Strategic_Awareness]: AI systems with models accurately representing power dynami
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE? *
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if

```

```

/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_TRUE if
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE?
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE i
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE i
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE i
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE i
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE i
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE i
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE i
/* What is the probability for Difficulty_Of_Alignment=difficulty_of_alignment_FALSE i
+ [Difficulty_Of_Alignment]: It is harder to build aligned systems than misaligned sys
/* What is the probability for Instrumental_Convergence=instrumental_convergence_TRU
/* What is the probability for Instrumental_Convergence=instrumental_convergence_FAL
+ [Instrumental_Convergence]: AI systems with misaligned objectives tend to seek pow
/* What is the probability for Problems_With_Proxies=problems_with_proxies_TRUE? */
/* What is the probability for Problems_With_Proxies=problems_with_proxies_FALSE? */
+ [Problems_With_Proxies]: Optimizing for proxy objectives breaks correlations with
/* What is the probability for Problems_With_Search=problems_with_search_TRUE? */
/* What is the probability for Problems_With_Search=problems_with_search_FALSE? */
+ [Problems_With_Search]: Search processes can yield systems pursuing different obje
/* What is the probability for Deployment_Decisions=deployment_decisions_DEPLOY? */
/* What is the probability for Deployment_Decisions=deployment_decisions_DEPLOY if Inc
/* What is the probability for Deployment_Decisions=deployment_decisions_DEPLOY if Inc
/* What is the probability for Deployment_Decisions=deployment_decisions_DEPLOY if Inc
/* What is the probability for Deployment_Decisions=deployment_decisions_DEPLOY if Inc
/* What is the probability for Deployment_Decisions=deployment_decisions_WITHHOLD? */
/* What is the probability for Deployment_Decisions=deployment_decisions_WITHHOLD if I
/* What is the probability for Deployment_Decisions=deployment_decisions_WITHHOLD if I
/* What is the probability for Deployment_Decisions=deployment_decisions_WITHHOLD if I
/* What is the probability for Deployment_Decisions=deployment_decisions_WITHHOLD if I
+ [Deployment_Decisions]: Decisions to deploy potentially misaligned AI systems. {"ins
/* What is the probability for Incentives_To_Build_APS=incentives_to_build_aps_STRON
/* What is the probability for Incentives_To_Build_APS=incentives_to_build_aps_STRON
/* What is the probability for Incentives_To_Build_APS=incentives_to_build_aps_STRON
/* What is the probability for Incentives_To_Build_APS=incentives_to_build_aps_STRON
/* What is the probability for Incentives_To_Build_APS=incentives_to_build_aps_STRON
/* What is the probability for Incentives_To_Build_APS=incentives_to_build_aps_WEAK?

```

```

/* What is the probability for Incentives_To_Build_APS=incentives_to_build_aps_WEAK? */
/* What is the probability for Incentives_To_Build_APS=incentives_to_build_aps_WEAK? */
/* What is the probability for Incentives_To_Build_APS=incentives_to_build_aps_WEAK? */
/* What is the probability for Incentives_To_Build_APS=incentives_to_build_aps_WEAK? */
+ [Incentives_To_Build_APS]: Strong incentives to build and deploy APS systems. {"instantia
/* What is the probability for Usefulness_Of_APS=usefulness_of_aps_HIGH? */
/* What is the probability for Usefulness_Of_APS=usefulness_of_aps_LOW? */
+ [Usefulness_Of_APS]: APS systems are very useful for many valuable tasks. {"inst
/* What is the probability for Competitive_Dynamics=competitive_dynamics_STRONG? */
/* What is the probability for Competitive_Dynamics=competitive_dynamics_WEAK? */
+ [Competitive_Dynamics]: Competitive pressures between AI developers. {"instantia
/* What is the probability for Deception_By_AI=deception_by_ai_TRUE? */
/* What is the probability for Deception_By_AI=deception_by_ai_FALSE? */
+ [Deception_By_AI]: AI systems deceiving humans about their true objectives. {"inst
/* What is the probability for Corrective_Feedback=corrective_feedback_EFFECTIVE? */
/* What is the probability for Corrective_Feedback=corrective_feedback_EFFECTIVE if Warr
/* What is the probability for Corrective_Feedback=corrective_feedback_EFFECTIVE if Warr
/* What is the probability for Corrective_Feedback=corrective_feedback_EFFECTIVE if Warr
/* What is the probability for Corrective_Feedback=corrective_feedback_EFFECTIVE if Warr
/* What is the probability for Corrective_Feedback=corrective_feedback_INEFFECTIVE? */
/* What is the probability for Corrective_Feedback=corrective_feedback_INEFFECTIVE if Wa
/* What is the probability for Corrective_Feedback=corrective_feedback_INEFFECTIVE if Wa
/* What is the probability for Corrective_Feedback=corrective_feedback_INEFFECTIVE if Wa
/* What is the probability for Corrective_Feedback=corrective_feedback_INEFFECTIVE if Wa
+ [Corrective_Feedback]: Human society implementing corrections after observing problems
/* What is the probability for Warning_Shots=warning_shots_OBSERVED? */
/* What is the probability for Warning_Shots=warning_shots_UNOBSERVED? */
+ [Warning_Shots]: Observable failures in weaker systems before catastrophic risks. {'
/* What is the probability for Rapid_Capability_Escalation=rapid_capability_escalation
/* What is the probability for Rapid_Capability_Escalation=rapid_capability_escalation
+ [Rapid_Capability_Escalation]: AI capabilities escalating very rapidly, allowing lit
/* What is the probability for Barriers_To_Understanding=barriers_to_understanding_HIGH? */
/* What is the probability for Barriers_To_Understanding=barriers_to_understanding_LOW? */
[Barriers_To_Understanding]: Difficulty in understanding the internal workings of advanced A
/* What is the probability for Adversarial_Dynamics=adversarial_dynamics_TRUE? */
/* What is the probability for Adversarial_Dynamics=adversarial_dynamics_FALSE? */
[Adversarial_Dynamics]: Potentially adversarial relationships between humans and power-seeki
/* What is the probability for Stakes_Of_Error=stakes_of_error_HIGH? */
/* What is the probability for Stakes_Of_Error=stakes_of_error_LOW? */
[Stakes_Of_Error]: The escalating impact of mistakes with power-seeking AI systems. {"instan

Generate BayesDown format
bayesdown_questions = extract_bayesdown_questions_fixed(

```

```

"ArgDown_WithQuestions.csv",
"BayesDownQuestions.md",
include_questions_as_comments=False
)

Display a preview of the format
print(

)
print(bayesdown_questions[:50000] + "...\\n")

```

Loading CSV from ArgDown\_WithQuestions.csv...

Successfully loaded CSV with 23 rows.

Generating BayesDown syntax with placeholder probabilities...

BayesDown Questions saved to BayesDownQuestions.md

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

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

+ [Scale\_Of\_Power\_Seeking]: Power-seeking by AI systems scaling to the point of permanent domination.

+ [Misaligned\_Power\_Seeking]: Deployed AI systems seeking power in unintended and high-risk ways.

+ [APS\_Systems]: AI systems with advanced capabilities, agentic planning, and strategic behavior.

+ [Advanced\_AI\_Capability]: AI systems that outperform humans on tasks that grant significant power.

+ [Agentic\_Planning]: AI systems making and executing plans based on world models to achieve goals.

+ [Strategic\_Awareness]: AI systems with models accurately representing power dynamics and human behavior.

+ [Difficulty\_Of\_Alignment]: It is harder to build aligned systems than misaligned systems.

+ [Instrumental\_Convergence]: AI systems with misaligned objectives tend to seek power as a means to their ends.

+ [Problems\_With\_Proxies]: Optimizing for proxy objectives breaks correlations with the actual goal.

+ [Problems\_With\_Search]: Search processes can yield systems pursuing different objectives than intended.

+ [Deployment\_Decisions]: Decisions to deploy potentially misaligned AI systems. {"instantiation": 0.5, "probability": 0.5}

+ [Incentives\_To\_Build\_APS]: Strong incentives to build and deploy APS systems. {"instantiation": 0.5, "probability": 0.5}

+ [Usefulness\_Of\_APS]: APS systems are very useful for many valuable tasks. {"instantiation": 0.5, "probability": 0.5}

+ [Competitive\_Dynamics]: Competitive pressures between AI developers. {"instantiation": 0.5, "probability": 0.5}

+ [Deception\_By\_AI]: AI systems deceiving humans about their true objectives. {"instantiation": 0.5, "probability": 0.5}

+ [Corrective\_Feedback]: Human society implementing corrections after observing problems. {"instantiation": 0.5, "probability": 0.5}

+ [Warning\_Shots]: Observable failures in weaker systems before catastrophic risks. {"instantiation": 0.5, "probability": 0.5}

+ [Rapid\_Capability\_Escalation]: AI capabilities escalating very rapidly, allowing little time for intervention. {"instantiation": 0.5, "probability": 0.5}

[Barriers\_To\_Understanding]: Difficulty in understanding the internal workings of advanced AI systems.

[Adversarial\_Dynamics]: Potentially adversarial relationships between humans and power-seeking AI systems.

[Stakes\_Of\_Error]: The escalating impact of mistakes with power-seeking AI systems. {"instantiation": 0.5, "probability": 0.5}

...

## 2.3 Generate BayesDown Probability Extraction Prompt

Generate 2nd Extraction Prompt for Probabilities based on the questions generated from the 'ArgDown.csv' extraction

### 2.3.0 BayesDown Format Specification

BayesDown extends ArgDown with probability data in a structured JSON format to represent Bayesian networks. This intermediate representation bridges the gap between natural language arguments and formal probabilistic models, preserving both narrative structure and quantitative relationships.

#### Core Structure

A BayesDown representation consists of:

1. **Nodes:** Variables or statements in brackets [Node\_Name] with descriptive text
2. **Relationships:** Hierarchical structure with indentation and + symbols
3. **Metadata:** JSON objects containing probability information:

```
{
 "instantiations": ["state_TRUE", "state_FALSE"], // Possible states of variable
 "priors": {
 "p(state_TRUE)": "0.7", // Unconditional probability of state_TRUE
 "p(state_FALSE)": "0.3" // Unconditional probability of state_FALSE
 },
 "posteriors": {
 "p(state_TRUE|condition1_TRUE,condition2_FALSE)": "0.9", // Conditional on parent state
 "p(state_TRUE|condition1_FALSE,condition2_TRUE)": "0.4" // Different parent configuration
 }
}
```

#### ### 2.3.1 Rain-Sprinkler-Lawn Example

```
[Grass_Wet]: Concentrated moisture on grass. {"instantiations": ["grass_wet_TRUE", "grass_wet_FALSE"],
"priors": {"p(grass_wet_TRUE)": "0.322", "p(grass_wet_FALSE)": "0.678"},
"posteriors": {"p(grass_wet_TRUE|sprinkler_TRUE,rain_TRUE)": "0.99",
"p(grass_wet_TRUE|sprinkler_TRUE,rain_FALSE)": "0.9",
"p(grass_wet_TRUE|sprinkler_FALSE,rain_TRUE)": "0.8",
"p(grass_wet_TRUE|sprinkler_FALSE,rain_FALSE)": "0.0"}}
+ [Rain]: Water falling from the sky. {"instantiations": ["rain_TRUE", "rain_FALSE"],
"priors": {"p(rain_TRUE)": "0.2", "p(rain_FALSE)": "0.8"}}
+ [Sprinkler]: Artificial watering system. {"instantiations": ["sprinkler_TRUE", "sprinkler_FALSE"],
"priors": {"p(sprinkler_TRUE)": "0.44838", "p(sprinkler_FALSE)": "0.55162"},
"posteriors": {"p(sprinkler_TRUE|rain_TRUE)": "0.01", "p(sprinkler_TRUE|rain_FALSE)": "0.4"}
+ [Rain]
```

In this example:

- + Grass\_Wet is the effect/outcome node
- + Rain and Sprinkler are parent nodes (causes)
- + Rain also influences Sprinkler (people tend not to use sprinklers when it's raining)

Role in AMTAIR

BayesDown serves as the critical intermediate representation in the AMTAIR extraction pipeline.

For full syntax details, see the BayesDownSyntax.md file in the repository.

### ### 2.3.2 Probability Extraction Process

The probability extraction pipeline follows these steps:

Identify variables and their possible states

Extract prior probability statements

Identify conditional relationships

Extract conditional probability statements

Format the data in BayesDown syntax

### ### 2.3.3 Implementation Steps

To extract probabilities and create BayesDown format:

Run the `extract_probabilities` function on ArgDown text

Process the results into a structured format

Validate the probability distributions (ensure they sum to 1)

Generate the enhanced BayesDown representation

### ### 2.3.4 Validation and Quality Control

The probability extraction process includes validation steps:

Ensuring coherent probability distributions

Checking for logical consistency in conditional relationships

Verifying that all required probability statements are present

Handling missing data with appropriate default values

## ## 2.4 Prepare 2nd API call

## ## 2.5 Make BayesDown Probability Extraction API Call



```
2.6 Save BayesDown with Probability Estimates (.csv)

2.7 Review & Verify BayesDown Probability Estimates

2.7.2 Check the Graph Structure with the ArgDown Sandbox Online
Copy and paste the BayesDown formatted ... in the ArgDown Sandbox below to quickly verify the

2.8 Extract BayesDown with Probability Estimates as Dataframe

3 Data Extraction: BayesDown (.md) to Database (.csv)

3.0 BayesDown to Structured Data: Network Construction

Extraction Pipeline Overview

This section implements the core extraction pipeline described in the AMTAIR project document.

1. **Input**: Text in BayesDown format (see Section 2.3.1)
2. **Parsing**: Extract nodes, relationships, and probability information
3. **Structuring**: Organize into a DataFrame with formal relationships
4. **Enhancement**: Add derived properties and network metrics
5. **Output**: Structured data ready for Bayesian network construction

Theoretical Foundation

This implementation follows the extraction algorithm outlined in the AMTAIR project description.

1. Get nodes: All premises and conclusions from the argument structure
2. Get edges: Parent-child relationships between nodes
3. Extract probability distributions: Prior and conditional probabilities
4. Calculate derived metrics: Network statistics and node classifications

The resulting structured data maintains the complete information needed to reconstruct the BayesDown.

Role in Thesis Research

This extraction pipeline represents a key contribution of the Master's thesis, demonstrating the feasibility of automated argument analysis.

The rain-sprinkler-lawn example serves as a simple but complete test case, demonstrating even the most complex cases.

3.0.0 ExtractBayesDown-Data_v1
Build data frame with extractable information from BayesDown
```

```

::: {.cell quarto-private-1='{ "key": "colab", "value": { "base_uri": "https://localhost:8080/", "h
``` {.python .cell-code}
# read sprinkler example -- Occam Colab Online
file_path_ex_rain = "https://raw.githubusercontent.com/SingularitySmith/AMTAIR_Prototype/main

# Use requests.get to fetch content from URL
response = requests.get(file_path_ex_rain)
response.raise_for_status() # Raise HTTPError for bad responses (4xx or 5xx)

# Read content from the response
md_content_ex_rain = response.text

md_content_ex_rain

```

```
'[Existential_Catastrophe]: The destruction of humanity\'s long-term potential due to AI sys
```

3.0.1 Test BayesDown Extraction

```
display(Markdown(md_content_ex_rain)) # view BayesDown file formatted as Markdown
```

[Existential_Catastrophe]: The destruction of humanity’s long-term potential due to AI systems we’ve lost control over. {“instantiations”: [“existential_catastrophe_TRUE”, “existential_catastrophe_FALSE”], “priors”: {“p(existential_catastrophe_TRUE)”: “0.05”, “p(existential_catastrophe_FALSE)”: “0.95”}, “posteriors”: {“p(existential_catastrophe_TRUE|human_disempowerment_FALSE)”: “0.95”, “p(existential_catastrophe_TRUE|human_disempowerment_TRUE)”: “0.0”, “p(existential_catastrophe_FALSE|human_disempowerment_FALSE)”: “0.05”, “p(existential_catastrophe_FALSE|human_disempowerment_TRUE)”: “1.0”}} - [Human_Disempowerment]: Permanent and collective disempowerment of humanity relative to AI systems. {“instantiations”: [“human_disempowerment_TRUE”, “human_disempowerment_FALSE”], “priors”: {“p(human_disempowerment_TRUE)”: “0.208”, “p(human_disempowerment_FALSE)”: “0.792”}, “posteriors”: {“p(human_disempowerment_TRUE|scale_of_power_seeking_FALSE)”: “1.0”, “p(human_disempowerment_TRUE|scale_of_power_seeking_TRUE)”: “0.0”, “p(human_disempowerment_FALSE|scale_of_power_seeking_FALSE)”: “0.0”, “p(human_disempowerment_FALSE|scale_of_power_seeking_TRUE)”: “1.0”}} - [Scale_Of_Power_Seeking]: Power-seeking by AI systems scaling to the point of permanently disempowering all of humanity. {“instantiations”: [“scale_of_power_seeking_TRUE”, “scale_of_power_seeking_FALSE”], “priors”: {“p(scale_of_power_seeking_TRUE)”: “0.208”, “p(scale_of_power_seeking_FALSE)”: “0.792”}, “posteriors”: {“p(scale_of_power_seeking_TRUE|misaligned_power_seeking_corrective_feedback_EFFECTIVE)”: “0.25”, “p(scale_of_power_seeking_TRUE|misaligned_power_seeking_corrective_feedback_INEFFECTIVE)”: “0.60”, “p(scale_of_power_seeking_TRUE|misaligned_power_seeking_corrective_feedback_EFFECTIVE)”: “0.0”, “p(scale_of_power_seeking_TRUE|misaligned_power_seeking_corrective_feedback_INEFFECTIVE)”: “0.0”, “p(scale_of_power_seeking_FALSE|misaligned_power_seeking_corrective_feedback_EFFECTIVE)”: “0.75”, “p(scale_of_power_seeking_FALSE|misaligned_power_seeking

corrective_feedback_INEFFECTIVE)": "0.40", "p(scale_of_power_seeking_FALSE|misaligned_power_seeking_FALSE|corrective_feedback_EFFECTIVE)": "1.0", "p(scale_of_power_seeking_FALSE|misaligned_power_seeking_FALSE|corrective_feedback_INEFFECTIVE)": "1.0"} - [Misaligned_Power_Seeking]: Deployed AI systems seeking power in unintended and high-impact ways due to problems with their objectives. {"instantiations": ["misaligned_power_seeking_TRUE", "misaligned_power_seeking_FALSE"], "priors": {"p(misaligned_power_seeking_TRUE)": "0.338", "p(misaligned_power_seeking_FALSE)": "0.662"}, "posteriors": {"p(misaligned_power_seeking_TRUE|difficulty_of_alignment_TRUE, deployment_decisions_DEPLOY)": "0.90", "p(misaligned_power_seeking_TRUE|difficulty_of_alignment_TRUE, deployment_decisions_WITHHOLD)": "0.10", "p(misaligned_power_seeking_FALSE|difficulty_of_alignment_FALSE, deployment_decisions_DEPLOY)": "0.25", "p(misaligned_power_seeking_FALSE|difficulty_of_alignment_FALSE, deployment_decisions_WITHHOLD)": "0.05", "p(misaligned_power_seeking_TRUE|difficulty_of_alignment_TRUE, deployment_decisions_DEPLOY)": "0.0", "p(misaligned_power_seeking_TRUE|difficulty_of_alignment_TRUE, deployment_decisions_WITHHOLD)": "0.0", "p(misaligned_power_seeking_FALSE|difficulty_of_alignment_FALSE, deployment_decisions_DEPLOY)": "0.0", "p(misaligned_power_seeking_FALSE|difficulty_of_alignment_FALSE, deployment_decisions_WITHHOLD)": "0.0", "p(misaligned_power_seeking_TRUE|difficulty_of_alignment_TRUE, deployment_decisions_DEPLOY)": "0.10", "p(misaligned_power_seeking_TRUE|difficulty_of_alignment_TRUE, deployment_decisions_WITHHOLD)": "0.90", "p(misaligned_power_seeking_FALSE|difficulty_of_alignment_FALSE, deployment_decisions_DEPLOY)": "0.75", "p(misaligned_power_seeking_FALSE|difficulty_of_alignment_FALSE, deployment_decisions_WITHHOLD)": "0.95", "p(misaligned_power_seeking_TRUE|difficulty_of_alignment_TRUE, deployment_decisions_DEPLOY)": "1.0", "p(misaligned_power_seeking_TRUE|difficulty_of_alignment_TRUE, deployment_decisions_WITHHOLD)": "1.0", "p(misaligned_power_seeking_FALSE|difficulty_of_alignment_FALSE, deployment_decisions_DEPLOY)": "1.0", "p(misaligned_power_seeking_FALSE|difficulty_of_alignment_FALSE, deployment_decisions_WITHHOLD)": "1.0"} - [APS_Systems]: AI systems with advanced capabilities, agentic planning, and strategic awareness. {"instantiations": ["aps_systems_TRUE", "aps_systems_FALSE"], "priors": {"p(aps_systems_TRUE)": "0.65", "p(aps_systems_FALSE)": "0.35"}, "posteriors": {"p(aps_systems_TRUE|advanced_ai_capability_TRUE, agentic_planning_TRUE, strategic_awareness_TRUE)": "1.0", "p(aps_systems_TRUE|advanced_ai_capability_TRUE, agentic_planning_TRUE, strategic_awareness_FALSE)": "0.0", "p(aps_systems_TRUE|advanced_ai_capability_TRUE, agentic_planning_FALSE, strategic_awareness_TRUE)": "0.0", "p(aps_systems_TRUE|advanced_ai_capability_TRUE, agentic_planning_FALSE, strategic_awareness_FALSE)": "0.0", "p(aps_systems_FALSE|advanced_ai_capability_TRUE, agentic_planning_TRUE, strategic_awareness_TRUE)": "0.0", "p(aps_systems_FALSE|advanced_ai_capability_TRUE, agentic_planning_TRUE, strategic_awareness_FALSE)": "0.0", "p(aps_systems_FALSE|advanced_ai_capability_TRUE, agentic_planning_FALSE, strategic_awareness_TRUE)": "0.0", "p(aps_systems_FALSE|advanced_ai_capability_TRUE, agentic_planning_FALSE, strategic_awareness_FALSE)": "0.0", "p(aps_systems_FALSE|advanced_ai_capability_FALSE, agentic_planning_TRUE, strategic_awareness_TRUE)": "0.0", "p(aps_systems_FALSE|advanced_ai_capability_FALSE, agentic_planning_TRUE, strategic_awareness_FALSE)": "1.0", "p(aps_systems_FALSE|advanced_ai_capability_FALSE, agentic_planning_FALSE, strategic_awareness_TRUE)": "1.0", "p(aps_systems_FALSE|advanced_ai_capability_FALSE, agentic_planning_FALSE, strategic_awareness_FALSE)": "1.0", "p(aps_systems_FALSE|advanced_ai_capability_FALSE, agentic_planning_TRUE, strategic_awareness_TRUE)": "1.0", "p(aps_systems_FALSE|advanced_ai_capability_FALSE, agentic_planning_TRUE, strategic_awareness_FALSE)": "1.0", "p(aps_systems_FALSE|advanced_ai_capability_FALSE, agentic_planning_FALSE, strategic_awareness_TRUE)": "1.0", "p(aps_systems_FALSE|advanced_ai_capability_FALSE, agentic_planning_FALSE, strategic_awareness_FALSE)": "1.0"} - [Advanced_AI_Capability]:

AI systems that outperform humans on tasks that grant significant power in the world. {"instantiations": ["advanced_ai_capability_TRUE", "advanced_ai_capability_FALSE"], "priors": {"p(advanced_ai_capability_TRUE)": "0.80", "p(advanced_ai_capability_FALSE)": "0.20"}} - [Agentic_Planning]: AI systems making and executing plans based on world models to achieve objectives. {"instantiations": ["agentic_planning_TRUE", "agentic_planning_FALSE"], "priors": {"p(agentic_planning_TRUE)": "0.85", "p(agentic_planning_FALSE)": "0.15"}} - [Strategic_Awareness]: AI systems with models accurately representing power dynamics with humans. {"instantiations": ["strategic_awareness_TRUE", "strategic_awareness_FALSE"], "priors": {"p(strategic_awareness_TRUE)": "0.75", "p(strategic_awareness_FALSE)": "0.25"}} - [Difficulty_Of_Alignment]: It is harder to build aligned systems than misaligned systems that are attractive to deploy. {"instantiations": ["difficulty_of_alignment_TRUE", "difficulty_of_alignment_FALSE"], "priors": {"p(difficulty_of_alignment_TRUE)": "0.40", "p(difficulty_of_alignment_FALSE)": "0.60"}, "posteriors": {"p(difficulty_of_alignment_TRUE|instrumental_problems_with_proxies_TRUE, problems_with_search_TRUE)": "0.85", "p(difficulty_of_alignment_TRUE|instrumental_problems_with_proxies_TRUE, problems_with_search_FALSE)": "0.70", "p(difficulty_of_alignment_TRUE|instrumental_problems_with_proxies_FALSE, problems_with_search_TRUE)": "0.60", "p(difficulty_of_alignment_TRUE|instrumental_problems_with_proxies_FALSE, problems_with_search_FALSE)": "0.40", "p(difficulty_of_alignment_FALSE|instrumental_problems_with_proxies_TRUE, problems_with_search_TRUE)": "0.55", "p(difficulty_of_alignment_FALSE|instrumental_problems_with_proxies_TRUE, problems_with_search_FALSE)": "0.40", "p(difficulty_of_alignment_FALSE|instrumental_problems_with_proxies_FALSE, problems_with_search_TRUE)": "0.30", "p(difficulty_of_alignment_FALSE|instrumental_problems_with_proxies_FALSE, problems_with_search_FALSE)": "0.10", "p(difficulty_of_alignment_FALSE|instrumental_problems_with_proxies_TRUE, problems_with_search_TRUE)": "0.15", "p(difficulty_of_alignment_FALSE|instrumental_problems_with_proxies_TRUE, problems_with_search_FALSE)": "0.30", "p(difficulty_of_alignment_FALSE|instrumental_problems_with_proxies_FALSE, problems_with_search_TRUE)": "0.40", "p(difficulty_of_alignment_FALSE|instrumental_problems_with_proxies_FALSE, problems_with_search_FALSE)": "0.60", "p(difficulty_of_alignment_FALSE|instrumental_problems_with_proxies_TRUE, problems_with_search_TRUE)": "0.45", "p(difficulty_of_alignment_FALSE|instrumental_problems_with_proxies_TRUE, problems_with_search_FALSE)": "0.60", "p(difficulty_of_alignment_FALSE|instrumental_problems_with_proxies_FALSE, problems_with_search_TRUE)": "0.70", "p(difficulty_of_alignment_FALSE|instrumental_problems_with_proxies_FALSE, problems_with_search_FALSE)": "0.90"}} - [Instrumental_Convergence]: AI systems with misaligned objectives tend to seek power as an instrumental goal. {"instantiations": ["instrumental_convergence_TRUE", "instrumental_convergence_FALSE"], "priors": {"p(instrumental_convergence_TRUE)": "0.75", "p(instrumental_convergence_FALSE)": "0.25"}} - [Problems_With_Proxies]: Optimizing for proxy objectives breaks correlations with intended goals. {"instantiations": ["problems_with_proxies_TRUE", "problems_with_proxies_FALSE"], "priors": {"p(problems_with_proxies_TRUE)": "0.80", "p(problems_with_proxies_FALSE)": "0.20"}} - [Problems_With_Search]: Search processes can yield systems pursuing different objectives than intended. {"instantiations": ["problems_with_search_TRUE", "problems_with_search_FALSE"], "priors": {"p(problems_with_search_TRUE)": "0.70", "p(problems_with_search_FALSE)": "0.30"}} - [Deployment_Decisions]: Decisions to deploy potentially misaligned AI systems. {"instantiations": ["deployment_decisions_DEPLOY", "deployment_decisions_WITHHOLD"], "priors": {"p(deployment_decisions_DEPLOY)": "0.70",

“p(deployment_decisions_WITHHOLD): “0.30”, “posteriors”: {“p(deployment_decisions_DEPLOY|incentives_to_build_aps_STRONG, deception_by_ai_TRUE): “0.90”, “p(deployment_decisions_DEPLOY|incentives_to_build_aps_STRONG, deception_by_ai_FALSE): “0.75”, “p(deployment_decisions_DEPLOY|incentives_to_build_aps_WEAK, deception_by_ai_TRUE): “0.60”, “p(deployment_decisions_DEPLOY|incentives_to_build_aps_WEAK, deception_by_ai_FALSE): “0.30”, “p(deployment_decisions_WITHHOLD|incentives_to_build_aps_STRONG, deception_by_ai_TRUE): “0.10”, “p(deployment_decisions_WITHHOLD|incentives_to_build_aps_STRONG, deception_by_ai_FALSE): “0.25”, “p(deployment_decisions_WITHHOLD|incentives_to_build_aps_WEAK, deception_by_ai_TRUE): “0.40”, “p(deployment_decisions_WITHHOLD|incentives_to_build_aps_WEAK, deception_by_ai_FALSE): “0.70”}} - [Incentives_To_Build_APS]: Strong incentives to build and deploy APS systems. {“instantiations”: [“incentives_to_build_aps_STRONG”, “incentives_to_build_aps_WEAK”], “priors”: {“p(incentives_to_build_aps_STRONG): “0.80”, “p(incentives_to_build_aps_WEAK): “0.20”}, “posteriors”: {“p(incentives_to_build_aps_STRONG|usefulness_of_aps_HIGH, competitive_dynamics_STRONG): “0.95”, “p(incentives_to_build_aps_STRONG|usefulness_of_aps_HIGH, competitive_dynamics_WEAK): “0.80”, “p(incentives_to_build_aps_STRONG|usefulness_of_aps_LOW, competitive_dynamics_STRONG): “0.70”, “p(incentives_to_build_aps_STRONG|usefulness_of_aps_LOW, competitive_dynamics_WEAK): “0.30”, “p(incentives_to_build_aps_WEAK|usefulness_of_aps_HIGH, competitive_dynamics_STRONG): “0.05”, “p(incentives_to_build_aps_WEAK|usefulness_of_aps_HIGH, competitive_dynamics_WEAK): “0.20”, “p(incentives_to_build_aps_WEAK|usefulness_of_aps_LOW, competitive_dynamics_STRONG): “0.30”, “p(incentives_to_build_aps_WEAK|usefulness_of_aps_LOW, competitive_dynamics_WEAK): “0.70”}} - [Usefulness_Of_APS]: APS systems are very useful for many valuable tasks. {“instantiations”: [“usefulness_of_aps_HIGH”, “usefulness_of_aps_LOW”], “priors”: {“p(usefulness_of_aps_HIGH): “0.85”, “p(usefulness_of_aps_LOW): “0.15”}} - [Competitive_Dynamics]: Competitive pressures between AI developers. {“instantiations”: [“competitive_dynamics_STRONG”, “competitive_dynamics_WEAK”], “priors”: {“p(competitive_dynamics_STRONG): “0.75”, “p(competitive_dynamics_WEAK): “0.25”}} - [Deception_By_AI]: AI systems deceiving humans about their true objectives. {“instantiations”: [“deception_by_ai_TRUE”, “deception_by_ai_FALSE”], “priors”: {“p(deception_by_ai_TRUE): “0.50”, “p(deception_by_ai_FALSE): “0.50”}} - [Corrective_Feedback]: Human society implementing corrections after observing problems. {“instantiations”: [“corrective_feedback_EFFECTIVE”, “corrective_feedback_INEFFECTIVE”], “priors”: {“p(corrective_feedback_EFFECTIVE): “0.60”, “p(corrective_feedback_INEFFECTIVE): “0.40”}, “posteriors”: {“p(corrective_feedback_EFFECTIVE|warning_shots_OBSERVED, rapid_capability_escalation_TRUE): “0.40”, “p(corrective_feedback_EFFECTIVE|warning_shots_OBSERVED, rapid_capability_escalation_FALSE): “0.80”, “p(corrective_feedback_EFFECTIVE|warning_shots_UNOBSERVED, rapid_capability_escalation_TRUE): “0.15”, “p(corrective_feedback_EFFECTIVE|warning_shots_UNOBSERVED, rapid_capability_escalation_FALSE): “0.50”, “p(corrective_feedback_INEFFECTIVE|warning_shots_OBSERVED, rapid_capability_escalation_TRUE): “0.60”, “p(corrective_feedback_INEFFECTIVE|warning_shots_OBSERVED, rapid_capability_escalation_FALSE): “0.20”, “p(corrective_feedback_INEFFECTIVE|warning_shots_UNOBSERVED, rapid_capability_escalation_TRUE): “0.85”, “p(corrective_feedback_INEFFECTIVE|warning_shots_UNOBSERVED, rapid_capability_escalation_FALSE): “0.50”}} - [Warning_Shots]: Observable failures in weaker systems before catastrophic risks. {“instantiations”: [“warning_shots_OBSERVED”, “warning_shots_UNOBSERVED”], “priors”: {“p(warning_shots_OBSERVED): “0.70”,

`"p(warning_shots_UNOBSERVED)": "0.30"} - [Rapid_Capability_Escalation]: AI capabilities escalating very rapidly, allowing little time for correction. {"instantiations": ["rapid_capability_escalation_TRUE", "rapid_capability_escalation_FALSE"], "priors": {"p(rapid_capability_escalation_TRUE)": "0.45", "p(rapid_capability_escalation_FALSE)": "0.55"}} [Barriers_To_Understanding]: Difficulty in understanding the internal workings of advanced AI systems. {"instantiations": ["barriers_to_understanding_HIGH", "barriers_to_understanding_LOW"], "priors": {"p(barriers_to_understanding_HIGH)": "0.70", "p(barriers_to_understanding_LOW)": "0.30"}, "posteriors": {"p(barriers_to_understanding_HIGH|misaligned_power_seeking_TRUE)": "0.85", "p(barriers_to_understanding_HIGH|misaligned_power_seeking_FALSE)": "0.60", "p(barriers_to_understanding_LOW|misaligned_power_seeking_TRUE)": "0.15", "p(barriers_to_understanding_LOW|misaligned_power_seeking_FALSE)": "0.40"}} - [Misaligned_Power_Seeking]: Deployed AI systems seeking power in unintended and high-impact ways due to problems with their objectives. {"instantiations": ["misaligned_power_seeking_TRUE", "misaligned_power_seeking_FALSE"], "priors": {"p(misaligned_power_seeking_TRUE)": "0.338", "p(misaligned_power_seeking_FALSE)": "0.662"}} [Adversarial_Dynamics]: Potentially adversarial relationships between humans and power-seeking AI. {"instantiations": ["adversarial_dynamics_TRUE", "adversarial_dynamics_FALSE"], "priors": {"p(adversarial_dynamics_TRUE)": "0.60", "p(adversarial_dynamics_FALSE)": "0.40"}, "posteriors": {"p(adversarial_dynamics_TRUE|misaligned_power_seeking_TRUE)": "0.95", "p(adversarial_dynamics_TRUE|misaligned_power_seeking_FALSE)": "0.10", "p(adversarial_dynamics_FALSE|misaligned_power_seeking_TRUE)": "0.05", "p(adversarial_dynamics_FALSE|misaligned_power_seeking_FALSE)": "0.90"}} - [Misaligned_Power_Seeking]: Deployed AI systems seeking power in unintended and high-impact ways due to problems with their objectives. {"instantiations": ["misaligned_power_seeking_TRUE", "misaligned_power_seeking_FALSE"], "priors": {"p(misaligned_power_seeking_TRUE)": "0.338", "p(misaligned_power_seeking_FALSE)": "0.662"}} [Stakes_Of_Error]: The escalating impact of mistakes with power-seeking AI systems. {"instantiations": ["stakes_of_error_HIGH", "stakes_of_error_LOW"], "priors": {"p(stakes_of_error_HIGH)": "0.85", "p(stakes_of_error_LOW)": "0.15"}, "posteriors": {"p(stakes_of_error_HIGH|misaligned_power_seeking_TRUE)": "0.95", "p(stakes_of_error_HIGH|misaligned_power_seeking_FALSE)": "0.50", "p(stakes_of_error_LOW|misaligned_power_seeking_TRUE)": "0.05", "p(stakes_of_error_LOW|misaligned_power_seeking_FALSE)": "0.50"}} - [Misaligned_Power_Seeking]: Deployed AI systems seeking power in unintended and high-impact ways due to problems with their objectives. {"instantiations": ["misaligned_power_seeking_TRUE", "misaligned_power_seeking_FALSE"], "priors": {"p(misaligned_power_seeking_TRUE)": "0.338", "p(misaligned_power_seeking_FALSE)": "0.662"}}`

3.0.2 Check the Graph Structure with the ArgDown Sandbox Online

Copy and paste the BayesDown formatted ... in the ArgDown Sandbox below to quickly verify that the network renders correctly.

3.1 Extraction

BayesDown Extraction Code already part of ArgDown extraction code, therefore just use same function “`parse_markdown_hierarchy(markdown_data)`” and ignore the extra argument (“ArgDown”) because it is automatically set to false and will by default extract BayesDown.

```
result_df = parse_markdown_hierarchy_fixed(md_content_ex_rain)
result_df
```

	Title	Description	line	line_number
0	Existential_Catastrophe	The destruction of humanity's long-term potent...	0	[0]
1	Human_Disempowerment	Permanent and collective disempowerment of hum...	1	[1]
2	Scale_Of_Power_Seeking	Power-seeking by AI systems scaling to the poi...	2	[2]
3	Misaligned_Power_Seeking	Deployed AI systems seeking power in unintende...	3	[3, 21, 23, 24]
4	APS_Systems	AI systems with advanced capabilities, agentic...	4	[4]
5	Advanced_AI_Capability	AI systems that outperform humans on tasks tha...	5	[5]
6	Agentic_Planning	AI systems making and executing plans based on...	6	[6]
7	Strategic_Awareness	AI systems with models accurately representing...	7	[7]
8	Difficulty_Of_Alignment	It is harder to build aligned systems than mis...	8	[8]
9	Instrumental_Convergence	AI systems with misaligned objectives tend to ...	9	[9]
10	Problems_With_Proxies	Optimizing for proxy objectives breaks correla...	10	[10]
11	Problems_With_Search	Search processes can yield systems pursuing di...	11	[11]
12	Deployment_Decisions	Decisions to deploy potentially misaligned AI ...	12	[12]
13	Incentives_To_Build_APS	Strong incentives to build and deploy APS syst...	13	[13]
14	Usefulness_Of_APS	APS systems are very useful for many valuable ...	14	[14]
15	Competitive_Dynamics	Competitive pressures between AI developers.	15	[15]
16	Deception_By_AI	AI systems deceiving humans about their true o...	16	[16]
17	Corrective_Feedback	Human society implementing corrections after o...	17	[17]
18	Warning_Shots	Observable failures in weaker systems before c...	18	[18]
19	Rapid_Capability_Escalation	AI capabilities escalating very rapidly, allow...	19	[19]
20	Barriers_To_Understanding	Difficulty in understanding the internal worki...	20	[20]
21	Adversarial_Dynamics	Potentially adversarial relationships between ...	22	[22]
22	Stakes_Of_Error	The escalating impact of mistakes with power-s...	24	[24]

3.2 Data-Post-Processing

Add rows to data frame that can be calculated from the extracted rows

```
# @title 3.2.0 Data Post-Processing Functions --- [data_post_processing_functions]

"""
BLOCK PURPOSE: Enhances the extracted BayesDown data with calculated metrics and network pro
```

3.2 Data-Probability Extractions: ArgDown (.csv) to BayesDown (.md + plugin JSON syntax)

This block provides functions to enrich the basic extracted data with additional calculated columns that are useful for analysis and visualization:

1. Joint probabilities - Calculating $P(A,B)$ from conditional and prior probabilities
2. Network metrics - Centrality measures that indicate importance of nodes in the network
3. Markov blanket - Identifying the minimal set of nodes that shield a node from the rest

These enhancements provide valuable context for understanding the network structure and the relationships between variables, enabling more advanced analysis and improving visualization.

DEPENDENCIES: networkx for graph calculations

INPUTS: DataFrame with basic extracted BayesDown data

OUTPUTS: Enhanced DataFrame with additional calculated columns

"""

```
def enhance_extracted_data(df):
```

```
    """
```

```
    Enhance the extracted data with calculated columns
```

```
    Args:
```

```
        df: DataFrame with extracted BayesDown data
```

```
    Returns:
```

```
        Enhanced DataFrame with additional columns
```

```
    """
```

```
    # Create a copy to avoid modifying the original
```

```
    enhanced_df = df.copy()
```

```
    # 1. Calculate joint probabilities -  $P(A,B) = P(A|B) * P(B)$ 
```

```
    enhanced_df['joint_probabilities'] = None
```

```
    for idx, row in enhanced_df.iterrows():
```

```
        title = row['Title']
```

```
        priors = row['priors'] if isinstance(row['priors'], dict) else {}
```

```
        posteriors = row['posteriors'] if isinstance(row['posteriors'], dict) else {}
```

```
        parents = row['Parents'] if isinstance(row['Parents'], list) else []
```

```
    # Skip if no parents or no priors
```

```
    if not parents or not priors:
```

```
        continue
```



```

# Initialize joint probabilities dictionary
joint_probs = {}

# Get instantiations
instantiations = row['instantiations']
if not isinstance(instantiations, list) or not instantiations:
    continue

# For each parent and child instantiation combination, calculate joint probability
for inst in instantiations:
    # Get this instantiation's prior probability
    inst_prior_key = f"p({inst})"
    if inst_prior_key not in priors:
        continue

    try:
        inst_prior = float(priors[inst_prior_key])
    except (ValueError, TypeError):
        continue

    # For each parent
    for parent in parents:
        parent_row = enhanced_df[enhanced_df['Title'] == parent]
        if parent_row.empty:
            continue

        parent_insts = parent_row.iloc[0]['instantiations']
        if not isinstance(parent_insts, list) or not parent_insts:
            continue

        for parent_inst in parent_insts:
            # Get conditional probability
            cond_key = f"p({inst}|{parent}={parent_inst})"
            if cond_key in posteriors:
                try:
                    cond_prob = float(posteriors[cond_key])

                    # Get parent's prior
                    parent_priors = parent_row.iloc[0]['priors']
                    if not isinstance(parent_priors, dict):
                        continue

```

```

        parent_prior_key = f"p({parent_inst})"
        if parent_prior_key not in parent_priors:
            continue

        try:
            parent_prior = float(parent_priors[parent_prior_key])

            # Calculate joint probability:
            #  $P(A,B) = P(A|B) * P(B)$ 
            joint_prob = cond_prob * parent_prior
            joint_key = f"p({inst},{parent}={parent_inst})"
            joint_probs[joint_key] = str(round(joint_prob, 4))
        except (ValueError, TypeError):
            joint_prob = cond_prob * parent_prior
            joint_key = f"p({inst},{parent}={parent_inst})"
            joint_probs[joint_key] = str(round(joint_prob, 4))
        except (ValueError, TypeError):
            continue
    except (ValueError, TypeError):
        continue

# Store joint probabilities in dataframe
enhanced_df.at[idx, 'joint_probabilities'] = joint_probs

# 2. Calculate network metrics
# Create a directed graph
import networkx as nx
G = nx.DiGraph()

# Add nodes
for idx, row in enhanced_df.iterrows():
    G.add_node(row['Title'])

# Add edges
for idx, row in enhanced_df.iterrows():
    child = row['Title']
    parents = row['Parents'] if isinstance(row['Parents'], list) else []

    for parent in parents:
        if parent in G.nodes():
            G.add_edge(parent, child)

```

```

# Calculate centrality measures
degree centrality = nx.degree centrality(G) # Overall connectedness
in_degree centrality = nx.in_degree centrality(G) # How many nodes affect this one
out_degree centrality = nx.out_degree centrality(G) # How many nodes this one affects

try:
    betweenness centrality = nx.betweenness centrality(G) # Node's role as a connector
except:
    betweenness centrality = {node: 0 for node in G.nodes()}

# Add metrics to dataframe
enhanced_df['degree centrality'] = None
enhanced_df['in_degree centrality'] = None
enhanced_df['out_degree centrality'] = None
enhanced_df['betweenness centrality'] = None

for idx, row in enhanced_df.iterrows():
    title = row['Title']
    enhanced_df.at[idx, 'degree centrality'] = degree centrality.get(title, 0)
    enhanced_df.at[idx, 'in_degree centrality'] = in_degree centrality.get(title, 0)
    enhanced_df.at[idx, 'out_degree centrality'] = out_degree centrality.get(title, 0)
    enhanced_df.at[idx, 'betweenness centrality'] = betweenness centrality.get(title, 0)

# 3. Add Markov blanket information (parents, children, and children's parents)
enhanced_df['markov blanket'] = None

for idx, row in enhanced_df.iterrows():
    title = row['Title']
    parents = row['Parents'] if isinstance(row['Parents'], list) else []
    children = row['Children'] if isinstance(row['Children'], list) else []

    # Get children's parents (excluding this node)
    childrens_parents = []
    for child in children:
        child_row = enhanced_df[enhanced_df['Title'] == child]
        if not child_row.empty:
            child_parents = child_row.iloc[0]['Parents']
            if isinstance(child_parents, list):
                childrens_parents.extend([p for p in child_parents if p != title])

    # Remove duplicates
    childrens_parents = list(set(childrens_parents))

```

```
# Combine to get Markov blanket
markov_blanket = list(set(parents + children + childrens_parents))
enhanced_df.at[idx, 'markov_blanket'] = markov_blanket

return enhanced_df
```

```
# @title 3.2.1 --- Enhance Extracted Data with Network Metrics --- [enhance_extracted_data_w
```

```
"""
```

BLOCK PURPOSE: Applies the post-processing functions to enhance the extracted data.

This block takes the basic extracted DataFrame from the BayesDown parsing step and enriches it with calculated metrics that provide deeper insight into the network structure and relationships. It:

1. Applies the enhancement functions defined previously
2. Displays summary information about key calculated metrics
3. Saves the enhanced data for further analysis and visualization

The enhanced DataFrame provides a richer representation of the Bayesian network, including measures of node importance and conditional relationships that are essential for effective analysis and visualization.

DEPENDENCIES: enhance_extracted_data function

INPUTS: DataFrame with basic extracted BayesDown data

OUTPUTS: Enhanced DataFrame with additional calculated columns, saved to CSV

```
"""
```

```
# Enhance the extracted dataframe with calculated columns
```

```
enhanced_df = enhance_extracted_data(result_df)
```

```
# Display the enhanced dataframe
```

```
print("Enhanced DataFrame with additional calculated columns:")
```

```
enhanced_df.head()
```

```
# Check some calculated metrics
```

```
print("\nJoint Probabilities Example:")
```

```
example_node = enhanced_df.loc[0, 'Title']
```

```
joint_probs = enhanced_df.loc[0, 'joint_probabilities']
```

```
print(f"Joint probabilities for {example_node}:")
```

```
print(joint_probs)
```

```
print("\nNetwork Metrics:")
for idx, row in enhanced_df.iterrows():
    print(f"{row['Title']}:")
    print(f"    Degree Centrality: {row['degree_centrality']:.3f}")
    print(f"    Betweenness Centrality: {row['betweenness_centrality']:.3f}")

# Save the enhanced dataframe
enhanced_df.to_csv('enhanced_extracted_data.csv', index=False)
print("\nEnhanced data saved to 'enhanced_extracted_data.csv'")
```

Enhanced DataFrame with additional calculated columns:

Joint Probabilities Example:

Joint probabilities for Existential_Catastrophe:

None

Network Metrics:

Existential_Catastrophe:

Degree Centrality: 0.000

Betweenness Centrality: 0.000

Human_Disempowerment:

Degree Centrality: 0.045

Betweenness Centrality: 0.000

Scale_Of_Power_Seeking:

Degree Centrality: 0.136

Betweenness Centrality: 0.037

Misaligned_Power_Seeking:

Degree Centrality: 0.182

Betweenness Centrality: 0.056

APS_Systems:

Degree Centrality: 0.182

Betweenness Centrality: 0.019

Advanced_AI_Capability:

Degree Centrality: 0.045

Betweenness Centrality: 0.000

Agentic_Planning:

Degree Centrality: 0.045

Betweenness Centrality: 0.000

Strategic_Awareness:

Degree Centrality: 0.045

Betweenness Centrality: 0.000

Difficulty_Of_Alignment:

Degree Centrality: 0.182
Betweenness Centrality: 0.019
Instrumental_Convergence:
Degree Centrality: 0.045
Betweenness Centrality: 0.000
Problems_With_Proxies:
Degree Centrality: 0.045
Betweenness Centrality: 0.000
Problems_With_Search:
Degree Centrality: 0.045
Betweenness Centrality: 0.000
Deployment_Decisions:
Degree Centrality: 0.136
Betweenness Centrality: 0.026
Incentives_To_Build_APS:
Degree Centrality: 0.136
Betweenness Centrality: 0.017
Usefulness_Of_APS:
Degree Centrality: 0.045
Betweenness Centrality: 0.000
Competitive_Dynamics:
Degree Centrality: 0.045
Betweenness Centrality: 0.000
Deception_By_AI:
Degree Centrality: 0.045
Betweenness Centrality: 0.000
Corrective_Feedback:
Degree Centrality: 0.136
Betweenness Centrality: 0.009
Warning_Shots:
Degree Centrality: 0.045
Betweenness Centrality: 0.000
Rapid_Capability_Escalation:
Degree Centrality: 0.045
Betweenness Centrality: 0.000
Barriers_To_Understanding:
Degree Centrality: 0.000
Betweenness Centrality: 0.000
Adversarial_Dynamics:
Degree Centrality: 0.000
Betweenness Centrality: 0.000
Stakes_Of_Error:

Degree Centrality: 0.000

Betweenness Centrality: 0.000

Enhanced data saved to 'enhanced_extracted_data.csv'

3.4 Download and save finished data frame as .csv file

```
# @title 3.4.0 --- Save Extracted Data for Further Processing --- [save_extracted_data_for_further_processing]

"""
BLOCK PURPOSE: Saves the extracted data to a CSV file for further processing.

This step is essential for:
1. Persisting the structured representation of the Bayesian network
2. Enabling further analysis in other tools or notebook sections
3. Creating a permanent record of the extraction results
4. Making the data available for the visualization pipeline

The CSV format provides a standardized, tabular representation of the network
that can be easily loaded and processed in subsequent analysis steps.

DEPENDENCIES: pandas DataFrame operations
INPUTS: Extracted DataFrame from the parsing step
OUTPUTS: CSV file containing the structured network data
"""

# Save the extracted data as a CSV file
result_df.to_csv('extracted_data.csv', index=False)

print(" Extracted data saved successfully to 'extracted_data.csv'")
print("Note: If using updated data in future steps, "
      + "the file must be pushed to the GitHub repository")
```

Extracted data saved successfully to 'extracted_data.csv'

Note: If using updated data in future steps, the file must be pushed to the GitHub repository

4 Analysis & Inference: Bayesian Network Visualization

4.0 Bayesian Network Visualization Approach

This section implements the visualization component of the AMTAIR project, transforming the structured data extracted from BayesDown into an interactive network visualization that makes complex probabilistic relationships accessible to human understanding.

Visualization Philosophy

A key challenge in AI governance is making complex probabilistic relationships understandable to diverse stakeholders. This visualization system addresses this challenge through:

1. **Visual Encoding of Probability:** Node colors reflect probability values (green for high probability, red for low)
2. **Structural Classification:** Border colors indicate node types (blue for root causes, purple for intermediate nodes, magenta for leaf nodes)
3. **Progressive Disclosure:** Basic information in tooltips, detailed probability tables in modal popups
4. **Interactive Exploration:** Draggable nodes, configurable physics, click interactions

Connection to AMTAIR Goals

This visualization approach directly supports the AMTAIR project's goal of improving coordination in AI governance by:

1. Making implicit models explicit through visual representation
2. Providing a common language for discussing probabilistic relationships
3. Enabling non-technical stakeholders to engage with formal models
4. Creating shareable artifacts that facilitate collaboration

Implementation Structure

The visualization system is implemented in four phases:

1. **Network Construction:** Creating a directed graph representation using NetworkX

2. **Node Classification:** Identifying node types based on network position
3. **Visual Enhancement:** Adding color coding, tooltips, and interactive elements
4. **Interactive Features:** Implementing click handling for detailed exploration

The resulting visualization serves as both an analytical tool for experts and a communication tool for broader audiences, bridging the gap between technical and policy domains in AI governance discussions.

4.1 Phase 1: Dependencies/Functions

```
# @title 4.1.0 --- Bayesian Network Visualization Functions --- [bayesian_network_visualizat

"""

BLOCK PURPOSE: Provides functions to create interactive Bayesian network
visualizations from DataFrame representations of ArgDown/BayesDown data.

This block implements the visualization pipeline described in the AMTAIR project,
transforming the structured DataFrame extracted from ArgDown/BayesDown into an
interactive network graph that displays nodes, relationships, and probability
information. The visualization leverages NetworkX for graph representation and
PyVis for interactive display.

Key visualization features:
1. Color-coding of nodes based on probability values
2. Border styling to indicate node types (root, intermediate, leaf)
3. Interactive tooltips with probability information
4. Modal popups with detailed conditional probability tables
5. Physics-based layout for intuitive exploration

DEPENDENCIES: networkx, pyvis, HTML display from IPython
INPUTS: DataFrame with node information, relationships, and probabilities
OUTPUTS: Interactive HTML visualization of the Bayesian network
"""

from pyvis.network import Network
import networkx as nx
from IPython.display import HTML
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import io
import base64
import colorsys
```

```

import json

def create_bayesian_network_with_probabilities(df):
    """
    Create an interactive Bayesian network visualization with enhanced
    probability visualization and node classification based on network structure.

    Args:
        df (pandas.DataFrame): DataFrame containing node information, relationships, and probabilities

    Returns:
        IPython.display.HTML: Interactive HTML visualization of the Bayesian network
    """
    # PHASE 1: Create a directed graph representation
    G = nx.DiGraph()

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

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

        # Add node with base information
        G.add_node(
            title,
            description=description,
            priors=priors,
            instantiations=instantiations,
            posteriors=get_posteriors(row)
        )

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

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

```

```

G.add_edge(parent, child)

# PHASE 2: Classify nodes based on network structure
classify_nodes(G)

# PHASE 3: Create interactive network visualization
net = Network(notebook=True, directed=True, cdn_resources="in_line", height="600px", width="1000px")

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

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

# PHASE 4: Enhance node appearance with probability information
for node in net.nodes:
    node_id = node['id']
    node_data = G.nodes[node_id]

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

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

    # Get proper state names
    instantiations = node_data.get('instantiations', ["TRUE", "FALSE"])
    true_state = instantiations[0] if len(instantiations) > 0 else "TRUE"
    false_state = instantiations[1] if len(instantiations) > 1 else "FALSE"

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

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

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

```

```

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

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

# PHASE 5: Setup interactive click handling
# Prepare data for click handler
setup_data = {
    'nodes_data': {node_id: {
        'expanded_content': json.dumps(G.nodes[node_id].get('expanded_content', '')),
        'description': G.nodes[node_id].get('description', ''),
        'priors': G.nodes[node_id].get('priors', {}),
        'posteriors': G.nodes[node_id].get('posteriors', {})
    } for node_id in G.nodes()}
}

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

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

        if (nodeInfo) {
            // Create a modal popup for expanded content
            var modal = document.createElement('div');
            modal.style.position = 'fixed';

```

```

        modal.style.left = '50%';
        modal.style.top = '50%';
        modal.style.transform = 'translate(-50%, -50%)';
        modal.style.backgroundColor = 'white';
        modal.style.padding = '20px';
        modal.style.borderRadius = '5px';
        modal.style.boxShadow = '0 0 10px rgba(0,0,0,0.5)';
        modal.style.zIndex = '1000';
        modal.style.maxWidth = '80%';
        modal.style.maxHeight = '80%';
        modal.style.overflow = 'auto';

        // Add expanded content
        modal.innerHTML = nodeInfo.expanded_content || 'No detailed information available';

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

        // Add modal to body
        document.body.appendChild(modal);
    }
}

});
""" % json.dumps(setup_data['nodes_data'])

# PHASE 6: Save the graph to HTML and inject custom click handling
html_file = "bayesian_network.html"
net.save_graph(html_file)

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

```

```

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

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

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

```

4.2 Phase 2: Node Classification and Styling Module

```

# @title 4.2.0 --- Node Classification and Styling Functions --- [node_classification_and_styling]
"""
BLOCK PURPOSE: Implements the visual classification and styling of nodes in the Bayesian network.

This module handles the identification of node types based on their position in
the network and provides appropriate visual styling for each type.
The functions:

1. Classify nodes as parents (causes), children (intermediate effects), or leaves (final effects)
2. Assign appropriate border colors to visually distinguish node types
3. Calculate background colors based on probability values
4. Extract relevant information from DataFrame rows in a robust manner

The visual encoding helps users understand both the structure of the network
and the probability distributions at a glance.

DEPENDENCIES: colorsys for color manipulation
INPUTS: Graph structure and node data
OUTPUTS: Classification and styling information for visualization
"""

def classify_nodes(G):
    """
    Classify nodes as parent, child, or leaf based on network structure

    Args:

```

```

G (networkx.DiGraph): Directed graph representation of the Bayesian network

Effects:
    Adds 'node_type' attribute to each node in the graph:
    - 'parent': Root node with no parents but has children (causal source)
    - 'child': Node with both parents and children (intermediate)
    - 'leaf': Node with parents but no children (final effect)
    - 'isolated': Node with no connections (rare in Bayesian networks)
    """

for node in G.nodes():
    predecessors = list(G.predecessors(node)) # Nodes pointing to this one (causes)
    successors = list(G.successors(node))      # Nodes this one points to (effects)

    if not predecessors: # No parents
        if successors: # Has children
            G.nodes[node]['node_type'] = 'parent' # Root cause
        else: # No children either
            G.nodes[node]['node_type'] = 'isolated' # Disconnected node
    else: # Has parents
        if not successors: # No children
            G.nodes[node]['node_type'] = 'leaf' # Final effect
        else: # Has both parents and children
            G.nodes[node]['node_type'] = 'child' # Intermediate node

def get_border_color(node_type):
    """
    Return border color based on node type

    Args:
        node_type (str): Type of node ('parent', 'child', 'leaf', or 'isolated')

    Returns:
        str: Hex color code for node border
    """
    if node_type == 'parent':
        return '#0000FF' # Blue for root causes
    elif node_type == 'child':
        return '#800080' # Purple for intermediate nodes
    elif node_type == 'leaf':
        return '#FF00FF' # Magenta for final effects
    else:
        return '#000000' # Default black for any other type

```



```

def get_probability_color(priors):
    """
    Create background color based on probability (red to green gradient)

    Args:
        priors (dict): Dictionary containing probability information

    Returns:
        str: Hex color code for node background, ranging from red (low probability)
            to green (high probability)
    """
    # Default to neutral color if no probability
    if not priors or 'true_prob' not in priors:
        return '#F8F8F8' # Light grey

    # Get probability value
    prob = priors['true_prob']

    # Create color gradient from red (0.0) to green (1.0)
    hue = 120 * prob # 0 = red, 120 = green (in HSL color space)
    saturation = 0.75
    lightness = 0.8 # Lighter color for better text visibility

    # Convert HSL to RGB
    r, g, b = colorsys.hls_to_rgb(hue/360, lightness, saturation)

    # Convert to hex format
    hex_color = "#{:02x}{:02x}{:02x}".format(int(r*255), int(g*255), int(b*255))

    return hex_color

def get_parents(row):
    """
    Extract parent nodes from row data, with safe handling for different data types

    Args:
        row (pandas.Series): Row from DataFrame containing node information

    Returns:
        list: List of parent node names
    """

```

```

if 'Parents' not in row:
    return []

parents_data = row['Parents']

# Handle NaN, None, or empty list
if isinstance(parents_data, float) and pd.isna(parents_data):
    return []

if parents_data is None:
    return []

# Handle different data types
if isinstance(parents_data, list):
    # Return a list with NaN and empty strings removed
    return [p for p in parents_data if not (isinstance(p, float) and pd.isna(p)) and p != '']

if isinstance(parents_data, str):
    if not parents_data.strip():
        return []

    # Remove brackets and split by comma, removing empty strings and NaN
    cleaned = parents_data.strip('[]\''')
    if not cleaned:
        return []

    return [p.strip(' "') for p in cleaned.split(',') if p.strip()]

# Default: empty list
return []

def get_instantiations(row):
    """
    Extract instantiations with safe handling for different data types

    Args:
        row (pandas.Series): Row from DataFrame containing node information

    Returns:
        list: List of possible instantiations (states) for the node
    """
    if 'instantiations' not in row:

```

```

        return ["TRUE", "FALSE"]

    inst_data = row['instantiations']

    # Handle NaN or None
    if isinstance(inst_data, float) and pd.isna(inst_data):
        return ["TRUE", "FALSE"]

    if inst_data is None:
        return ["TRUE", "FALSE"]

    # Handle different data types
    if isinstance(inst_data, list):
        return inst_data if inst_data else ["TRUE", "FALSE"]

    if isinstance(inst_data, str):
        if not inst_data.strip():
            return ["TRUE", "FALSE"]

        # Remove brackets and split by comma
        cleaned = inst_data.strip('[]\''')
        if not cleaned:
            return ["TRUE", "FALSE"]

        return [i.strip(' "') for i in cleaned.split(',') if i.strip()]

    # Default
    return ["TRUE", "FALSE"]

def get_priors(row):
    """
    Extract prior probabilities with safe handling for different data types

    Args:
        row (pandas.Series): Row from DataFrame containing node information

    Returns:
        dict: Dictionary of prior probabilities with 'true_prob' added for convenience
    """
    if 'priors' not in row:
        return {}

```

```

priors_data = row['priors']

# Handle NaN or None
if isinstance(priors_data, float) and pd.isna(priors_data):
    return {}

if priors_data is None:
    return {}

result = {}

# Handle dictionary
if isinstance(priors_data, dict):
    result = priors_data
# Handle string representation of dictionary
elif isinstance(priors_data, str):
    if not priors_data.strip() or priors_data == '{}':
        return {}

    try:
        # Try to evaluate as Python literal
        import ast
        result = ast.literal_eval(priors_data)
    except:
        # Simple parsing for items like {'p(TRUE)': '0.2', 'p(FALSE)': '0.8'}
        if '{' in priors_data and '}' in priors_data:
            content = priors_data[priors_data.find('{')+1:priors_data.rfind('}')]
            items = [item.strip() for item in content.split(',')]

            for item in items:
                if ':' in item:
                    key, value = item.split(':', 1)
                    key = key.strip(' \\'')
                    value = value.strip(' \\'')
                    result[key] = value

# Extract main probability for TRUE state
instantiations = get_instantiations(row)
true_state = instantiations[0] if instantiations else "TRUE"
true_key = f"p({true_state})"

if true_key in result:

```

```

        try:
            result['true_prob'] = float(result[true_key])
        except:
            pass

    return result

def get_posteriors(row):
    """
    Extract posterior probabilities with safe handling for different data types

    Args:
        row (pandas.Series): Row from DataFrame containing node information

    Returns:
        dict: Dictionary of conditional probabilities
    """
    if 'posteriors' not in row:
        return {}

    posteriors_data = row['posteriors']

    # Handle NaN or None
    if isinstance(posteriors_data, float) and pd.isna(posteriors_data):
        return {}

    if posteriors_data is None:
        return {}

    result = {}

    # Handle dictionary
    if isinstance(posteriors_data, dict):
        result = posteriors_data
    # Handle string representation of dictionary
    elif isinstance(posteriors_data, str):
        if not posteriors_data.strip() or posteriors_data == '{}':
            return {}

    try:
        # Try to evaluate as Python literal
        import ast

```

```

        result = ast.literal_eval(posterior_data)
    except:
        # Simple parsing
        if '{' in posterior_data and '}' in posterior_data:
            content = posterior_data[posterior_data.find('{')+1:posterior_data.rfind('}')]
            items = [item.strip() for item in content.split(',')]

            for item in items:
                if ':' in item:
                    key, value = item.split(':', 1)
                    key = key.strip(' \\'')
                    value = value.strip(' \\'')
                    result[key] = value

    return result

```

4.3 Phase 3: HTML Content Generation Module

```

# @title 4.3.0 --- HTML Content Generation Functions --- [html_content_generation_functions]

"""
BLOCK PURPOSE: Creates rich HTML content for the interactive Bayesian network visualization.

This module generates the HTML components that enhance the Bayesian network
visualization:

1. Probability bars - Visual representation of probability distributions
2. Node tooltips - Rich information displayed on hover
3. Expanded content - Detailed probability information shown when clicking nodes

These HTML components make the mathematical concepts of Bayesian networks more
intuitive and accessible to users without requiring deep statistical knowledge.
The visual encoding of probabilities (colors, bars) and the progressive
disclosure of information (hover, click) help users build understanding at
their own pace.

DEPENDENCIES: HTML generation capabilities
INPUTS: Node data from the Bayesian network
OUTPUTS: HTML content for visualization components
"""

def create_probability_bar(true_prob, false_prob, height="15px", show_values=True, value_pre
    """

```

Creates a reusable HTML component to visualize probability distribution

Args:

true_prob (float): Probability of the true state (0.0-1.0)
 false_prob (float): Probability of the false state (0.0-1.0)
 height (str): CSS height of the bar
 show_values (bool): Whether to display numerical values
 value_prefix (str): Prefix to add before values (e.g., "p=")

Returns:

str: HTML for a horizontal bar showing probabilities

"""

Prepare display labels if showing values

true_label = f"{value_prefix}{true_prob:.3f}" if show_values else ""

false_label = f"{value_prefix}{false_prob:.3f}" if show_values else ""

Create the HTML for a horizontal stacked bar

html = f"""

<div style="width:100%; height:{height}; display:flex; border:1px solid #ccc; overflow:h

<div style="flex-basis:{true_prob*100}%; background:linear-gradient(to bottom, rgba

{true_l

</div>

<div style="flex-basis:{false_prob*100}%; background:linear-gradient(to bottom, rgba

{false_

</div>

</div>

"""

return html

def create_tooltip(node_id, node_data):

"""

Create rich HTML tooltip with probability information

Args:

node_id (str): Identifier of the node

node_data (dict): Node attributes including probabilities

Returns:

str: HTML content for tooltip displayed on hover

"""

Extract node information

description = node_data.get('description', '')

```

priors = node_data.get('priors', {})
instantiations = node_data.get('instantiations', ["TRUE", "FALSE"])

# Start building the HTML tooltip
html = f"""
<div style="max-width:350px; padding:10px; background-color:#f8f9fa; border-radius:5px;
    <h3 style="margin-top:0; color:#202124;">{node_id}</h3>
    <p style="font-style:italic;">{description}</p>
    """

# Add prior probabilities section
if priors and 'true_prob' in priors:
    true_prob = priors['true_prob']
    false_prob = 1.0 - true_prob

    # Get proper state names
    true_state = instantiations[0] if len(instantiations) > 0 else "TRUE"
    false_state = instantiations[1] if len(instantiations) > 1 else "FALSE"

    html += f"""
    <div style="margin-top:10px; background-color:#fff; padding:8px; border-radius:4px;
        <h4 style="margin-top:0; font-size:14px;">Prior Probabilities:</h4>
        <div style="display:flex; justify-content:space-between; margin-bottom:4px;">
            <div style="font-size:12px;">{true_state}: {true_prob:.3f}</div>
            <div style="font-size:12px;">{false_state}: {false_prob:.3f}</div>
        </div>
        {create_probability_bar(true_prob, false_prob, "20px", True)}
    </div>
    """

# Add click instruction
html += """
<div style="margin-top:8px; font-size:12px; color:#666; text-align:center;">
    Click node to see full probability details
</div>
</div>
    """

return html

def create_expanded_content(node_id, node_data):
    """

```


Create expanded content shown when a node is clicked

Args:

node_id (str): Identifier of the node
 node_data (dict): Node attributes including probabilities

Returns:

str: HTML content for detailed view displayed on click

```

"""
# Extract node information
description = node_data.get('description', '')
priors = node_data.get('priors', {})
posteriors = node_data.get('posteriors', {})
instantiations = node_data.get('instantiations', ["TRUE", "FALSE"])

# Get proper state names
true_state = instantiations[0] if len(instantiations) > 0 else "TRUE"
false_state = instantiations[1] if len(instantiations) > 1 else "FALSE"

# Extract probabilities
true_prob = priors.get('true_prob', 0.5)
false_prob = 1.0 - true_prob

# Start building the expanded content
html = f"""
<div style="max-width:500px; padding:15px; font-family:Arial, sans-serif;">
  <h2 style="margin-top:0; color:#333;">{node_id}</h2>
  <p style="font-style:italic; margin-bottom:15px;">{description}</p>

  <div style="margin-bottom:20px; padding:12px; border:1px solid #ddd; background-color:#f9f9f9;">
    <h3 style="margin-top:0; color:#333;">Prior Probabilities</h3>
    <div style="display:flex; justify-content:space-between; margin-bottom:5px;">
      <div><strong>{true_state}</strong> {true_prob:.3f}</div>
      <div><strong>{false_state}</strong> {false_prob:.3f}</div>
    </div>
    {create_probability_bar(true_prob, false_prob, "25px", True)}
  </div>
</div>
"""

# Add conditional probability table if available
if posteriors:
  html += """

```

```

<div style="padding:12px; border:1px solid #ddd; background-color:#f9f9f9; border-ra
    <h3 style="margin-top:0; color:#333;">Conditional Probabilities</h3>
    <table style="width:100%; border-collapse:collapse; font-size:13px;">
        <tr style="background-color:#eee;">
            <th style="padding:8px; text-align:left; border:1px solid #ddd;">Conditio
            <th style="padding:8px; text-align:center; border:1px solid #ddd; width:
            <th style="padding:8px; text-align:center; border:1px solid #ddd;">Visua
        </tr>
    </table>
    """

# Sort posteriors to group by similar conditions
posterior_items = list(posteriors.items())
posterior_items.sort(key=lambda x: x[0])

# Add rows for conditional probabilities
for key, value in posterior_items:
    try:
        # Try to parse probability value
        prob_value = float(value)
        inv_prob = 1.0 - prob_value

        # Add row with probability visualization
        html += f"""
        <tr>
            <td style="padding:8px; border:1px solid #ddd;">{key}</td>
            <td style="padding:8px; text-align:center; border:1px solid #ddd;">{prob
            <td style="padding:8px; border:1px solid #ddd;">
                {create_probability_bar(prob_value, inv_prob, "20px", False)}
            </td>
        </tr>
        """
    except:
        # Fallback for non-numeric values
        html += f"""
        <tr>
            <td style="padding:8px; border:1px solid #ddd;">{key}</td>
            <td style="padding:8px; text-align:center; border:1px solid #ddd;" colspan=
        </tr>
        """

html += """
</table>

```

```

    </div>
    """

    html += "</div>"

    return html

```

4.4 Phase 4: Main Visualization Function

```

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

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

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

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

        # Add node with base information
        G.add_node(
            title,
            description=description,
            priors=priors,
            instantiations=instantiations,
            posteriors=get_posteriors(row)
        )

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

```

```

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

# Classify nodes based on network structure
classify_nodes(G)

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

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

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

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

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

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

    # Get proper state names
    instantiations = node_data.get('instantiations', ["TRUE", "FALSE"])
    true_state = instantiations[0] if len(instantiations) > 0 else "TRUE"
    false_state = instantiations[1] if len(instantiations) > 1 else "FALSE"

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

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

```

```

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

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

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

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

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

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

        if (nodeInfo) {
            // Create a modal popup for expanded content

```

```

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

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

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

        // Add modal to body
        document.body.appendChild(modal);
    }
}

});
""" % json.dumps(setup_data['nodes_data'])

# Save the graph to HTML
html_file = "bayesian_network.html"

```

```
net.save_graph(html_file)

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

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

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

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


5 Quick check HTML Outputs

```
# @title 5.0 --- Quick check HTML Outputs--- [html_graph_visualization]

create_bayesian_network_with_probabilities(result_df)
```

Quick check HTML Outputs

```
# Use the function to create and display the visualization

print(result_df)
```

	Title \
0	Existential_Catastrophe
1	Human_Disempowerment
2	Scale_Of_Power_Seeking
3	Misaligned_Power_Seeking
4	APS_Systems
5	Advanced_AI_Capability
6	Agentic_Planning
7	Strategic_Awareness
8	Difficulty_Of_Alignment
9	Instrumental_Convergence
10	Problems_With_Proxies
11	Problems_With_Search
12	Deployment_Decisions
13	Incentives_To_Build_APS
14	Usefulness_Of_APS
15	Competitive_Dynamics
16	Deception_By_AI
17	Corrective_Feedback
18	Warning_Shots
19	Rapid_Capability_Escalation
20	Barriers_To_Understanding
21	Adversarial_Dynamics
22	Stakes_Of_Error

	Description	line	line_numbers \
0	The destruction of humanity's long-term potent...	0	[0]
1	Permanent and collective disempowerment of hum...	1	[1]
2	Power-seeking by AI systems scaling to the poi...	2	[2]
3	Deployed AI systems seeking power in unintende...	3	[3, 21, 23, 25]
4	AI systems with advanced capabilities, agentic...	4	[4]
5	AI systems that outperform humans on tasks tha...	5	[5]
6	AI systems making and executing plans based on...	6	[6]
7	AI systems with models accurately representing...	7	[7]
8	It is harder to build aligned systems than mis...	8	[8]
9	AI systems with misaligned objectives tend to ...	9	[9]
10	Optimizing for proxy objectives breaks correla...	10	[10]
11	Search processes can yield systems pursuing di...	11	[11]
12	Decisions to deploy potentially misaligned AI ...	12	[12]
13	Strong incentives to build and deploy APS syst...	13	[13]
14	APS systems are very useful for many valuable ...	14	[14]
15	Competitive pressures between AI developers.	15	[15]
16	AI systems deceiving humans about their true o...	16	[16]
17	Human society implementing corrections after o...	17	[17]
18	Observable failures in weaker systems before c...	18	[18]
19	AI capabilities escalating very rapidly, allow...	19	[19]
20	Difficulty in understanding the internal worki...	20	[20]
21	Potentially adversarial relationships between ...	22	[22]
22	The escalating impact of mistakes with power-s...	24	[24]

	indentation	indentation_levels \
0	0	[0]
1	0	[0]
2	4	[4]
3	8	[8, 0, 0, 0]
4	12	[12]
5	16	[16]
6	16	[16]
7	16	[16]
8	12	[12]
9	16	[16]
10	16	[16]
11	16	[16]
12	12	[12]
13	16	[16]
14	20	[20]

15	20	[20]
16	16	[16]
17	8	[8]
18	12	[12]
19	12	[12]
20	0	[0]
21	0	[0]
22	0	[0]

	Parents \
0	[]
1	[Scale_Of_Power_Seeking]
2	[Misaligned_Power_Seeking, Corrective_Feedback]
3	[APS_Systems, Difficulty_Of_Alignment, Deploym...
4	[Advanced_AI_Capability, Agentic_Planning, Str...
5	[]
6	[]
7	[]
8	[Instrumental_Convergence, Problems_With_Proxi...
9	[]
10	[]
11	[]
12	[Incentives_To_Build_APS, Deception_By_AI]
13	[Usefulness_Of_APS, Competitive_Dynamics]
14	[]
15	[]
16	[]
17	[Warning_Shots, Rapid_Capability_Escalation]
18	[]
19	[]
20	[]
21	[]
22	[]

	Children \
0	[]
1	[]
2	[Human_Disempowerment]
3	[Scale_Of_Power_Seeking]
4	[Misaligned_Power_Seeking]
5	[APS_Systems]
6	[APS_Systems]

```

7             [APS_Systems]
8  [Misaligned_Power_Seeking]
9    [Difficulty_Of_Alignment]
10   [Difficulty_Of_Alignment]
11   [Difficulty_Of_Alignment]
12  [Misaligned_Power_Seeking]
13    [Deployment_Decisions]
14   [Incentives_To_Build_APS]
15   [Incentives_To_Build_APS]
16    [Deployment_Decisions]
17   [Scale_Of_Power_Seeking]
18    [Corrective_Feedback]
19    [Corrective_Feedback]
20
21
22

```

```

                                instantiations \
0  [existential_catastrophe_TRUE, existential_cat...
1  [human_disempowerment_TRUE, human_disempowerme...
2  [scale_of_power_seeking_TRUE, scale_of_power_s...
3  [misaligned_power_seeking_TRUE, misaligned_pow...
4      [aps_systems_TRUE, aps_systems_FALSE]
5  [advanced_ai_capability_TRUE, advanced_ai_capa...
6      [agentic_planning_TRUE, agentic_planning_FALSE]
7  [strategic_awareness_TRUE, strategic_awareness...
8  [difficulty_of_alignment_TRUE, difficulty_of_a...
9  [instrumental_convergence_TRUE, instrumental_c...
10 [problems_with_proxies_TRUE, problems_with_pro...
11 [problems_with_search_TRUE, problems_with_sear...
12 [deployment_decisions_DEPLOY, deployment_decis...
13 [incentives_to_build_aps_STRONG, incentives_to...
14     [usefulness_of_aps_HIGH, usefulness_of_aps_LOW]
15 [competitive_dynamics_STRONG, competitive_dyna...
16     [deception_by_ai_TRUE, deception_by_ai_FALSE]
17 [corrective_feedback_EFFECTIVE, corrective_fee...
18 [warning_shots_OBSERVED, warning_shots_UNOBSER...
19 [rapid_capability_escalation_TRUE, rapid_capab...
20 [barriers_to_understanding_HIGH, barriers_to_u...
21 [adversarial_dynamics_TRUE, adversarial_dynami...
22     [stakes_of_error_HIGH, stakes_of_error_LOW]

```

```

                                priors \
0  {'p(existential_catastrophe_TRUE)': '0.05', 'p...
1  {'p(human_disempowerment_TRUE)': '0.208', 'p(h...
2  {'p(scale_of_power_seeking_TRUE)': '0.208', 'p...
3  {'p(misaligned_power_seeking_TRUE)': '0.338', ...
4  {'p(aps_systems_TRUE)': '0.65', 'p(aps_systems...
5  {'p(advanced_ai_capability_TRUE)': '0.80', 'p(...
6  {'p(agentive_planning_TRUE)': '0.85', 'p(agenti...
7  {'p(strategic_awareness_TRUE)': '0.75', 'p(str...
8  {'p(difficulty_of_alignment_TRUE)': '0.40', 'p...
9  {'p(instrumental_convergence_TRUE)': '0.75', '...
10 {'p(problems_with_proxies_TRUE)': '0.80', 'p(p...
11 {'p(problems_with_search_TRUE)': '0.70', 'p(pr...
12 {'p(deployment_decisions_DEPLOY)': '0.70', 'p(...
13 {'p(incentives_to_build_aps_STRONG)': '0.80', ...
14 {'p(usefulness_of_aps_HIGH)': '0.85', 'p(usefu...
15 {'p(competitive_dynamics_STRONG)': '0.75', 'p(...
16 {'p(deception_by_ai_TRUE)': '0.50', 'p(decepti...
17 {'p(corrective_feedback_EFFECTIVE)': '0.60', '...
18 {'p(warning_shots_OBSERVED)': '0.70', 'p(warni...
19 {'p(rapid_capability_escalation_TRUE)': '0.45'...
20 {'p(barriers_to_understanding_HIGH)': '0.70', ...
21 {'p(adversarial_dynamics_TRUE)': '0.60', 'p(ad...
22 {'p(stakes_of_error_HIGH)': '0.85', 'p(stakes_...

```

```

                                posteriors  No_Parent  No_Children \
0  {'p(existential_catastrophe_TRUE|human_disempo...      True      True
1  {'p(human_disempowerment_TRUE|scale_of_power_s...     False      True
2  {'p(scale_of_power_seeking_TRUE|misaligned_pow...     False     False
3  {'p(misaligned_power_seeking_TRUE|aps_systems_...     False     False
4  {'p(aps_systems_TRUE|advanced_ai_capability_TR...     False     False
5                                     {}           True     False
6                                     {}           True     False
7                                     {}           True     False
8  {'p(difficulty_of_alignment_TRUE|instrumental_...     False     False
9                                     {}           True     False
10                                    {}           True     False
11                                    {}           True     False
12 {'p(deployment_decisions_DEPLOY|incentives_to_...     False     False
13 {'p(incentives_to_build_aps_STRONG|usefulness_...     False     False
14                                    {}           True     False
15                                    {}           True     False

```

16	{}	True	False
17	{'p(corrective_feedback_EFFECTIVE warning_shot...	False	False
18	{}	True	False
19	{}	True	False
20	{'p(barriers_to_understanding_HIGH misaligned...	True	True
21	{'p(adversarial_dynamics_TRUE misaligned_power...	True	True
22	{'p(stakes_of_error_HIGH misaligned_power_seek...	True	True

```

                                parent_instantiations
0                                []
1  [[scale_of_power_seeking_TRUE, scale_of_power_...
2  [[misaligned_power_seeking_TRUE, misaligned_po...
3  [[aps_systems_TRUE, aps_systems_FALSE], [diffi...
4  [[advanced_ai_capability_TRUE, advanced_ai_cap...
5                                []
6                                []
7                                []
8  [[instrumental_convergence_TRUE, instrumental_...
9                                []
10                               []
11                               []
12  [[incentives_to_build_aps_STRONG, incentives_t...
13  [[usefulness_of_aps_HIGH, usefulness_of_aps_LO...
14                               []
15                               []
16                               []
17  [[warning_shots_OBSERVED, warning_shots_UNOBSE...
18                               []
19                               []
20                               []
21                               []
22                               []

```

```
# @title 5.1 --- File Import --- Load Files [file_import]
```

```

import requests
import io
from IPython.display import HTML, display

def load_and_display_html_from_github(repo_url, relative_path):
    """
    Loads an HTML file from a public GitHub repository and displays it.

```

```
Args:
    repo_url (str): The base URL of the GitHub repository (raw content).
    relative_path (str): The path to the HTML file relative to the repo_url.
"""
file_url = f"{repo_url}/{relative_path}"
print(f"Attempting to load HTML from: {file_url}")

try:
    # Fetch the file content from GitHub
    response = requests.get(file_url)

    # Check for successful response
    response.raise_for_status()

    # Read the content
    html_content = io.StringIO(response.text).read()

    print(f" Successfully loaded {relative_path}.")

    # Render the HTML content directly in the notebook
    display(HTML(html_content))

except requests.exceptions.RequestException as e:
    print(f" Error loading HTML file: {e}")
    print("Please check the URL and your internet connection.")
except Exception as e:
    print(f" An unexpected error occurred: {e}")

# Specify the base repository URL and the relative path to the HTML file
repo_base_url = "https://raw.githubusercontent.com/SingularitySmith/AMTAIR_Prototype/main/data"
html_relative_path = "runtime_created_data/bayesian_network.html"

# Load and display the HTML file
load_and_display_html_from_github(repo_base_url, html_relative_path)
```

Attempting to load HTML from: https://raw.githubusercontent.com/SingularitySmith/AMTAIR_Prototype/main/data/runtime_created_data/bayesian_network.html
Successfully loaded runtime_created_data/bayesian_network.html.

File Import

```
# @title 5.2 --- File Import --- Load Files [file_import2]

import requests
import io
```

```

from IPython.display import HTML, display

def load_and_display_html_from_github(repo_url, relative_path):
    """
    Loads an HTML file from a public GitHub repository and displays it.

    Args:
        repo_url (str): The base URL of the GitHub repository (raw content).
        relative_path (str): The path to the HTML file relative to the repo_url.
    """
    file_url = f"{repo_url}/{relative_path}"
    print(f"Attempting to load HTML from: {file_url}")

    try:
        # Fetch the file content from GitHub
        response = requests.get(file_url)

        # Check for successful response
        response.raise_for_status()

        # Read the content
        html_content = io.StringIO(response.text).read()

        print(f" Successfully loaded {relative_path}.")

        # Render the HTML content directly in the notebook
        display(HTML(html_content))

    except requests.exceptions.RequestException as e:
        print(f" Error loading HTML file: {e}")
        print("Please check the URL and your internet connection.")
    except Exception as e:
        print(f" An unexpected error occurred: {e}")

# Specify the base repository URL and the relative path to the HTML file
repo_base_url = "https://raw.githubusercontent.com/VJMeyer/submission/refs/heads/main/AMTAIF"
html_relative_path = "runtime_created_data/bayesian_network.html"

# Load and display the HTML file
load_and_display_html_from_github(repo_base_url, html_relative_path)

```

Attempting to load HTML from: https://raw.githubusercontent.com/VJMeyer/submission/refs/heads/main/AMTAIF/runtime_created_data/bayesian_network.html

Successfully loaded runtime_created_data/bayesian_network.html.

File Import 2

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

```
IFrame(src="https://singularitysmith.github.io/AMTAIR_Prototype/bayesian_network_carlsmith.h
```

```
<IPython.lib.display.IFrame at 0x7f04d69f0d90>
```

Dynamic Html Rendering of the Carlsmith Bayesian Network/DAG Visualization

Conclusion: From Prototype to Production

Summary of Achievements

This notebook has successfully demonstrated the core AMTAIR extraction pipeline, transforming structured argument representations into interactive Bayesian network visualizations through the following steps:

1. **Environment Setup:** Established a reproducible environment with necessary libraries and data access
2. **Argument Extraction:** Processed structured ArgDown representations preserving the hierarchical relationships
3. **Probability Integration:** Enhanced arguments with probability information to create BayesDown
4. **Data Transformation:** Converted BayesDown into structured DataFrame representation
5. **Visualization & Analysis:** Created interactive Bayesian network visualizations with probability encoding

The rain-sprinkler-lawn example, though simple, demonstrates all the key components of the extraction pipeline that can be applied to more complex AI safety arguments.

Limitations and Future Work

While this prototype successfully demonstrates the core pipeline, several limitations and opportunities for future work remain:

1. **LLM Extraction:** The current implementation focuses on processing pre-formatted ArgDown rather than performing extraction directly from unstructured text. Future work will integrate LLM-powered extraction.
2. **Scalability:** The system has been tested on small examples; scaling to larger, more complex arguments will require additional optimization and handling of computational complexity.
3. **Policy Evaluation:** The current implementation focuses on representation and visual-

ization; future work will add policy evaluation capabilities by implementing intervention modeling.

4. **Prediction Market Integration:** Future versions will integrate with forecasting platforms to incorporate live data into the models.

Connection to AMTAIR Project

This prototype represents just one component of the broader AMTAIR project described in the project documentation (see PY_AMTAIRDescription and PY_AMTAIR_SoftwareToolsNMilestones). The full project includes:

1. **AI Risk Pathway Analyzer (ARPA):** The core extraction and visualization system demonstrated in this notebook
2. **Worldview Comparator:** Tools for comparing different perspectives on AI risk
3. **Policy Impact Evaluator:** Systems for evaluating intervention effects across scenarios
4. **Strategic Intervention Generator:** Tools for identifying robust governance strategies

Together, these components aim to address the coordination crisis in AI governance by providing computational tools that make implicit models explicit, identify cruxes of disagreement, and evaluate policy impacts across diverse worldviews.

By transforming unstructured text into formal, analyzable representations, the AMTAIR project helps bridge the gaps between technical researchers, policy specialists, and other stakeholders, enabling more effective coordination in addressing existential risks from advanced AI.

6 Save Outputs

6.0 Saving and Exporting Results

This section provides tools for saving the notebook results and visualizations in various formats:

1. **HTML Export:** Creates a self-contained HTML version of the notebook with all visualizations
2. **Markdown Export:** Generates documentation-friendly Markdown version of the notebook
3. **PDF Export:** Creates a PDF document for formal sharing (requires LaTeX installation)

These exports are essential for: - Sharing analysis results with colleagues and stakeholders - Including visualizations in presentations and reports - Creating documentation for the AMTAIR project - Preserving results for future reference

The different formats serve different purposes, from interactive exploration (HTML) to documentation (Markdown) to formal presentation (PDF).

Instruction:

Download the ipynb, which you want to convert, on your local computer. Run the code below to upload the ipynb.

The html version will be downloaded automatically on your local machine. Enjoy it!

```
# @title 6.0.0 --- Save Visualization and Notebook Outputs as .HTML--- [save_visualization_a

"""
BLOCK PURPOSE: Provides tools for saving the notebook results in various formats.

This block offers functions to:
1. Convert the notebook to HTML for easy sharing and viewing
2. Convert the notebook to Markdown for documentation purposes
3. Save the visualization outputs for external use

These tools are essential for preserving the analysis results and making them
accessible outside the notebook environment, supporting knowledge transfer
and integration with other AMTAIR project components.
```

```

DEPENDENCIES: nbformat, nbconvert modules
INPUTS: Current notebook state
OUTPUTS: HTML, Markdown, or other format versions of the notebook
"""

import nbformat
from nbconvert import HTMLExporter
import os

# Repository URL variable for file access
repo_url = "https://raw.githubusercontent.com/SingularitySmith/AMTAIR_Prototype/main/data/example_notebook.ipynb"
notebook_name = "AMTAIR_Prototype_example_carlsmith" # Change when working with different example notebooks

# Download the notebook file
!wget {repo_url}{notebook_name}.ipynb -O {notebook_name}.ipynb

# Load the notebook
try:
    with open(f"{notebook_name}.ipynb") as f:
        nb = nbformat.read(f, as_version=4)
        print(f" Successfully loaded notebook: {notebook_name}.ipynb")
except FileNotFoundError:
    print(f" Error: File '{notebook_name}.ipynb' not found. Please check if it was downloaded correctly.")

# Initialize the HTML exporter
exporter = HTMLExporter()

# Convert the notebook to HTML
try:
    (body, resources) = exporter.from_notebook_node(nb)

    # Save the HTML to a file
    with open(f"{notebook_name}IPYNB.html", "w") as f:
        f.write(body)
    print(f" Successfully saved HTML version to: {notebook_name}IPYNB.html")
except Exception as e:
    print(f" Error converting notebook to HTML: {str(e)}")

```

```

--2025-05-24 20:09:38-- https://raw.githubusercontent.com/SingularitySmith/AMTAIR_Prototype/main/data/example_notebook.ipynb
Resolving raw.githubusercontent.com (raw.githubusercontent.com)... 185.199.108.133, 185.199.108.133, 185.199.108.133, 185.199.108.133
Connecting to raw.githubusercontent.com (raw.githubusercontent.com)|185.199.108.133|:443...
HTTP request sent, awaiting response... 200 OK

```

Length: 1689816 (1.6M) [text/plain]

Saving to: 'AMTAIR_Prototype_example_carlsmith.ipynb'

AMTAIR_Prototype_ex 100%[=====>] 1.61M 6.36MB/s in 0.3s

2025-05-24 20:09:38 (6.36 MB/s) - 'AMTAIR_Prototype_example_carlsmith.ipynb' saved [1689816/

Successfully loaded notebook: AMTAIR_Prototype_example_carlsmith.ipynb

Successfully saved HTML version to: AMTAIR_Prototype_example_carlsmithIPYNB.html

6.1 Convert .ipynb Notebook to Markdown

```
# @title 6.1.0 --- Convert .ipynb Notebook to Markdown --- [convert_notebook_to_markdown]

import nbformat
from nbconvert import MarkdownExporter
import os

# repo_url = "https://raw.githubusercontent.com/SingularitySmith/AMTAIR_Prototype/main/data/"
notebook_name = "AMTAIR_Prototype_example_carlsmith" #Change Notebook name and path when wo

# Download the notebook file
!wget {repo_url}{notebook_name}.ipynb -O {notebook_name}.ipynb # Corrected line

# Load the notebook
# add error handling for file not found
try:
    with open(f"{notebook_name}.ipynb") as f:
        nb = nbformat.read(f, as_version=4)
except FileNotFoundError:
    print(f"Error: File '{notebook_name}.ipynb' not found. Please check if it was downloaded c

# Initialize the Markdown exporter
exporter = MarkdownExporter(exclude_output=True) # Correct initialization

# Convert the notebook to Markdown
(body, resources) = exporter.from_notebook_node(nb)

# Save the Markdown to a file
with open(f"{notebook_name}IPYNB.md", "w") as f:
    f.write(body)
```

```
--2025-05-24 20:09:47-- https://raw.githubusercontent.com/SingularitySmith/AMTAIR_Prototype
Resolving raw.githubusercontent.com (raw.githubusercontent.com)... 185.199.108.133, 185.199.
Connecting to raw.githubusercontent.com (raw.githubusercontent.com)|185.199.108.133|:443...
HTTP request sent, awaiting response... 200 OK
Length: 1689816 (1.6M) [text/plain]
Saving to: 'AMTAIR_Prototype_example_carlsmith.ipynb'
```

```
AMTAIR_Prototype_ex 100%[=====>] 1.61M 5.38MB/s in 0.3s
```

```
2025-05-24 20:09:48 (5.38 MB/s) - 'AMTAIR_Prototype_example_carlsmith.ipynb' saved [1689816/
```

6.2 Convert Notebook to Markdown Documentation

```
# @title 6.2.0 --- Convert Notebook to Markdown Documentation --- [convert_notebook_to_markdown]

"""
BLOCK PURPOSE: Converts the notebook to Markdown format for documentation purposes.

Markdown is a lightweight markup language that is widely used for documentation
and is easily readable in both plain text and rendered formats. This conversion:

1. Preserves the structure and content of the notebook
2. Creates a format suitable for inclusion in documentation systems
3. Excludes code outputs to focus on the process and methodology
4. Supports version control and collaboration on GitHub

The resulting Markdown file can be used in project documentation, GitHub wikis,
or as a standalone reference guide to the AMTAIR extraction pipeline.

DEPENDENCIES: nbformat, nbconvert.MarkdownExporter modules
INPUTS: Current notebook state
OUTPUTS: Markdown version of the notebook
"""

import nbformat
from nbconvert import MarkdownExporter
import os

# Repository URL variable for file access
# repo_url = "https://raw.githubusercontent.com/SingularitySmith/AMTAIR_Prototype/main/data/"
notebook_name = "AMTAIR_Prototype_example_carlsmith" # Change when working with different e
```



```

# Download the notebook file
!wget {repo_url}{notebook_name}.ipynb -O {notebook_name}.ipynb

# Load the notebook
try:
    with open(f"{notebook_name}.ipynb") as f:
        nb = nbformat.read(f, as_version=4)
        print(f" Successfully loaded notebook: {notebook_name}.ipynb")
except FileNotFoundError:
    print(f" Error: File '{notebook_name}.ipynb' not found. Please check "
          + "if it was downloaded correctly.")

# Initialize the Markdown exporter
exporter = MarkdownExporter(exclude_output=True) # Exclude outputs for cleaner documentation

# Convert the notebook to Markdown
try:
    (body, resources) = exporter.from_notebook_node(nb)

    # Save the Markdown to a file
    with open(f"{notebook_name}IPYNB.md", "w") as f:
        f.write(body)
    print(f" Successfully saved Markdown version to: {notebook_name}IPYNB.md")
except Exception as e:
    print(f" Error converting notebook to Markdown: {str(e)}")

```

```

--2025-05-24 20:09:53--  https://raw.githubusercontent.com/SingularitySmith/AMTAIR_Prototype
Resolving raw.githubusercontent.com (raw.githubusercontent.com)... 185.199.108.133, 185.199.
Connecting to raw.githubusercontent.com (raw.githubusercontent.com)|185.199.108.133|:443...
HTTP request sent, awaiting response... 200 OK
Length: 1689816 (1.6M) [text/plain]
Saving to: 'AMTAIR_Prototype_example_carlsmith.ipynb'

```

```

AMTAIR_Prototype_ex 100%[=====>]    1.61M  5.78MB/s    in 0.3s

```

```

2025-05-24 20:09:53 (5.78 MB/s) - 'AMTAIR_Prototype_example_carlsmith.ipynb' saved [1689816/

```

```

Successfully loaded notebook: AMTAIR_Prototype_example_carlsmith.ipynb

```

```

Successfully saved Markdown version to: AMTAIR_Prototype_example_carlsmithIPYNB.md

```

6.3 Create PDF and Latex

```
# @title 6.3.0 --- PDF and Latex--- [pdf_and_latex]

import nbformat
from nbconvert import PDFExporter
import os
import subprocess
import re

def escape_latex_special_chars(text):
    """Escapes special LaTeX characters in a string."""
    latex_special_chars = ['&', '%', '#', '_', '{', '}', '~', '^', '\\']
    replacement_patterns = [
        (char, '\\' + char) for char in latex_special_chars
    ]

    # Escape reserved characters
    for original, replacement in replacement_patterns:
        text = text.replace(original, replacement) # This is the fix
    return text

# Function to check if a command is available
def is_command_available(command):
    try:
        subprocess.run([command], capture_output=True, check=True)
        return True
    except (subprocess.CalledProcessError, FileNotFoundError):
        return False

# Check if xelatex is installed, and install if necessary
if not is_command_available("xelatex"):
    print("Installing necessary TeX packages...")
    !apt-get install -y texlive-xetex texlive-fonts-recommended texlive-plain-generic
    print("TeX packages installed successfully.")
else:
    print("xelatex is already installed. Skipping installation.")

# repo_url = "https://raw.githubusercontent.com/SingularitySmith/AMTAIR_Prototype/main/data/
notebook_name = "AMTAIR_Prototype_example_carlsmith" # Change Notebook name
# and path when working on different examples
```

```

# Download the notebook file
!wget {repo_url}{notebook_name}.ipynb -O {notebook_name}.ipynb # Corrected line

# Load the notebook
# add error handling for file not found
try:
    with open(f"{notebook_name}.ipynb") as f:
        nb = nbformat.read(f, as_version=4)
except FileNotFoundError:
    print(f"Error: File '{notebook_name}.ipynb' not found. Please check if it was downloaded c

# Initialize the PDF exporter
exporter = PDFExporter(exclude_output=True) # Changed to PDFExporter

# Sanitize notebook cell titles to escape special LaTeX characters like '&'
for cell in nb.cells:
    if 'cell_type' in cell and cell['cell_type'] == 'markdown':
        if 'source' in cell and isinstance(cell['source'], str):
            # Replace '&' with '\protect&' in markdown cell titles AND CONTENT
            # Updated to use escape_latex_special_chars function
            cell['source'] = escape_latex_special_chars(cell['source'])
            # Additionally, escape special characters in headings
            cell['source'] = re.sub(r'(\#+)\s*(.*)', lambda m: m.group(1) + ' ' + escape_latex

# Convert the notebook to PDF
(body, resources) = exporter.from_notebook_node(nb)

# Save the PDF to a file
with open(f"{notebook_name}IPYNB.pdf", "wb") as f: # Changed to 'wb' for binary writing
    f.write(body)

```

Installing necessary TeX packages...

Reading package lists... Done

Building dependency tree... Done

Reading state information... Done

The following additional packages will be installed:

dvisvgm fonts-droid-fallback fonts-lato fonts-lmodern fonts-noto-mono

fonts-texgyre fonts-urw-base35 libapache-pom-java libcommons-logging-java

```
libcommons-parent-java libfontbox-java libgs9 libgs9-common libidn12
libijs-0.35 libjbig2dec0 libkpathsea6 libpdfbox-java libptexenc1 libruby3.0
libsynctex2 libteckit0 libtexlua53 libtexluajit2 libwoff1 libzip-0-13
lmodern poppler-data preview-latex-style rake ruby ruby-net-telnet
ruby-rubygems ruby-webrick ruby-xmlrpc ruby3.0 rubygems-integration t1utils
teckit tex-common tex-gyre texlive-base texlive-binaries texlive-latex-base
texlive-latex-extra texlive-latex-recommended texlive-pictures tipa
xfonts-encodings xfonts-utils
```

Suggested packages:

```
fonts-noto fonts-freefont-otf | fonts-freefont-ttf libavalon-framework-java
libcommons-logging-java-doc libexcalibur-logkit-java liblog4j1.2-java
poppler-utils ghostscript fonts-japanese-mincho | fonts-ipafont-mincho
fonts-japanese-gothic | fonts-ipafont-gothic fonts-arphic-ukai
fonts-arphic-uming fonts-nanum ri ruby-dev bundler debhelper gv
| postscript-viewer perl-tk xpdf | pdf-viewer xzdec
texlive-fonts-recommended-doc texlive-latex-base-doc python3-pygments
icc-profiles libfile-which-perl libspreadsheet-parseexcel-perl
texlive-latex-extra-doc texlive-latex-recommended-doc texlive-luatex
texlive-pstricks dot2tex prerex texlive-pictures-doc vprerex
default-jre-headless tipa-doc
```

The following NEW packages will be installed:

```
dvisvgm fonts-droid-fallback fonts-lato fonts-lmodern fonts-noto-mono
fonts-texgyre fonts-urw-base35 libapache-pom-java libcommons-logging-java
libcommons-parent-java libfontbox-java libgs9 libgs9-common libidn12
libijs-0.35 libjbig2dec0 libkpathsea6 libpdfbox-java libptexenc1 libruby3.0
libsynctex2 libteckit0 libtexlua53 libtexluajit2 libwoff1 libzip-0-13
lmodern poppler-data preview-latex-style rake ruby ruby-net-telnet
ruby-rubygems ruby-webrick ruby-xmlrpc ruby3.0 rubygems-integration t1utils
teckit tex-common tex-gyre texlive-base texlive-binaries
texlive-fonts-recommended texlive-latex-base texlive-latex-extra
texlive-latex-recommended texlive-pictures texlive-plain-generic
texlive-xetex tipa xfonts-encodings xfonts-utils
```

0 upgraded, 53 newly installed, 0 to remove and 34 not upgraded.

Need to get 182 MB of archives.

After this operation, 571 MB of additional disk space will be used.

Get:1 <http://archive.ubuntu.com/ubuntu> jammy/main amd64 fonts-droid-fallback all 1:6.0.1r16-

Get:2 <http://archive.ubuntu.com/ubuntu> jammy/main amd64 fonts-lato all 2.0-2.1 [2,696 kB]

Get:3 <http://archive.ubuntu.com/ubuntu> jammy/main amd64 poppler-data all 0.4.11-1 [2,171 kB]

Get:4 <http://archive.ubuntu.com/ubuntu> jammy/universe amd64 tex-common all 6.17 [33.7 kB]

Get:5 <http://archive.ubuntu.com/ubuntu> jammy/main amd64 fonts-urw-base35 all 20200910-1 [6,3

Get:6 <http://archive.ubuntu.com/ubuntu> jammy-updates/main amd64 libgs9-common all 9.55.0~dfs

Get:7 <http://archive.ubuntu.com/ubuntu> jammy-updates/main amd64 libidn12 amd64 1.38-4ubuntu1

```

Get:8 http://archive.ubuntu.com/ubuntu jammy/main amd64 libijs-0.35 amd64 0.35-15build2 [16.
Get:9 http://archive.ubuntu.com/ubuntu jammy/main amd64 libjbig2dec0 amd64 0.19-3build2 [64.
Get:10 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 libgs9 amd64 9.55.0~dfsg1-0
Get:11 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 libkpathsea6 amd64 2021.202
Get:12 http://archive.ubuntu.com/ubuntu jammy/main amd64 libwoff1 amd64 1.0.2-1build4 [45.2
Get:13 http://archive.ubuntu.com/ubuntu jammy/universe amd64 dvisvgm amd64 2.13.1-1 [1,221 k
Get:14 http://archive.ubuntu.com/ubuntu jammy/universe amd64 fonts-lmodern all 2.004.5-6.1
Get:15 http://archive.ubuntu.com/ubuntu jammy/main amd64 fonts-noto-mono all 20201225-1buil
Get:16 http://archive.ubuntu.com/ubuntu jammy/universe amd64 fonts-texgyre all 20180621-3.1
Get:17 http://archive.ubuntu.com/ubuntu jammy/universe amd64 libapache-pom-java all 18-1 [4,
Get:18 http://archive.ubuntu.com/ubuntu jammy/universe amd64 libcommons-parent-java all 43-1
Get:19 http://archive.ubuntu.com/ubuntu jammy/universe amd64 libcommons-logging-java all 1.2
Get:20 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 libptexenc1 amd64 2021.2021
Get:21 http://archive.ubuntu.com/ubuntu jammy/main amd64 rubygems-integration all 1.18 [5,33
Get:22 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 ruby3.0 amd64 3.0.2-7ubunt
Get:23 http://archive.ubuntu.com/ubuntu jammy/main amd64 ruby-rubygems all 3.3.5-2 [228 kB]
Get:24 http://archive.ubuntu.com/ubuntu jammy/main amd64 ruby amd64 1:3.0~exp1 [5,100 B]
Get:25 http://archive.ubuntu.com/ubuntu jammy/main amd64 rake all 13.0.6-2 [61.7 kB]
Get:26 http://archive.ubuntu.com/ubuntu jammy/main amd64 ruby-net-telnet all 0.1.1-2 [12.6 k
Get:27 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 ruby-webrick all 1.7.0-3ubu
Get:28 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 ruby-xmlrpc all 0.3.2-1ubun
Get:29 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 libruby3.0 amd64 3.0.2-7ubu
Get:30 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 libsincTeX2 amd64 2021.2021
Get:31 http://archive.ubuntu.com/ubuntu jammy/universe amd64 libteckit0 amd64 2.5.11+ds1-1
Get:32 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 libtexlua53 amd64 2021.2021
Get:33 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 libtexluaajit2 amd64 2021.20
Get:34 http://archive.ubuntu.com/ubuntu jammy/universe amd64 libzzip-0-13 amd64 0.13.72+dfsg
Get:35 http://archive.ubuntu.com/ubuntu jammy/main amd64 xfonts-encodings all 1:1.0.5-0ubunt
Get:36 http://archive.ubuntu.com/ubuntu jammy/main amd64 xfonts-utils amd64 1:7.7+6build2 [9
Get:37 http://archive.ubuntu.com/ubuntu jammy/universe amd64 lmodern all 2.004.5-6.1 [9,471
Get:38 http://archive.ubuntu.com/ubuntu jammy/universe amd64 preview-latex-style all 12.2-1u
Get:39 http://archive.ubuntu.com/ubuntu jammy/main amd64 t1utils amd64 1.41-4build2 [61.3 kB]
Get:40 http://archive.ubuntu.com/ubuntu jammy/universe amd64 teckit amd64 2.5.11+ds1-1 [699
Get:41 http://archive.ubuntu.com/ubuntu jammy/universe amd64 tex-gyre all 20180621-3.1 [6,20
Get:42 http://archive.ubuntu.com/ubuntu jammy-updates/universe amd64 texlive-binaries amd64
Get:43 http://archive.ubuntu.com/ubuntu jammy/universe amd64 texlive-base all 2021.20220204-
Get:44 http://archive.ubuntu.com/ubuntu jammy/universe amd64 texlive-fonts-recommended all 2
Get:45 http://archive.ubuntu.com/ubuntu jammy/universe amd64 texlive-latex-base all 2021.202
Get:46 http://archive.ubuntu.com/ubuntu jammy/universe amd64 libfontbox-java all 1:1.8.16-2
Get:47 http://archive.ubuntu.com/ubuntu jammy/universe amd64 libpdfbox-java all 1:1.8.16-2
Get:48 http://archive.ubuntu.com/ubuntu jammy/universe amd64 texlive-latex-recommended all 2
Get:49 http://archive.ubuntu.com/ubuntu jammy/universe amd64 texlive-pictures all 2021.20220

```

```
Get:50 http://archive.ubuntu.com/ubuntu jammy/universe amd64 texlive-latex-extra all 2021.20200910-1 [2,967 kB]
Get:51 http://archive.ubuntu.com/ubuntu jammy/universe amd64 texlive-plain-generic all 2021.20200910-1 [2,967 kB]
Get:52 http://archive.ubuntu.com/ubuntu jammy/universe amd64 tipa all 2:1.3-21 [2,967 kB]
Get:53 http://archive.ubuntu.com/ubuntu jammy/universe amd64 texlive-xetex all 2021.20220204-1 [2,967 kB]
Fetched 182 MB in 12s (15.2 MB/s)
Extracting templates from packages: 100%
Preconfiguring packages ...
Selecting previously unselected package fonts-droid-fallback.
(Reading database ... 126327 files and directories currently installed.)
Preparing to unpack .../00-fonts-droid-fallback_1%3a6.0.1r16-1.1build1_all.deb ...
Unpacking fonts-droid-fallback (1:6.0.1r16-1.1build1) ...
Selecting previously unselected package fonts-lato.
Preparing to unpack .../01-fonts-lato_2.0-2.1_all.deb ...
Unpacking fonts-lato (2.0-2.1) ...
Selecting previously unselected package poppler-data.
Preparing to unpack .../02-poppler-data_0.4.11-1_all.deb ...
Unpacking poppler-data (0.4.11-1) ...
Selecting previously unselected package tex-common.
Preparing to unpack .../03-tex-common_6.17_all.deb ...
Unpacking tex-common (6.17) ...
Selecting previously unselected package fonts-urw-base35.
Preparing to unpack .../04-fonts-urw-base35_20200910-1_all.deb ...
Unpacking fonts-urw-base35 (20200910-1) ...
Selecting previously unselected package libgs9-common.
Preparing to unpack .../05-libgs9-common_9.55.0~dfsg1-0ubuntu5.11_all.deb ...
Unpacking libgs9-common (9.55.0~dfsg1-0ubuntu5.11) ...
Selecting previously unselected package libidn12:amd64.
Preparing to unpack .../06-libidn12_1.38-4ubuntu1_amd64.deb ...
Unpacking libidn12:amd64 (1.38-4ubuntu1) ...
Selecting previously unselected package libijs-0.35:amd64.
Preparing to unpack .../07-libijs-0.35_0.35-15build2_amd64.deb ...
Unpacking libijs-0.35:amd64 (0.35-15build2) ...
Selecting previously unselected package libjbig2dec0:amd64.
Preparing to unpack .../08-libjbig2dec0_0.19-3build2_amd64.deb ...
Unpacking libjbig2dec0:amd64 (0.19-3build2) ...
Selecting previously unselected package libgs9:amd64.
Preparing to unpack .../09-libgs9_9.55.0~dfsg1-0ubuntu5.11_amd64.deb ...
Unpacking libgs9:amd64 (9.55.0~dfsg1-0ubuntu5.11) ...
Selecting previously unselected package libkpathsea6:amd64.
Preparing to unpack .../10-libkpathsea6_2021.20210626.59705-1ubuntu0.2_amd64.deb ...
Unpacking libkpathsea6:amd64 (2021.20210626.59705-1ubuntu0.2) ...
Selecting previously unselected package libwoff1:amd64.
```

```
Preparing to unpack .../11-libwoff1_1.0.2-1build4_amd64.deb ...
Unpacking libwoff1:amd64 (1.0.2-1build4) ...
Selecting previously unselected package dvisvgm.
Preparing to unpack .../12-dvisvgm_2.13.1-1_amd64.deb ...
Unpacking dvisvgm (2.13.1-1) ...
Selecting previously unselected package fonts-lmodern.
Preparing to unpack .../13-fonts-lmodern_2.004.5-6.1_all.deb ...
Unpacking fonts-lmodern (2.004.5-6.1) ...
Selecting previously unselected package fonts-noto-mono.
Preparing to unpack .../14-fonts-noto-mono_20201225-1build1_all.deb ...
Unpacking fonts-noto-mono (20201225-1build1) ...
Selecting previously unselected package fonts-texgyre.
Preparing to unpack .../15-fonts-texgyre_20180621-3.1_all.deb ...
Unpacking fonts-texgyre (20180621-3.1) ...
Selecting previously unselected package libapache-pom-java.
Preparing to unpack .../16-libapache-pom-java_18-1_all.deb ...
Unpacking libapache-pom-java (18-1) ...
Selecting previously unselected package libcommons-parent-java.
Preparing to unpack .../17-libcommons-parent-java_43-1_all.deb ...
Unpacking libcommons-parent-java (43-1) ...
Selecting previously unselected package libcommons-logging-java.
Preparing to unpack .../18-libcommons-logging-java_1.2-2_all.deb ...
Unpacking libcommons-logging-java (1.2-2) ...
Selecting previously unselected package libptexenc1:amd64.
Preparing to unpack .../19-libptexenc1_2021.20210626.59705-1ubuntu0.2_amd64.deb ...
Unpacking libptexenc1:amd64 (2021.20210626.59705-1ubuntu0.2) ...
Selecting previously unselected package rubygems-integration.
Preparing to unpack .../20-rubygems-integration_1.18_all.deb ...
Unpacking rubygems-integration (1.18) ...
Selecting previously unselected package ruby3.0.
Preparing to unpack .../21-ruby3.0_3.0.2-7ubuntu2.10_amd64.deb ...
Unpacking ruby3.0 (3.0.2-7ubuntu2.10) ...
Selecting previously unselected package ruby-rubygems.
Preparing to unpack .../22-ruby-rubygems_3.3.5-2_all.deb ...
Unpacking ruby-rubygems (3.3.5-2) ...
Selecting previously unselected package ruby.
Preparing to unpack .../23-ruby_1%3a3.0~exp1_amd64.deb ...
Unpacking ruby (1:3.0~exp1) ...
Selecting previously unselected package rake.
Preparing to unpack .../24-rake_13.0.6-2_all.deb ...
Unpacking rake (13.0.6-2) ...
Selecting previously unselected package ruby-net-telnet.
```

```
Preparing to unpack .../25-ruby-net-telnet_0.1.1-2_all.deb ...
Unpacking ruby-net-telnet (0.1.1-2) ...
Selecting previously unselected package ruby-webrick.
Preparing to unpack .../26-ruby-webrick_1.7.0-3ubuntu0.1_all.deb ...
Unpacking ruby-webrick (1.7.0-3ubuntu0.1) ...
Selecting previously unselected package ruby-xmlrpc.
Preparing to unpack .../27-ruby-xmlrpc_0.3.2-1ubuntu0.1_all.deb ...
Unpacking ruby-xmlrpc (0.3.2-1ubuntu0.1) ...
Selecting previously unselected package libruby3.0:amd64.
Preparing to unpack .../28-libruby3.0_3.0.2-7ubuntu2.10_amd64.deb ...
Unpacking libruby3.0:amd64 (3.0.2-7ubuntu2.10) ...
Selecting previously unselected package libsyntax2:amd64.
Preparing to unpack .../29-libsyntax2_2021.20210626.59705-1ubuntu0.2_amd64.deb ...
Unpacking libsyntax2:amd64 (2021.20210626.59705-1ubuntu0.2) ...
Selecting previously unselected package libteckit0:amd64.
Preparing to unpack .../30-libteckit0_2.5.11+ds1-1_amd64.deb ...
Unpacking libteckit0:amd64 (2.5.11+ds1-1) ...
Selecting previously unselected package libtexlua53:amd64.
Preparing to unpack .../31-libtexlua53_2021.20210626.59705-1ubuntu0.2_amd64.deb ...
Unpacking libtexlua53:amd64 (2021.20210626.59705-1ubuntu0.2) ...
Selecting previously unselected package libtexluaajit2:amd64.
Preparing to unpack .../32-libtexluaajit2_2021.20210626.59705-1ubuntu0.2_amd64.deb ...
Unpacking libtexluaajit2:amd64 (2021.20210626.59705-1ubuntu0.2) ...
Selecting previously unselected package libzip-0-13:amd64.
Preparing to unpack .../33-libzip-0-13_0.13.72+dfsg.1-1.1_amd64.deb ...
Unpacking libzip-0-13:amd64 (0.13.72+dfsg.1-1.1) ...
Selecting previously unselected package xfonts-encodings.
Preparing to unpack .../34-xfonts-encodings_1%3a1.0.5-0ubuntu2_all.deb ...
Unpacking xfonts-encodings (1:1.0.5-0ubuntu2) ...
Selecting previously unselected package xfonts-utils.
Preparing to unpack .../35-xfonts-utils_1%3a7.7+6build2_amd64.deb ...
Unpacking xfonts-utils (1:7.7+6build2) ...
Selecting previously unselected package lmodern.
Preparing to unpack .../36-lmodern_2.004.5-6.1_all.deb ...
Unpacking lmodern (2.004.5-6.1) ...
Selecting previously unselected package preview-latex-style.
Preparing to unpack .../37-preview-latex-style_12.2-1ubuntu1_all.deb ...
Unpacking preview-latex-style (12.2-1ubuntu1) ...
Selecting previously unselected package t1utils.
Preparing to unpack .../38-t1utils_1.41-4build2_amd64.deb ...
Unpacking t1utils (1.41-4build2) ...
Selecting previously unselected package teckit.
```



```
Preparing to unpack .../39-teckit_2.5.11+ds1-1_amd64.deb ...
Unpacking teckit (2.5.11+ds1-1) ...
Selecting previously unselected package tex-gyre.
Preparing to unpack .../40-tex-gyre_20180621-3.1_all.deb ...
Unpacking tex-gyre (20180621-3.1) ...
Selecting previously unselected package texlive-binaries.
Preparing to unpack .../41-texlive-binaries_2021.20210626.59705-1ubuntu0.2_amd64.deb ...
Unpacking texlive-binaries (2021.20210626.59705-1ubuntu0.2) ...
Selecting previously unselected package texlive-base.
Preparing to unpack .../42-texlive-base_2021.20220204-1_all.deb ...
Unpacking texlive-base (2021.20220204-1) ...
Selecting previously unselected package texlive-fonts-recommended.
Preparing to unpack .../43-texlive-fonts-recommended_2021.20220204-1_all.deb ...
Unpacking texlive-fonts-recommended (2021.20220204-1) ...
Selecting previously unselected package texlive-latex-base.
Preparing to unpack .../44-texlive-latex-base_2021.20220204-1_all.deb ...
Unpacking texlive-latex-base (2021.20220204-1) ...
Selecting previously unselected package libfontbox-java.
Preparing to unpack .../45-libfontbox-java_1%3a1.8.16-2_all.deb ...
Unpacking libfontbox-java (1:1.8.16-2) ...
Selecting previously unselected package libpdfbox-java.
Preparing to unpack .../46-libpdfbox-java_1%3a1.8.16-2_all.deb ...
Unpacking libpdfbox-java (1:1.8.16-2) ...
Selecting previously unselected package texlive-latex-recommended.
Preparing to unpack .../47-texlive-latex-recommended_2021.20220204-1_all.deb ...
Unpacking texlive-latex-recommended (2021.20220204-1) ...
Selecting previously unselected package texlive-pictures.
Preparing to unpack .../48-texlive-pictures_2021.20220204-1_all.deb ...
Unpacking texlive-pictures (2021.20220204-1) ...
Selecting previously unselected package texlive-latex-extra.
Preparing to unpack .../49-texlive-latex-extra_2021.20220204-1_all.deb ...
Unpacking texlive-latex-extra (2021.20220204-1) ...
Selecting previously unselected package texlive-plain-generic.
Preparing to unpack .../50-texlive-plain-generic_2021.20220204-1_all.deb ...
Unpacking texlive-plain-generic (2021.20220204-1) ...
Selecting previously unselected package tipa.
Preparing to unpack .../51-tipa_2%3a1.3-21_all.deb ...
Unpacking tipa (2:1.3-21) ...
Selecting previously unselected package texlive-xetex.
Preparing to unpack .../52-texlive-xetex_2021.20220204-1_all.deb ...
Unpacking texlive-xetex (2021.20220204-1) ...
Setting up fonts-lato (2.0-2.1) ...
```

```

Setting up fonts-noto-mono (20201225-1build1) ...
Setting up libwoff1:amd64 (1.0.2-1build4) ...
Setting up libtexlua53:amd64 (2021.20210626.59705-1ubuntu0.2) ...
Setting up libijs-0.35:amd64 (0.35-15build2) ...
Setting up libtexluajit2:amd64 (2021.20210626.59705-1ubuntu0.2) ...
Setting up libfontbox-java (1:1.8.16-2) ...
Setting up rubygems-integration (1.18) ...
Setting up libzip-0-13:amd64 (0.13.72+dfsg.1-1.1) ...
Setting up fonts-urw-base35 (20200910-1) ...
Setting up poppler-data (0.4.11-1) ...
Setting up tex-common (6.17) ...
update-language: texlive-base not installed and configured, doing nothing!
Setting up libjbig2dec0:amd64 (0.19-3build2) ...
Setting up libteckit0:amd64 (2.5.11+ds1-1) ...
Setting up libapache-pom-java (18-1) ...
Setting up ruby-net-telnet (0.1.1-2) ...
Setting up xfonts-encodings (1:1.0.5-0ubuntu2) ...
Setting up t1utils (1.41-4build2) ...
Setting up libidn12:amd64 (1.38-4ubuntu1) ...
Setting up fonts-texgyre (20180621-3.1) ...
Setting up libkpathsea6:amd64 (2021.20210626.59705-1ubuntu0.2) ...
Setting up ruby-webrick (1.7.0-3ubuntu0.1) ...
Setting up fonts-lmodern (2.004.5-6.1) ...
Setting up fonts-droid-fallback (1:6.0.1r16-1.1build1) ...
Setting up ruby-xmlrpc (0.3.2-1ubuntu0.1) ...
Setting up libsynchronet2:amd64 (2021.20210626.59705-1ubuntu0.2) ...
Setting up libgs9-common (9.55.0~dfsg1-0ubuntu5.11) ...
Setting up teckit (2.5.11+ds1-1) ...
Setting up libpdfbox-java (1:1.8.16-2) ...
Setting up libgs9:amd64 (9.55.0~dfsg1-0ubuntu5.11) ...
Setting up preview-latex-style (12.2-1ubuntu1) ...
Setting up libcommons-parent-java (43-1) ...
Setting up dvisvgm (2.13.1-1) ...
Setting up libcommons-logging-java (1.2-2) ...
Setting up xfonts-utils (1:7.7+6build2) ...
Setting up libptexenc1:amd64 (2021.20210626.59705-1ubuntu0.2) ...
Setting up texlive-binaries (2021.20210626.59705-1ubuntu0.2) ...
update-alternatives: using /usr/bin/xdvi-xaw to provide /usr/bin/xdvi.bin (xdvi.bin) in auto
update-alternatives: using /usr/bin/bibtex.original to provide /usr/bin/bibtex (bibtex) in a
Setting up lmodern (2.004.5-6.1) ...
Setting up texlive-base (2021.20220204-1) ...
/usr/bin/ucfr

```

```

/usr/bin/ucfr
/usr/bin/ucfr
/usr/bin/ucfr
mktexlsr: Updating /var/lib/texmf/ls-R-TEXLIVEDIST...
mktexlsr: Updating /var/lib/texmf/ls-R-TEXMFMAIN...
mktexlsr: Updating /var/lib/texmf/ls-R...
mktexlsr: Done.
tl-paper: setting paper size for dvips to a4: /var/lib/texmf/dvips/config/config-paper.ps
tl-paper: setting paper size for dvipdfmx to a4: /var/lib/texmf/dvipdfmx/dvipdfmx-paper.cfg
tl-paper: setting paper size for xdvi to a4: /var/lib/texmf/xdvi/XDvi-paper
tl-paper: setting paper size for pdftex to a4: /var/lib/texmf/tex/generic/tex-ini-files/pdft
Setting up tex-gyre (20180621-3.1) ...
Setting up texlive-plain-generic (2021.20220204-1) ...
Setting up texlive-latex-base (2021.20220204-1) ...
Setting up texlive-latex-recommended (2021.20220204-1) ...
Setting up texlive-pictures (2021.20220204-1) ...
Setting up texlive-fonts-recommended (2021.20220204-1) ...
Setting up tipa (2:1.3-21) ...
Setting up texlive-latex-extra (2021.20220204-1) ...
Setting up texlive-xetex (2021.20220204-1) ...
Setting up rake (13.0.6-2) ...
Setting up libruby3.0:amd64 (3.0.2-7ubuntu2.10) ...
Setting up ruby3.0 (3.0.2-7ubuntu2.10) ...
Setting up ruby (1:3.0~exp1) ...
Setting up ruby-rubygems (3.3.5-2) ...
Processing triggers for man-db (2.10.2-1) ...
Processing triggers for mailcap (3.70+nmulubuntu1) ...
Processing triggers for fontconfig (2.13.1-4.2ubuntu5) ...
Processing triggers for libc-bin (2.35-0ubuntu3.8) ...
/sbin/ldconfig.real: /usr/local/lib/libumf.so.0 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtcm_debug.so.1 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtbbbbind.so.3 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libhwloc.so.15 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtbbbbind_2_5.so.3 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtbb.so.12 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libur_adapter_level_zero.so.0 is not a symbolic link

```

```

/sbin/ldconfig.real: /usr/local/lib/libtbbmalloc.so.2 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtbbbind_2_0.so.3 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtcm.so.1 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libur_adapter_opengl.so.0 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtbbmalloc_proxy.so.2 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libur_loader.so.0 is not a symbolic link

Processing triggers for tex-common (6.17) ...
Running updmap-sys. This may take some time... done.
Running mktexlsr /var/lib/texmf ... done.
Building format(s) --all.
    This may take some time... done.
TeX packages installed successfully.
--2025-05-24 20:12:59-- https://raw.githubusercontent.com/SingularitySmith/AMTAIR_Prototype
Resolving raw.githubusercontent.com (raw.githubusercontent.com)... 185.199.108.133, 185.199.
Connecting to raw.githubusercontent.com (raw.githubusercontent.com)|185.199.108.133|:443...
HTTP request sent, awaiting response... 200 OK
Length: 1689816 (1.6M) [text/plain]
Saving to: 'AMTAIR_Prototype_example_carlsmith.ipynb'

AMTAIR_Prototype_ex 100%[=====>] 1.61M 5.31MB/s in 0.3s

2025-05-24 20:12:59 (5.31 MB/s) - 'AMTAIR_Prototype_example_carlsmith.ipynb' saved [1689816/

```

Bibliography

- [1] Dario Amodei et al. *Concrete Problems in AI Safety*. July 25, 2016. DOI: 10.48550/arXiv.1606.06565. arXiv: 1606.06565 [cs]. URL: <http://arxiv.org/abs/1606.06565> (visited on 05/25/2025). Pre-published.
- [2] Terence J. Anderson. “Visualization Tools and Argument Schemes: A Question of Standpoint”. In: *Law, Prob. & Risk* 6 (2007), p. 97. URL: https://heinonline.org/hol-cgi-bin/get_pdf.cgi?handle=hein.journals/lawprisk6§ion=9 (visited on 05/25/2025).
- [3] Stuart Armstrong, Nick Bostrom, and Carl Shulman. “Racing to the Precipice: A Model of Artificial Intelligence Development”. In: *AI & SOCIETY* 31.2 (May 1, 2016), pp. 201–206. ISSN: 1435-5655. DOI: 10.1007/s00146-015-0590-y. URL: <https://doi.org/10.1007/s00146-015-0590-y> (visited on 05/26/2025).
- [4] Amanda Askell et al. *A General Language Assistant as a Laboratory for Alignment*. Dec. 9, 2021. DOI: 10.48550/arXiv.2112.00861. arXiv: 2112.00861 [cs]. URL: <http://arxiv.org/abs/2112.00861> (visited on 05/25/2025). Pre-published.
- [5] Nikolay Babakov et al. “Reusability of Bayesian Networks Case Studies: A Survey”. In: *Applied Intelligence* 55.6 (Feb. 7, 2025), p. 417. ISSN: 1573-7497. DOI: 10.1007/s10489-025-06289-5. URL: <https://doi.org/10.1007/s10489-025-06289-5> (visited on 05/15/2025).
- [6] Taiyu Ban et al. *Causal Structure Learning Supervised by Large Language Model*. Nov. 20, 2023. DOI: 10.48550/arXiv.2311.11689. arXiv: 2311.11689 [cs]. URL: <http://arxiv.org/abs/2311.11689> (visited on 05/26/2025). Pre-published.
- [7] Neil Benn and Ann Macintosh. “Argument Visualization for eParticipation: Towards a Research Agenda and Prototype Tool”. In: *Electronic Participation*. Ed. by Efthimios Tambouris, Ann Macintosh, and Hans De Bruijn. Red. by David Hutchison et al. Vol. 6847. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 60–73. ISBN: 978-3-642-23332-6 978-3-642-23333-3. DOI: 10.1007/978-3-642-23333-3_6. URL: http://link.springer.com/10.1007/978-3-642-23333-3_6 (visited on 05/25/2025).
- [8] Steven John Bethard. “Finding Event, Temporal and Causal Structure in Text: A Machine Learning Approach”. PhD thesis. University of Colorado at Boulder, 2007. URL: <https://search.proquest.com/openview/405fe32503123d9b5f4836dc3be4c011/1?pq-origsite=gscholar&cbl=18750> (visited on 05/26/2025).
- [9] Nick Bostrom. *Superintelligence: Paths, Strategies, Dangers*. Oxford: Oxford University Press, 2014. ISBN: 978-0-19-967811-2. URL: <https://scholar.dominican.edu/cynthia-stokes-brown-books-big-history/47>.

- [10] George E. P. Box. “Science and Statistics”. In: *Journal of the American Statistical Association* 71.356 (Dec. 1976), pp. 791–799. ISSN: 0162-1459, 1537-274X. DOI: 10.1080/01621459.1976.10480949. URL: <http://www.tandfonline.com/doi/abs/10.1080/01621459.1976.10480949> (visited on 05/26/2025).
- [11] Miles Brundage et al. *The Malicious Use of Artificial Intelligence: Forecasting, Prevention, and Mitigation*. Version 1. 2018. DOI: 10.48550/ARXIV.1802.07228. URL: <https://arxiv.org/abs/1802.07228> (visited on 11/13/2024). Pre-published.
- [12] Miles Brundage et al. “The Malicious Use of Artificial Intelligence: Forecasting, Prevention, and Mitigation”. 2018. arXiv: 1802.07228.
- [13] Joseph Carlsmith. *Is Power-Seeking AI an Existential Risk?* 2021. DOI: 10.48550/arXiv.2206.13353. URL: <https://arxiv.org/abs/2206.13353>. Pre-published.
- [14] Joseph Carlsmith. “Is Power-Seeking AI an Existential Risk?” 2022. arXiv: 2206.13353.
- [15] Joseph Carlsmith. *Is Power-Seeking AI an Existential Risk?* Aug. 13, 2024. DOI: 10.48550/arXiv.2206.13353. arXiv: 2206.13353 [cs]. URL: <http://arxiv.org/abs/2206.13353> (visited on 05/25/2025). Pre-published.
- [16] Lu Chen et al. “Inducing Causal Structure for Abstractive Text Summarization”. In: *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*. CIKM ’23: The 32nd ACM International Conference on Information and Knowledge Management. Birmingham United Kingdom: ACM, Oct. 21, 2023, pp. 213–223. ISBN: 979-8-4007-0124-5. DOI: 10.1145/3583780.3614934. URL: <https://dl.acm.org/doi/10.1145/3583780.3614934> (visited on 05/26/2025).
- [17] Paul F. Christiano. “What Failure Looks Like”. In: (Mar. 17, 2019). URL: <https://www.alignmentforum.org/posts/HBxe6wdjxK239zajf/what-failure-looks-like> (visited on 05/25/2025).
- [18] Sam Clarke et al. *Modeling Transformative AI Risks (MTAIR) Project – Summary Report*. Version 1. 2022. DOI: 10.48550/ARXIV.2206.09360. URL: <https://arxiv.org/abs/2206.09360> (visited on 11/13/2024). Pre-published.
- [19] Francesca Cuomo, Christine Mallin, and Alessandro Zattoni. “Corporate Governance Codes: A Review and Research Agenda”. In: *Corporate governance: an international review* 24.3 (2016), pp. 222–241. URL: <https://ueaeprints.uea.ac.uk/id/eprint/57664/> (visited on 05/26/2025).
- [20] Allan Dafoe. “AI Governance: A Research Agenda”. In: *Governance of AI Program, Future of Humanity Institute, University of Oxford: Oxford, UK* 1442 (2018), p. 1443. URL: <http://www.fhi.ox.ac.uk/wp-content/uploads/GovAI-Agenda.pdf> (visited on 05/25/2025).
- [21] Allan Dafoe. *AI Governance: A Research Agenda*. 2021. URL: <https://www.fhi.ox.ac.uk/wp-content/uploads/GovAI-Agenda.pdf> (visited on 05/25/2025). Pre-published.
- [22] Charl De Villiers and Ruth Dimes. “Determinants, Mechanisms and Consequences of Corporate Governance Reporting: A Research Framework”. In: *Journal of Management and Governance* 25.1 (Mar. 2021), pp. 7–26. ISSN: 1385-3457, 1572-963X. DOI: 10.1007/s10997-020-09530-0. URL: <https://link.springer.com/10.1007/s10997-020-09530-0> (visited on 05/26/2025).

- [23] Istemi Demirag, Sudi Sudarsanam, and MIKE WRIGHT. “Corporate Governance: Overview and Research Agenda”. In: *The British Accounting Review* 32.4 (2000), pp. 341–354. URL: <https://www.academia.edu/download/49469624/bare.2000.014620161009-3955-1dt4aq5.pdf> (visited on 05/26/2025).
- [24] Jackie Di Vito and Kim Trottier. “A Literature Review on Corporate Governance Mechanisms: Past, Present, and Future*”. In: *Accounting Perspectives* 21.2 (June 2022), pp. 207–235. ISSN: 1911-382X, 1911-3838. DOI: 10.1111/1911-3838.12279. URL: <https://onlinelibrary.wiley.com/doi/10.1111/1911-3838.12279> (visited on 05/26/2025).
- [25] K. Eric Drexler. “Reframing Superintelligence: Comprehensive AI Services as General Intelligence”. In: (2019).
- [26] KE Drexler. *Reframing Superintelligence: Comprehensive AI Services as General Intelligence*. Technical Report. Future of Humanity Institute, 2019. URL: https://owainevans.github.io/pdfs/Reframing_Superintelligence_FHI-TR-2019.pdf.
- [27] Pierre Maurice Marie Duhem. *The Aim and Structure of Physical Theory*. 1. Princeton University Press, 1954. 85–87.
- [28] Union European. *The Act Texts / EU Artificial Intelligence Act*. 2024. URL: <https://artificialintelligenceact.eu/the-act/> (visited on 05/25/2025).
- [29] Irving John Good. “Speculations Concerning the First Ultraintelligent Machine”. In: *Advances in Computers* (1966), p. 31. DOI: 10.1016/S0065-2458(08)60418-0. URL: <https://www.sciencedirect.com/science/article/abs/pii/S0065245808604180>.
- [30] Jakub Growiec. “Existential Risk from Transformative AI: An Economic Perspective”. In: *Technological and Economic Development of Economy* 30.6 (2024), pp. 1682–1708.
- [31] Ross Gruetzmacher. “Bayesian Networks vs. Conditional Trees for Creating Questions for Forecasting Tournaments”. In: (2022).
- [32] Stéphane Hallegatte et al. “Investment Decision-Making under Deep Uncertainty-Application to Climate Change”. In: *Policy research working paper* 6193 (2012). URL: <https://enpc.hal.science/hal-00802049/document> (visited on 05/25/2025).
- [33] Christina Heinze-Deml, Marloes H. Maathuis, and Nicolai Meinshausen. “Causal Structure Learning”. In: *Annual Review of Statistics and Its Application* 5.1 (Mar. 7, 2018), pp. 371–391. ISSN: 2326-8298, 2326-831X. DOI: 10.1146/annurev-statistics-031017-100630. URL: <https://www.annualreviews.org/doi/10.1146/annurev-statistics-031017-100630> (visited on 05/26/2025).
- [34] Dan Hendrycks et al. *Unsolved Problems in ML Safety*. Version 5. 2021. DOI: 10.48550/ARXIV.2109.13916. URL: <https://arxiv.org/abs/2109.13916> (visited on 11/13/2024). Pre-published.
- [35] Dan Hendrycks et al. “Unsolved Problems in ML Safety”. 2021. arXiv: 2109.13916.
- [36] Tam Hunt. *The Insane “Logic” of the AI Arms Race*. Medium. Mar. 3, 2025. URL: <https://tamhunt.medium.com/the-insane-logic-of-the-ai-arms-race-45a5f79f4c0e> (visited on 05/26/2025).
- [37] Edwin T Jaynes. *Probability Theory: The Logic of Science*. Cambridge university press, 2003.

- [38] Kawaljit Kaur. “Corporate Governance and Legal Accountability: A Critical Review of Global Practices”. In: *Journal of Law* 2.6 (2024), pp. 1–7. URL: <https://joi.shodhsagar.org/index.php/SSJOI/article/view/16> (visited on 05/26/2025).
- [39] D. Khartabil et al. “Design and Evaluation of Visualization Techniques to Facilitate Argument Exploration”. In: *Computer Graphics Forum* 40.6 (Sept. 2021), pp. 447–465. ISSN: 0167-7055, 1467-8659. DOI: 10.1111/cgf.14389. URL: <https://onlinelibrary.wiley.com/doi/10.1111/cgf.14389> (visited on 05/25/2025).
- [40] Dana Khartabil. “Visualisation Techniques to Facilitate Argument Exploration”. PhD thesis. 2020. URL: <https://napier-repository.worktribe.com/output/2694675> (visited on 05/25/2025).
- [41] Daphne Koller and Nir Friedman. *Probabilistic Graphical Models: Principles and Techniques*. MIT press, 2009. URL: [https://books.google.ca/books?hl=en&lr=&id=7dzpHCHzNQ4C&oi=fnd&pg=PR9&dq=Koller,+D.,+%26+Friedman,+N.+\(2009\).+Probabilistic+Graphical+Models&ots=py2HAh0VAL&sig=gpaID3x6-TY8x5SOopuXpZDXfzs](https://books.google.ca/books?hl=en&lr=&id=7dzpHCHzNQ4C&oi=fnd&pg=PR9&dq=Koller,+D.,+%26+Friedman,+N.+(2009).+Probabilistic+Graphical+Models&ots=py2HAh0VAL&sig=gpaID3x6-TY8x5SOopuXpZDXfzs) (visited on 05/25/2025).
- [42] Ram Shankar Siva Kumar et al. *Failure Modes in Machine Learning Systems*. Version 1. 2019. DOI: 10.48550/ARXIV.1911.11034. URL: <https://arxiv.org/abs/1911.11034> (visited on 11/13/2024). Pre-published.
- [43] Ram Shankar Siva Kumar et al. “Failure Modes in Machine Learning Systems”. 2019. arXiv: 1911.11034.
- [44] Robert J Lempert, Steven W Popper, and Steven C Bankes. *Shaping the next One Hundred Years: New Methods for Quantitative, Long-Term Policy Analysis*. RAND Corporation, 2003.
- [45] Dennis V. Lindley. *Understanding Uncertainty*. John Wiley & Sons, 2013. URL: [https://books.google.ca/books?hl=en&lr=&id=Tfk8AgAAQBAJ&oi=fnd&pg=PR11&dq=Lindley,+D.+\(2006\).+Understanding+Uncertainty&ots=55HS6lTOVP&sig=0gKCDvRu5rUKhuPyJqhOzW23upU](https://books.google.ca/books?hl=en&lr=&id=Tfk8AgAAQBAJ&oi=fnd&pg=PR11&dq=Lindley,+D.+(2006).+Understanding+Uncertainty&ots=55HS6lTOVP&sig=0gKCDvRu5rUKhuPyJqhOzW23upU) (visited on 05/26/2025).
- [46] Christian List and Philip Pettit. *Group Agency: The Possibility, Design, and Status of Corporate Agents*. Oxford University Press, 2011.
- [47] Nestor Maslej. “Artificial Intelligence Index Report 2025”. In: *Artificial Intelligence* (2025).
- [48] Tegan McCaslin et al. “Conditional Trees: A Method for Generating Informative Questions about Complex Topics”. In: *Forecasting Research Institute* (2024). URL: <https://static1.squarespace.com/static/635693acf15a3e2a14a56a4a/t/66ba37a144f1d6095de467df/1723479995772/AIConditionalTrees.pdf>.
- [49] Valentin Jakob Meyer. “A Structure of Knowledge & the Process of Science”. In: *Philosophy of the Social Sciences First Course Paper* (University Bayreuth 2022). DOI: <https://www.vjmeyer.com/papers/essays>.
- [50] Andrea Miotti et al. *A Narrow Path*. 2024. URL: <https://www.narrowpath.co/> (visited on 05/19/2025).
- [51] Roger B. Nelson. *An Introduction to Copulas*. Springer Series in Statistics. New York, NY: Springer New York, 2006. ISBN: 978-0-387-28659-4. DOI: 10.1007/0-387-28678-0. URL: <http://link.springer.com/10.1007/0-387-28678-0> (visited on 05/25/2025).

- [52] Berty Nsolly Ngajie et al. “Investigating the Effects of a Systematic and Model-Based Design of Computer-Supported Argument Visualization on Critical Thinking”. In: *Thinking Skills and Creativity* 38 (2020), p. 100742. URL: <https://www.sciencedirect.com/science/article/pii/S1871187120302169> (visited on 05/25/2025).
- [53] Paul. *The elephAInt – Are We All like the Six Blind Men When It Comes to AI? / PRISMAGuard LLC*. 2023. URL: <https://www.prismaguard.com/the-elephaint-are-we-all-like-the-six-blind-men-when-it-comes-to-ai/> (visited on 05/25/2025).
- [54] Judea Pearl. *Causality: Models, Reasoning and Inference*. 2nd ed. Cambridge University Press, 2009.
- [55] Judea Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge, U.K. ; New York: Cambridge University Press, 2000. 384 pp. ISBN: 978-0-521-89560-6 978-0-521-77362-1.
- [56] Judea Pearl. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Elsevier, 2014. URL: [https://books.google.ca/books?hl=en&lr=&id=mn2jBQAAQBAJ&oi=fnd&pg=PP1&dq=Pearl,+J.+\(1988\).+Probabilistic+Reasoning+in+Intelligent+Systems&ots=4tEX2A4Ha8&sig=lgUs_RCoeXEEuGwM5xMEoyJy4HI](https://books.google.ca/books?hl=en&lr=&id=mn2jBQAAQBAJ&oi=fnd&pg=PP1&dq=Pearl,+J.+(1988).+Probabilistic+Reasoning+in+Intelligent+Systems&ots=4tEX2A4Ha8&sig=lgUs_RCoeXEEuGwM5xMEoyJy4HI) (visited on 05/25/2025).
- [57] John L. Pollock. *Cognitive Carpentry: A Blueprint for How to Build a Person*. Mit Press, 1995. URL: [https://books.google.ca/books?hl=en&lr=&id=JAfHrHTqswAC&oi=fnd&pg=PA1&dq=Pollock,+J.+\(1995\).+Cognitive+Carpentry&ots=rq-qSCBcxV&sig=aAfHGSGUosxl_1-JuxIEA7C2QO4](https://books.google.ca/books?hl=en&lr=&id=JAfHrHTqswAC&oi=fnd&pg=PA1&dq=Pollock,+J.+(1995).+Cognitive+Carpentry&ots=rq-qSCBcxV&sig=aAfHGSGUosxl_1-JuxIEA7C2QO4) (visited on 05/25/2025).
- [58] D. E. Prokudin, E. N. Lisanyuk, and I. R. Baymuratov. “Visualization Functions in Argumentation Representation Software”. In: *Scientific Visualization* 16.3 (2024). URL: <https://sv-journal.org/2024-3/11/en.pdf> (visited on 05/25/2025).
- [59] Iskander Rehman. *The Battle for Brilliant Minds: From the Nuclear Age to AI*. War on the Rocks. Jan. 13, 2025. URL: <https://warontherocks.com/2025/01/the-battle-for-brilliant-minds-from-the-nuclear-age-to-ai/> (visited on 05/25/2025).
- [60] Stuart Russell et al. “Research Priorities for Robust and Beneficial Artificial Intelligence: An Open Letter”. In: *AI Magazine* 36.4 (2015), pp. 3–4. DOI: 10.1609/aimag.v36i4.2621. URL: <https://ojs.aaai.org/aimagazine/index.php/aimagazine/article/view/2621>.
- [61] Veronika Samborska. “Scaling up: How Increasing Inputs Has Made Artificial Intelligence More Capable”. In: *Our World in Data* (Jan. 20, 2025). URL: <https://ourworldindata.org/scaling-up-ai> (visited on 05/25/2025).
- [62] Sigal Samuel. *AI Is a “Tragedy of the Commons.” We’ve Got Solutions for That*. Vox. July 7, 2023. URL: <https://www.vox.com/future-perfect/2023/7/7/23787011/ai-arms-race-tragedy-commons-risk-safety> (visited on 05/26/2025).
- [63] Thomas C. Schelling. “1960. The Strategy of Conflict”. In: *Cambridge, Mass* (1960).
- [64] Oliver Scheuer et al. “Computer-Supported Argumentation: A Review of the State of the Art”. In: *International Journal of Computer-Supported Collaborative Learning* 5.1 (Mar. 2010), pp. 43–102. ISSN: 1556-1607, 1556-1615. DOI: 10.1007/s11412-009-9080-x. URL: <http://link.springer.com/10.1007/s11412-009-9080-x> (visited on 05/25/2025).
- [65] Jill Solomon. *Corporate Governance and Accountability*. John Wiley & Sons, 2020. URL: <https://books.google.ca/books?hl=en&lr=&id=JAX9DwAAQBAJ&oi=fnd&pg=PR1&d>

- q=review+of+the+effects+of+liability+frameworks+on+corporate+governance+&ots=ny23_vd-U0&sig=3LuNNhvSWXriEeg-ipAdDIQGAgO (visited on 05/26/2025).
- [66] Chandler Squires and Caroline Uhler. “Causal Structure Learning: A Combinatorial Perspective”. In: *Foundations of Computational Mathematics* 23.5 (Oct. 2023), pp. 1781–1815. ISSN: 1615-3375, 1615-3383. DOI: 10.1007/s10208-022-09581-9. URL: <https://link.springer.com/10.1007/s10208-022-09581-9> (visited on 05/26/2025).
- [67] Max Tegmark. *Asilomar AI Principles*. Future of Life Institute. 2024. URL: <https://futureoflife.org/open-letter/ai-principles/> (visited on 05/25/2025).
- [68] Philip E. Tetlock and Dan Gardner. *Superforecasting: The Art and Science of Prediction*. First paperback edition. New York: Broadway Books, 2015. 340 pp. ISBN: 978-0-8041-3671-6.
- [69] Benjamin Todd. *It Looks like There Are Some Good Funding Opportunities in AI Safety Right Now*. Benjamin Todd. Dec. 21, 2024. URL: <https://benjamintodd.substack.com/p/looks-like-there-are-some-good-funding> (visited on 05/25/2025).
- [70] Christian Voigt. *Christianvoigt/Argdown*. May 23, 2025. URL: <https://github.com/christianvoigt/argdown> (visited on 05/25/2025).
- [71] Douglas Walton. “Argument Visualization Tools for Corroborative Evidence”. In: *Proc. of the 2nd International Conference on Evidence Law and Forensic Science*. 2009, pp. 32–49. URL: <https://www.academia.edu/download/37718171/09ArguVis.pdf> (visited on 05/25/2025).
- [72] Nick Wilson et al. “The Need for Long-Term Thinking—Especially for Preventing Catastrophic Risks”. In: *Public Health Expert Briefing* (2023).
- [73] Jie Yang, Soyeon Caren Han, and Josiah Poon. “A Survey on Extraction of Causal Relations from Natural Language Text”. In: *Knowledge and Information Systems* 64.5 (May 2022), pp. 1161–1186. ISSN: 0219-1377, 0219-3116. DOI: 10.1007/s10115-022-01665-w. URL: <https://link.springer.com/10.1007/s10115-022-01665-w> (visited on 05/26/2025).
- [74] Eliezer Yudkowsky. “Artificial Intelligence as a Positive and Negative Factor in Global Risk”. In: *Global Catastrophic Risks*. Oxford University Press, July 3, 2008. ISBN: 978-0-19-857050-9 978-0-19-191810-0. DOI: 10.1093/oso/9780198570509.003.0021. URL: <https://academic.oup.com/book/40615/chapter/348239228> (visited on 11/15/2024).



UNIVERSITÄT
BAYREUTH

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

Affidavit

Declaration of Academic Honesty

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

Automating the Modelling of Transformative Artificial Intelligence Risks

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

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

BAYREUTH on the
May 25, 2025

VALENTIN MEYER