

AMTAIR – Automating the Modelling of Transformative Artificial Intelligence Risks

1. Introduction

- **Coordination crisis in AI governance:** ~~Open with a striking example of how misaligned priorities between AI researchers, policymakers, and ethicists lead to fragmented efforts. Highlight the paradox of unprecedented AI investment coexisting with fundamental coordination failures, resulting in safety gaps and misallocated resources (e.g. redundant projects, conflicting regulations).~~ Emphasize that these divisions **increase existential risk** by preventing unified responses to AI threats.
- Explain how **different stakeholders** operate on incompatible assumptions (technical metrics vs. ethical principles vs. geopolitical interests), undermining a shared strategy. Draw parallels to historical coordination failures in nuclear safety or climate change to illustrate stakes.
- Stress the **urgency**: rapid AI capability gains compress the time available to resolve disagreements, making a coordinated approach ever more critical.
- **Research question & scope:** Clearly state the thesis question: *“How can we use frontier AI techniques to automate the extraction of probabilistic world models from AI safety arguments, in order to better predict the impacts of policy interventions?”* Break down key terms:
- Define *“frontier AI”* (cutting-edge language models) and *“probabilistic world models”* (formal structures capturing assumptions and uncertainties about AI risks).
- Clarify scope boundaries: the focus is on existential risk from **misaligned AI** (power-seeking AI scenario), not general AI ethics or all governance problems. Note that the approach centers on **knowledge representation and risk modeling**, intersecting philosophy (argumentation theory) and economics (policy evaluation).
- Justify this scope: tackling this specific modeling problem can yield insights for the broader AI governance challenge without overextending into unrelated debates.
- **Proposed solution overview:** Introduce **AMTAIR** as a novel approach addressing the above coordination crisis. Summarize the core idea: an **automated pipeline** that turns textual arguments from AI safety literature into **Bayesian network models**. Discuss the *“multiplicative benefits”* of combining:
 - **Automated argument extraction** (using AI to parse structured arguments from text),
 - **Probabilistic modeling** (Bayesian networks quantifying beliefs and uncertainties),
 - **Policy impact evaluation** (tools to simulate interventions and outcomes on these networks).
- The **synergy** of these components should be emphasized: by making implicit expert reasoning explicit and quantitative, we gain the ability to run “what-if” scenarios, compare differing worldviews side by side, and identify leverage points for interventions. For example, integrating **forecasting or prediction-market data** into the network could continually refine risk estimates, multiplying the insight each component alone would provide.

- (Figure: Conceptual diagram of how automated extraction, probabilistic reasoning, and policy evaluation interact, illustrating feedback loops and synergy.)
- Underscore the **original contribution**: This thesis proposes a unique combination of argument mapping and AI tools to create a decision-support framework for AI governance that hasn't been explored in prior literature.
- **Thesis structure**: Provide a roadmap of the chapters to orient the reader:
- *Chapter 2 (Context & Background)* establishes theoretical foundations (existential risk models, Bayesian networks) and the methodological basis (argument mapping with ArgDown/BayesDown), situating the project in existing research.
- *Chapter 3 (AMTAIR Implementation)* details the software implementation of the proposed approach, including the step-by-step pipeline (demonstrated with a simple **Rain-Sprinkler-Lawn** example) and the results of applying it to a complex real-world case (**Carlsmith's AI risk model**).
- *Chapter 4 (Discussion)* interprets the findings, discussing how this approach addresses the coordination problem, examining limitations, and considering counterarguments (e.g. concerns about oversimplification or misuse of formal models).
- *Chapter 5 (Conclusion)* summarizes the contributions, reflects on implications for AI governance, acknowledges remaining challenges, and suggests directions for future research.
- [] *Ensure the introduction clearly motivates the problem and states the thesis, providing a compelling motivation for the reader (Rubric: Introduction). Include a concise roadmap for structure clarity.*

2. Context & Background

2.1 Theoretical Foundations

- **Carlsmith's existential risk model**: Introduce **Joseph Carlsmith's 2022 report** as a landmark structured analysis of AI x-risk. Explain how Carlsmith breaks down the probability of AI-driven catastrophe into **six key premises**, each with an estimated probability ¹. These premises cover stages like AI development, misalignment likelihood, deployment decisions, etc., which multiply into an overall ~5% chance of doom.
- Detail *each premise* and Carlsmith's original estimates (e.g. probability advanced AI is misaligned, probability of power-seeking behavior if misaligned, etc.), conveying the logic of his model. Emphasize how this approach brings clarity by quantifying uncertain beliefs in a chain.
- **Composite risk calculation**: Show how multiplying the premises yields ~5% existential risk ², and note Carlsmith's interpretation of this figure. Explain that while the absolute number is debated, the structured approach itself is a major contribution, setting a precedent for rigor in AI risk assessment.
- **Significance**: Highlight why this model is pivotal: it provides a *transparent, modular framework* for discussing P(doom) that others can critique or modify. This makes it an ideal candidate for formal modeling—its logical structure is already explicit ³. By formally modeling it, we can test the coherence of the argument and explore variations (e.g., what if one premise changes?).
- **Formalization potential**: Argue that Carlsmith's model, being one of the most explicit quantitative arguments in AI safety, is "low-hanging fruit" for tools like AMTAIR. If we can automate transforming such structured arguments into a Bayesian network, it validates our approach and could handle other, less explicit arguments as well.
- **CODE EXAMPLE**: *Diagram or pseudo-code illustrating Carlsmith's probability chain.* For instance, show a simple product formula of the six premises or a small tree diagram where leaves are premises and root is the final probability. *(This will help readers visually grasp how the ~5% is obtained.)*
- (Carlsmith 2022 on P(doom) ~5% overall risk)

- **Bayesian networks (BN) as knowledge representations:** Provide a gentle introduction to Bayesian networks as the theoretical backbone for probabilistic reasoning. Define a BN formally: a **Directed Acyclic Graph (DAG)** where nodes = random variables and directed edges = conditional dependencies ⁴. Explain the components:
- **Nodes and states:** Each node represents a proposition or variable (e.g. “*Grass is wet*” or “*AI is misaligned*”) and can take on certain states (T/F or more complex states).
- **Edges and conditional probabilities:** An edge from A to B means A is assumed to directly influence B. With each node comes a **Conditional Probability Table (CPT)** specifying $P(\text{Node} | \text{Parents})$ for all parent state combinations ⁴. Clarify that no edge implies conditional independence, a crucial assumption derived from the argument’s structure ⁵.
- **Inference:** Outline how once a BN is built, we can compute any **posterior probability** given evidence using algorithms like variable elimination. For example, we can ask $P(\text{Rain}=T \mid \text{GrassWet}=T)$ in a weather network ⁶. This is the kind of reasoning we want to automate for AI risk models.
- **Rain–Sprinkler–Lawn example:** Introduce the *canonical* BN example of rain causing wet grass, possibly confounded by a sprinkler ⁷. Use this simple three-node network to illustrate BN concepts:
 - Nodes: *Rain*, *Sprinkler*, *Grass wet*. Edges: $\text{Rain} \rightarrow \text{GrassWet}$, $\text{Sprinkler} \rightarrow \text{GrassWet}$ (and perhaps $\text{Rain} \rightarrow \text{Sprinkler}$ if we assume rain might stop sprinkler usage, depending on the variant).
 - Provide a small CPT example: e.g. $P(\text{GrassWet}=T | \text{Rain}=T, \text{Sprinkler}=T) = 0.99$ (grass almost certainly wet if both causes happen) ⁸, etc., and $P(\text{GrassWet}=T | \text{Rain}=F, \text{Sprinkler}=F) = 0$ (dry if no rain or sprinkler). These illustrate how probabilistic assumptions are encoded.
 - *Include a diagram* showing this network structure (three circles with arrows) to visually reinforce how a DAG looks. Also mention how we can compute outcomes: e.g. if grass is wet, a BN can update belief in whether it rained.
 - Discuss **cognitive advantage**: Because BNs force us to specify relationships clearly and numerically, they prevent certain logical fallacies and highlight assumptions. They are excellent for **uncertainty management**, as they allow all pieces of evidence to be consistently combined via probability theory. This is highly relevant to AI risk, where uncertainty is large.
- **Application to AI risk:** Connect BN theory back to our domain: Carlsmith’s premises can be seen as nodes in a causal chain (e.g. *AI developed*, *AI misaligned*, *AI escapes oversight*, *causes catastrophe*). A BN can capture these dependencies and their probabilities. By formalizing an AI risk argument as a BN, we can rigorously ask questions like “*Which premise’s uncertainty contributes most to overall risk?*” or “*How would improved oversight (a parent node) reduce the probability of catastrophe?*”. This justifies using BNs for modeling transformative AI risks – they can mirror the logical structure of arguments and quantify them.
- **CODE EXAMPLE:** *Simple Python snippet for the Rain–Sprinkler BN.* For instance, show pseudo-code using a library (like `pgmpy`) to define the structure and CPTs for the three nodes, then perform an inference query. This grounds the abstract concept in a tangible form (but without getting lost in code details). It demonstrates how easily a BN can answer queries once the structure and probabilities are set.
- **Epistemic challenges in policy evaluation:** Analyze why evaluating AI policies (e.g. a moratorium on advanced AI development) is especially hard without formal models. Identify challenges:

- **Deep uncertainty:** Many AI policy decisions must be made under unprecedented uncertainty (no historical frequency for AI singularity scenarios). Traditional policy analysis tools struggle when probabilities are subjective or expert-elicited.
- **Complex causality:** Policies can have long causal chains (e.g. funding interpretability → better alignment techniques → reduced accident risk), which are hard to intuit. Without a structured model, important indirect effects or feedback loops might be overlooked.
- **Multidisciplinary inputs:** AI governance combines technical facts, ethical considerations, and strategic behavior. Informal reasoning might miss how these interact. We need a representation that different experts can all contribute to, and that can integrate qualitative judgments with quantitative analysis.
- **Existing approaches fall short:** Note that typical approaches (like purely qualitative scenario analysis or simplistic quantitative risk models) either lack rigor or flexibility. *For example*, qualitative arguments are rich but hard to combine or compare systematically; conversely, pure quantitative models (like econometric models) require simplifying assumptions that may ignore crucial qualitative insights. This gap highlights the need for a **hybrid approach** – precisely what AMTAIR aims to provide by marrying argument mapping with probabilistic modeling.
- Emphasize that without improvement in our epistemic tools, policy debates remain mired in vague analogies or talking past one another. This sets the stage for why the methods in the next section (argument mapping, etc.) are necessary.
- [] *Demonstrate research depth by integrating foundational concepts (Bayesian inference, risk modeling) and showing understanding of core literature (Rubric: Background). Use this section to convince readers that you grasp both the philosophical context (AI x-risk debates) and the technical tools (BNs), establishing credibility.*

2.2 Methodology

- **Argument mapping and structured representations:** Present **argument maps** as a bridge between natural-language debates and formal models. Introduce *ArgDown*, a Markdown-like syntax for capturing argument structure (claims, premises, supports/attacks) in text form.
- Explain the basic elements of ArgDown: statements (claims or premises) and relationships (support or conflict). For example, a claim “*AI will not be controlled*” might be supported by a premise “*Advanced AI can bypass all known security measures*”. ArgDown uses simple indentations or markers to denote such structure.
- **ArgDown example:** Use the Rain–Sprinkler context to give a trivial ArgDown snippet. E.g., in ArgDown:

```
(1) If it rains, the grass gets wet.
(2) It rained.
(3) Therefore, the grass is wet.
```

This shows a conclusion (3) supported by premise (1) and fact (2). Such a format makes the logical structure explicit.

- *(Figure/Diagram:)* Possibly show a small argument map diagram corresponding to the ArgDown example above (nodes for statements, arrows for inference), illustrating how text translates to a graph.

- Emphasize how **ArgDown helps clarity**: it forces authors to lay out assumptions and inferences clearly, which is invaluable for complicated AI risk arguments. However, ArgDown on its own is qualitative; this leads to the next step.
- **CODE EXAMPLE**: Basic ArgDown syntax highlighting hierarchical structure and references. (For instance, show how a premise and conclusion would be written and linked in ArgDown format in a monospace block.)
- **Extending to BayesDown**: Introduce **BayesDown** as a custom extension of ArgDown developed in this project (or in related work) that adds probabilistic information to argument maps. Describe how BayesDown builds on an ArgDown structure by annotating nodes with:
 - **Instantiations**: the possible states of a proposition (e.g. TRUE/FALSE, or categorical levels).
 - **Priors**: a priori probabilities for each instantiation (if the node has no parents in the argument graph).
 - **Posteriors (conditional probabilities)**: probabilities of the node's instantiations given various parent configurations ⁹. Essentially, BayesDown allows one to attach a mini-CPT to each argument link.
- Explain with a simple example: take the statement “grass is wet” supported by “rain” and “sprinkler.” In BayesDown, we would specify something like:
 - Instantiations: grass_wet_TRUE, grass_wet_FALSE.
 - Prior: $P(\text{grass_wet_TRUE}) = 0.0$ (if no causes).
 - Posterior: $P(\text{grass_wet_TRUE} \mid \text{rain_TRUE}, \text{sprinkler_TRUE}) = 0.99$, etc. ⁸. This turns the argument “*Rain or sprinkler cause wet grass*” into a quantitative rule set.
- Emphasize **design principles**: BayesDown aims to remain **human-readable** (so domain experts can write or verify it) while being **machine-parsable** into a BN. It's a hybrid format bridging the gap between prose and mathematics.
- **CODE EXAMPLE**: Show a snippet of BayesDown syntax for the Rain–Sprinkler–Lawn example. For instance:

```
[GrassWet]: Grass is wet. {"instantiations": ["TRUE","FALSE"],
  "priors": {"P(TRUE)": 0.1},
  "posteriors": {"P(TRUE | Rain=TRUE, Sprinkler=TRUE)": 0.99, ...} }
```

This demonstrates how probabilities are embedded directly in the structured argument text.

- **Two-stage extraction process**: Outline the core methodology of the AMTAIR pipeline, which operates in two distinct stages:
- **Stage 1 – Structure Extraction**: Use natural language processing to convert raw texts or PDFs from AI safety literature into **ArgDown** format (i.e., identify arguments, conclusions, premises, their support/attack relations). Describe how this might be done:
 - The process involves parsing the text for argumentative structure. We leverage a combination of rules and **AI (LLM) assistance** to identify claims and premises. For instance, a function `parse_markdown_hierarchy()` reads an indented outline of arguments and produces a structured representation ¹⁰. If using an LLM, it could be prompted to output ArgDown given a section of a paper.
 - Mention any specific tools: perhaps a custom parser or fine-tuned model is used to extract argument graphs. Key challenges include correctly capturing logical relationships and referencing evidence.
 - The output of Stage 1 is an ArgDown (or intermediate graph) representation of the document's arguments, which is then easily transformable.

- **CODE EXAMPLE:** Reference a key function (e.g. a Python function or pseudocode) responsible for ArgDown parsing. This could be a placeholder like `parse_arguments_to_argdown(text)` showcasing how stage 1 is implemented in code.
- *Visualization:* Indicate that we can log or visualize this extracted argument structure (e.g., as a tree or graph) to verify that the narrative has been correctly captured before adding probabilities.
- **Stage 2 – Probability Integration:** Take the structured argument and augment it with probabilistic data to create BayesDown.
 - **Process details:** Explain that this stage attaches numbers to the qualitative skeleton. Sometimes the source text (like Carlsmith’s report) contains explicit probabilities – these can be pulled in automatically (e.g., regex for percentages or a structured data table in the text). In other cases, domain experts or external sources provide the priors and likelihoods.
 - **Question generation:** If probabilities aren’t directly given, the system or analyst might generate targeted questions for experts or a language model: e.g. *“Given premise A and B, what is the likelihood of conclusion C?”*. This step ensures every relationship in the ArgDown map gets quantified.
 - Use Carlsmith’s model as illustration: Stage 2 would extract the “~5%” from the text for the final node, and the individual premise probabilities (like 50%, 40%, etc.) for each node, plugging them into BayesDown format.
 - **CODE EXAMPLE:** Highlight a function or method that adds probabilities, e.g. `augment_with_probabilities(argdown) -> bayesdown`. Show a snippet where an ArgDown node is updated with a prior or a conditional probability entry.
 - *Visualization:* Mention that after Stage 2, we have a fully specified BayesDown (essentially a BN in text form), which can be visualized or converted to a standard `.bn` or `.net` format. For example, one could generate a graph visualization at this point to inspect the model structure with probabilities annotated.
- **Previous work (MTAIR) and rationale for automation:** Briefly describe the **Modeling Transformative AI Risks (MTAIR)** initiative (if any precedent exists) or the general state-of-the-art in AI risk modeling prior to this work:
- **Key innovations in prior work:** The MTAIR project (as referenced) attempted to manually map out AI risk scenarios and perhaps build models (e.g., using Analytica or causal diagrams). Note any contributions, such as structured frameworks or community engagement it achieved ¹¹.
- **Limitations of manual approach:** Point out why a manual, non-automated approach hits a wall: it’s **labor-intensive**, slow to update, and potentially inconsistent if different people do it. In Carlsmith’s case, one person structured one argument – but scaling that to dozens of papers or dynamically updating it as knowledge evolves is impractical. Also, manual models may not be easily reproducible or sharable in detail.
- These limitations directly motivate **automation**: by using AI to assist, we can dramatically speed up model construction, ensure consistency, and update models as new data arrives. Automation also opens the door to non-experts contributing (with the machine handling complexity under the hood).
- Summarize that our methodology leverages cutting-edge NLP to overcome these limitations, representing a step-change in how such risk models can be built and maintained.
- [] *Ensure the methodology section clearly explains how the research is conducted (Rubric: Methodology/Structure). Define the innovative format (BayesDown) and pipeline so that the approach is transparent. Highlight the originality of using an LLM-driven pipeline for argument-to-BN conversion (Rubric: Originality).*

3. AMTAIR Implementation

3.1 Software Implementation (Including Colab Notebook)

- **System architecture overview:** Describe the overall architecture of the AMTAIR prototype system and how data flows through it. Outline the main components and their interactions ¹² :
- *Text ingestion & preprocessing:* Tools to input AI safety documents (PDFs, text) and clean/structure them for analysis.
- *LLM-powered extraction pipeline:* Module that performs Stage 1 (argument extraction) and Stage 2 (probability integration), possibly calling external APIs or using local models. It takes raw or lightly structured text and produces ArgDown/BayesDown.
- *Bayesian network construction:* A component (using a library like `pgmpy` or similar) that reads the BayesDown data and instantiates a BN object (nodes, edges, CPTs in code).
- *Visualization interface:* A sub-system to visualize the constructed Bayesian network for the user. This could be an interactive HTML (e.g. using d3.js or pyvis network graph) embedded in a Jupyter/Colab notebook.
- *Analysis & inference engine:* Functions that perform probabilistic inference, sensitivity analysis, etc., on the BN. Also includes routines for exporting results (figures, data) and archiving the model.
- Present a **data flow diagram** (in narrative or figure) showing how an input document moves through these components to yield an interactive model. For example: PDF/Text → (ArgDown extractor) → ArgDown → (Probability integrator) → BayesDown → (BN constructor) → Bayesian Network → (Visualization/Analysis).
- **Implementation technologies:** List the key technologies: Python is used as the primary language; libraries such as `pgmpy` for Bayesian networks, `networkx` or visualization libs for graphs, possibly `pandas` for intermediate data frames, and an LLM API (like OpenAI GPT-4) or NLP library for parsing text. Also mention the use of **Jupyter/Colab** environment for an interactive demo and rapid prototyping.
- **Design principles:** Note any important choices (e.g., modular design so each stage can be improved independently; using open formats like JSON or CSV to intermediate store the BayesDown data; ensuring transparency by logging intermediate outputs for verification).
- **CODE EXAMPLE:** Provide a high-level code outline or module list to illustrate the structure. For instance:

```
amtair/  
├─ extract.py    # Stage 1 extraction functions  
├─ quantify.py  # Stage 2 probability integration  
├─ build_bn.py   # BN construction utilities using pgmpy  
├─ visualize.py # Functions to generate network graphs  
└─ analyze.py   # Inference and analysis functions
```

This modular breakdown shows how the implementation was organized, aligning with the described architecture.

- **Pipeline demonstration with Rain-Sprinkler-Lawn:** Before tackling the complex AI risk model, the implementation was validated on the simple Rain-Sprinkler example introduced earlier ¹³ ¹⁴ . Walk the reader through how the pipeline handles this example step by step (this is mirrored in the provided Colab Notebook demonstration):

- **Example overview:** Reiterate that this toy example has a known correct model, so it's used to verify each stage of AMTAIR is working. The input is a short ArgDown description of the rain-sprinkler situation (maybe prepared manually) along with some probabilities.
- **Stage 1 – ArgDown parsing:** Show how the system reads the ArgDown text for the example. It should detect nodes (Rain, Sprinkler, GrassWet) and their relations (Rain and Sprinkler both support GrassWet becoming true).
 - *Process explanation:* The extraction function processes the structured text lines, creates internal representations for each statement and link. In this simple case, it's straightforward since we hand-crafted the ArgDown.
 - **CODE EXAMPLE:** Insert a snippet of the ArgDown input for this example or how the parser code represents the Rain node. For instance:

```
(1) Rain -> (3) GrassWet
(2) Sprinkler -> (3) GrassWet
```

Then show the Python object or JSON output after parsing, e.g., a list of nodes and edges derived.

- Verify that the structure matches the expected DAG (two parents feeding into one child, no cycles).
- **Stage 2 – BayesDown enhancement:** Next, the system augments the extracted map with probabilities. For the rain-sprinkler example, these probabilities might be hardcoded or provided in a small JSON. The system assigns:
 - Priors like $P(\text{Rain}=\text{TRUE}) = \text{e.g. } 0.2$, $P(\text{Sprinkler}=\text{TRUE}) = 0.1$ (for instance), and
 - Conditional probabilities like $P(\text{GrassWet}=\text{TRUE} \mid \text{Rain, Sprinkler}) = 0.99$, etc.
 - *Process explanation:* Show that the pipeline either reads these from an extended input file or uses default values. In practice, the example data might be stored in a BayesDown markdown file that the notebook loads ¹⁵.
 - **CODE EXAMPLE:** Provide a snippet of the BayesDown representation after adding probabilities, similar to what was given in the methodology section but now as actually processed by the code. E.g.:

```
{
  "node": "GrassWet",
  "parents": ["Rain", "Sprinkler"],
  "P(true|Rain=true,Sprinkler=true)": 0.99,
  "P(true|Rain=true,Sprinkler=false)": 0.8, ...
}
```

(The exact format may differ, but illustrate that the quantitative info is now attached.)

- **Stage 3 – BN construction:** The BayesDown data is then fed into the BN builder. The implementation creates a `BayesianNetwork` object, adds nodes and edges, and populates CPTs for each node using the given probabilities.
 - *Process explanation:* Mention that we use a library for this (like creating a `pgmpy.DAG` and adding `TabularCPD` objects for each CPT). The outcome is a fully functional Bayesian network in memory.

- **CODE EXAMPLE:** Show a brief code snippet of constructing the BN (e.g., `model = BayesianNetwork([('Rain', 'GrassWet'), ('Sprinkler', 'GrassWet')])` followed by definitions of CPDs and adding them to the model).
- **Visual result:** Once built, the notebook produces a visualization. For example, it might generate a graph image or interactive widget. Indicate that a **figure of the rain-sprinkler network** is shown, confirming the structure (a small graph with the three nodes and arrows).
- If interactive, note features like clicking on nodes to see probabilities. (In the actual Colab, this is likely done with pyvis or a custom HTML.)
- **Inference demonstration:** Demonstrate that the BN works by querying something. For instance, calculate the probability of grass being wet given evidence:
 - Query: $P(\text{GrassWet}=T)$ (with no evidence) should equal the weighted combination of scenarios (and indeed the system can compute that – e.g., it might output ~ 0.27 if using sample probabilities).
 - Or query $P(\text{Rain}=T \mid \text{GrassWet}=T)$ to see how likely it is that it rained given wet grass (which should be higher than the prior for rain).
 - Show that the notebook runs these queries with a function (like `model.infer(query='GrassWet':True, evidence={})`).
 - The results should match manual calculation, giving confidence the pipeline preserves logical correctness.
- **Validation:** Highlight that the outcome from the automated pipeline for this example can be compared to textbook results or manual calculation, and they match. For instance, if we manually compute $P(\text{GrassWet}=T)$ with the assumed CPTs and the system's result is identical, that validates the implementation. This step ensures that before applying to complex arguments, the pipeline's mechanics are sound.
- **Interface and usage:** Mention that all these steps were executed in an interactive **Google Colab Notebook** (linked in the thesis), providing an accessible demonstration. The notebook is structured in numbered sections (as described in the context) ¹⁶, and a user can rerun it to see each stage's output, modify inputs, etc. This emphasizes reproducibility and transparency of the implementation.
- [] *Ensure technical clarity here: each part of the implementation should be explained without assuming too much prior knowledge (Rubric: Clarity). Use diagrams or pseudo-code to clarify complex processes. Also, underscore the original engineering effort: this custom pipeline is a key contribution (Rubric: Originality/Technical Achievement).*

3.2 Results

- **Rain-Sprinkler test results:** Summarize the results from the simple example as a baseline:
- The pipeline successfully reproduced the expected **Rain-Sprinkler-Lawn network** and basic inferences. For example, given the input probabilities, it computed $P(\text{GrassWet}=T) \approx 0.27$ (hypothetical) which is consistent with manual calculation. This sanity check builds confidence that the automation is logically correct.
- This result, while unsurprising, demonstrates that even a simple causal argument can be fully formalized and queried with minimal human intervention, illustrating the promise of the approach.
- (If any minor issues were encountered, like needing to tweak the syntax or CPT format, mention that they were resolved at this stage, ensuring the pipeline is robust for more complex input.)
- **Formalizing Carlsmith's AI risk model:** Present the outcomes of applying the entire AMTAIR pipeline to Carlsmith's six-premise existential risk model from Section 2. This is the primary result of the project:

- **Model construction:** The system processed Carlsmith's structured argument (provided in a semi-structured form, possibly manually curated ArgDown with his premises). It produced a Bayesian network capturing the causal chain of events leading to AI-driven catastrophe, with on the order of *dozens of nodes* (each premise, sub-premise, and intermediate factor becomes a node).
 - The BN includes Carlsmith's main nodes (e.g., *Advanced AI is developed by year X, It is misaligned, It causes an unrecoverable catastrophe*, etc.) and possibly additional sub-nodes where Carlsmith's text had sub-arguments or assumptions.
 - **Structural validation:** The resulting DAG reflects the logical structure of Carlsmith's essay. For instance, nodes corresponding to each premise feed into the node representing existential catastrophe, mirroring the multiplication of probabilities in the original model. If Carlsmith had conditional reasoning (like "If premise 1 and 2, then 3"), the network structure captures those dependencies appropriately. We ensure no spurious connections were added – every edge has a justification in the text.
 - (Figure:) Include or describe an illustrative **graph of the Carlsmith model BN**. This could be a complex graph visualization showing clusters of nodes (perhaps grouping by premise) with directed links. We will highlight a portion of it in the thesis for readability, but the full interactive version is available digitally. This figure demonstrates that a lengthy prose argument can indeed be converted into a rigorous graphical model.
- **Probability calibration:** The BN's parameters were set using Carlsmith's estimates:
 - Each of the six top-level premises has a prior equal to Carlsmith's probability for that premise being true. For example, if Carlsmith estimated a 20% chance that advanced AI will be misaligned, that becomes $P(\text{node_Misalignment} = \text{TRUE}) = 0.2$ (assuming a binary TRUE/FALSE node).
 - The relationships (edges) carry conditional probabilities reflecting Carlsmith's conditional claims. For instance, if he implied "If AI is developed and misaligned, there is a 40% chance of catastrophe," that becomes a conditional probability in the network.
 - Where Carlsmith multiplied independent probabilities, the network treats those nodes as independent parents of the outcome node, reproducing the multiplication in a graphical way.
 - **Resulting P(doom):** We confirm that the BN yields approximately the same final probability of existential catastrophe as Carlsmith's original calculation. Indeed, when we set all premises as per his scenario, the network's output for $P(\text{Catastrophe} = \text{T})$ was **~4.98%**, essentially reproducing the ~5% figure ¹⁷. This is a crucial validation: it shows that our formal representation has not diverged from the source material in aggregate.
 - (This close match with Carlsmith's ~5% demonstrates fidelity of the model to the original argument.)
 - **CODE EXAMPLE:** (If appropriate) show a snippet of querying the final network for the probability of catastrophe with all premises "on." For example, a pseudo-call like `P_catastrophe = model.predict_proba({Premise1:True, Premise2:True, ...})` yielding ~0.05. This would illustrate how one interacts with the final model to get results.
- **Structural insights:** With the model built, we can analyze its structure for insights that are harder to see in the prose:
 - Compute graph metrics: e.g., identify which nodes have the highest **in-degree** (they rely on many factors) or **out-degree** (they influence many other nodes). In Carlsmith's model, perhaps one premise has multiple sub-premises feeding into it, indicating a particularly complex part of the argument.

- Identify **key dependencies**: The model might show, for example, that *Misalignment* and *Power-Seeking tendency* both have to be true to get to catastrophe, creating a conjunctive bottleneck. Recognizing such “AND” nodes (where multiple parents must align) is useful for strategy – it might hint that mitigating one factor alone isn’t enough.
- Check for any independencies: The BN format might reveal that certain premises were effectively independent in Carlsmith’s reasoning (no overlapping nodes), which validates his multiplication method. Alternatively, if our model introduced any linkage Carlsmith didn’t consider, that’s a spot for discussion or refinement.
- *(These structural results will be further interpreted in the Discussion, but they are listed here as outcomes of the implementation exercise.)*
- **Quantitative analysis**: Beyond replicating Carlsmith’s final number, the BN allows more granular results:
 - **Sensitivity analysis**: We varied each premise’s probability to see which one influences the final outcome most. For instance, if increasing the probability of *misalignment* from 20% to 30% raises P(doom) significantly more than a similar increase in *deployment without oversight*, that tells us misalignment is the more sensitive parameter. Early results indicate one or two premises (to be identified in text) have outsized impact on P(doom), suggesting where reducing uncertainty would most change the bottom line.
 - **Posterior updates**: If new evidence came in (say a breakthrough in alignment techniques making misalignment less likely), the model can instantly recompute P(doom). We demonstrate this by adjusting one premise probability within the BN and observing the new outcome (perhaps showing a drop to ~3% if misalignment risk is halved, as an illustrative example).
 - **Policy scenario testing**: As a preliminary result, we show how one could incorporate a policy effect. For example, introduce a hypothetical node “*Global safety standards enforced*” which directly reduces the probability of catastrophe (or of some intermediate node). By toggling this node (policy on vs off) in the model, we can quantify the difference in existential risk (e.g., with standards, P(doom) might go down to 3%, without it stays at 5%). This showcases the model’s capability to evaluate interventions, though these numbers are only as good as the assumptions provided.
 - **CODE EXAMPLE**: Provide pseudo-code for a sensitivity calculation: e.g., loop through premises, nudging their probability, and logging outcome probability – to illustrate how one would programmatically derive the sensitivity results. Or show a snippet for adding a policy node and computing the new risk.
- **Performance and reliability**: Note how the system performed with this real-world example:
 - **Execution time**: The extraction and model-building for Carlsmith’s text (which is lengthy) was on the order of minutes. The biggest time sink was the LLM parsing step for ArgDown. Once the structure was obtained, computing the BN was fast (millisecond-scale for inference queries). This suggests the approach is feasible for documents of similar complexity, though scaling to dozens of documents would multiply the extraction time.
 - **Accuracy of extraction**: Since Carlsmith’s text was structured, the automated extraction had a high success rate. We cross-checked the extracted ArgDown against the original text and found it captured the key premises and their logical flow correctly. A few minor points (like nuanced caveats or conditional statements) required manual clarification, pointing to areas to improve in the NLP step.
 - **Expert validation**: (If applicable) mention that the formal model was reviewed by an AI safety expert or the thesis advisor to ensure it correctly represents Carlsmith’s argument. Their

feedback confirmed that the model aligns with the intended interpretation of each premise, lending credibility to the formalization.

- **Limitations encountered:** Acknowledge any issues: e.g., “Premise 4 was not explicitly quantified by Carlsmith, so we had to assume a probability for it” – this introduces some uncertainty. Or if some dependencies were ambiguous in text, we had to make a modeling judgment (noting that these decisions will be discussed later as potential limitations).
- [] *In presenting results, maintain clarity: use figures or tables to summarize key results (like sensitivity rankings) so the reader isn’t lost in numbers (Rubric: Clarity). Also, connect results back to the thesis motivation – e.g., show that having these results exemplifies why the approach is useful (improving insight into AI risk). Ensure the depth of analysis is evident by not just reporting what the model did, but extracting meaningful patterns or confirmations (Rubric: Research Depth).*

4. Discussion

- **Addressing the coordination problem:** Begin by reflecting on how the AMTAIR approach, as demonstrated, contributes to solving the AI governance coordination issues outlined in the introduction.
- **Shared language:** Now that arguments like Carlsmith’s can be turned into BNs, diverse stakeholders can **literally see** the same model of the problem. Discuss how an interactive BN (with visual nodes and adjustable assumptions) can act as a common reference. For example, a policy-maker might not follow all the technical jargon in AI research papers, but they can interact with a risk model and ask “what if” questions, making the discourse more concrete.
- **Bridging qualitative and quantitative:** Emphasize the philosophical point that AMTAIR bridges *interpretive, narrative reasoning* (favored in philosophy and ethics) with *quantitative, predictive reasoning* (favored in technical fields). This hybrid epistemic approach is well-suited to AI governance, which requires both clarity of assumptions (a philosophical strength) and rigor of probabilistic prediction (an engineering strength).
- **Epistemic infrastructure:** Connect to the broader idea of creating an “**epistemic infrastructure**” for AI governance ¹⁸. By systematically formalizing arguments, we lay groundwork for a knowledge base of AI risks and interventions that can be built upon collaboratively. This can improve information flows between academia, industry, and government, helping align efforts.
- **Motivation reemphasized:** Tie back to the motivating example from the intro – now one can see how a miscommunication or disagreement among stakeholders (e.g., over the likelihood of misalignment) could be resolved by examining the model and adjusting that parameter to see implications. The discussion should make it clear that the project’s value lies not just in one model, but in introducing a new **paradigm for discourse** (from debates to collaborative modeling).
- **Interpretation of key findings:** Discuss what we learned from the results in Chapter 3, beyond the numbers:
 - Carlsmith’s ~5% risk being reproduced by the BN lends confidence to his breakdown, but also validates our tool. It shows that careful qualitative reasoning and quantitative modeling are consistent, which is encouraging for using formal methods in this domain.
 - The **sensitivity analysis** result (mention which premise was most impactful) has an important implication: it suggests priority areas for research or policy. For instance, if the premise “Misaligned AI will be deployed” drives the risk most, then governance efforts should perhaps focus on ensuring unsafe AI is never deployed (through regulation or norms).
 - Unexpected insight: maybe the formal model revealed a subtle point, like two premises were almost redundant (correlated) or one premise’s effect was mitigated by another. These are the kinds of

discussions to surface: how formal modeling can reveal when our intuitions about an argument might be incomplete or when certain risk factors are overlapping.

- **Policy leverage:** Interpret the policy scenario example: If adding “Global safety standards” node dropped risk by X%, discuss what that suggests. It might mean that international coordination on safety could be extremely valuable (if it significantly lowers multiple risk pathways simultaneously). This ties the model back to real governance implications, showing how the formal approach can inform strategy (which traditional narrative analysis might not quantify).
- Reiterate that these findings are illustrative and depend on model assumptions, but the point is we now have a **tool to generate such insights** systematically.
- **Limitations of the approach:** Provide a critical evaluation of AMTAIR's current limitations, to honestly appraise the work and identify areas for caution:
 - **Formalization oversimplification:** A major concern is that by forcing complex reality into a structured model, we might oversimplify. Discuss that many nuances in Carlsmith's reasoning (and AI risk in general) may not fit neatly into binary nodes or static probabilities. *Counterpoint:* However, any model is a simplification, and the goal is to simplify *wisely*. By preserving a lot of structure (via ArgDown) and explicitly noting uncertainties, we arguably lose less nuance than ad-hoc mental models do.
 - **Objection:** “Formal models can’t capture the full complexity of socio-technical systems; relying on them could give a false sense of security.”
 - **Response:** Acknowledge the simplification but argue that formal models are complements to human intuition, not replacements. They enforce consistency and highlight assumptions, which actually *prevents* some forms of oversimplification (like ignoring a factor entirely). The key is to continuously refine models and combine them with qualitative insights.
- **False precision and uncertainty:** With a numeric model, there is a temptation to take the output (e.g. 4.98%) too seriously. The numbers might project an illusion of precision.
 - **Objection:** “Attaching exact probabilities to one-time events (like an AI catastrophe) is speculative; doing so might engender overconfidence in those numbers.”
 - **Response:** Emphasize that the model explicitly shows uncertainty ranges and can be used to perform **uncertainty analysis**. Rather than saying “it’s 5%”, one should say “it’s 5% given these premises and their probabilities.” This actually fosters epistemic humility: it breaks the problem into parts so we can discuss which parts we’re uncertain about. Also mention that we could incorporate probability distributions or confidence intervals in future, instead of point estimates, to better represent uncertainty.

- **Stakeholder inclusion and transparency:** Another limitation is the barrier to entry – these models could be opaque or intimidating to non-technical stakeholders, potentially **excluding voices** from the discussion.

- **Objection:** “Turning policy debates into complex graphs and equations might sideline those without technical training, concentrating influence in the hands of modelers.”
- **Response:** This is a valid concern; however, the project has actively tried to mitigate it by focusing on **visual and interactive** outputs. The idea is not to hide arguments in math, but to make them navigable: e.g., an interactive interface where anyone can toggle assumptions and see effects can demystify the analysis. Also, because the source format (ArgDown/BayesDown) is text-based and relatively straightforward, stakeholders could be taught to read or even write parts of it. The goal is a **tool for communication**, not a black box. We also propose involving domain experts and diverse stakeholders in the model-building process, so the resulting network reflects a blend of perspectives.

- **Scalability and feasibility:** Discuss practical challenges in scaling this approach up:

- Technical: If we tried to process an entire library of AI safety literature, would the system handle it? The LLM-based extraction might become costly or require significant compute. There’s also maintenance – as new arguments arise, models must be updated.
- Institutional: Would policymakers and researchers actually use such a system? There might be resistance to adopting a new workflow, or skepticism of its outputs until it’s proven.
- Data availability: Some arguments don’t have any numbers to begin with, so an automated system has limits without expert input. There’s also the risk of **Garbage In, Garbage Out** – if source estimates are wildly speculative, the model’s output is equally speculative.
- **Objection:** “This is a neat academic exercise, but integrating it into real policy decision-making may not be realistic given resource and complexity constraints.”
- **Response:** Outline an incremental adoption scenario: we don’t expect governments to take Bayesian nets from day one to decide policy. Instead, researchers can start by using AMTAIR internally to test assumptions and clarify debates. Over time, as models prove useful (e.g., by identifying a previously overlooked risk factor or resolving a contentious point empirically), confidence in the approach will grow. Also, as AI itself advances, tools like these can be made more user-friendly and integrated into standard analysis pipelines (just as data science tools became common in policymaking over the last decades).

- We propose developing **user-friendly interfaces** and training programs for policy analysts to become comfortable with these models. Also emphasize that results should complement, not replace, narrative reports; one can always accompany a BN analysis with a natural-language explanation of the insights.

- **Counterarguments and rebuttals:** (Integrated above as objections/responses.) Summarize the key philosophical counterpoints to our thesis and our replies in a concise way:

- Formal models vs. complexity (oversimplification vs. necessary abstraction),
- Numeric uncertainty vs. humility (false precision vs. explicit uncertainty),
- Technocracy vs. inclusion (expert-driven models vs. accessible visualization),
- Feasibility (idealized solution vs. incremental implementation).
- Show that we have thought deeply about each and have either solutions or at least a path to mitigating these issues.
- **Alignment with Philosophy & Economics perspective:** Reflect on how this work sits at the intersection of philosophy and economics:
 - Philosophical: it deals with epistemology (how we know and reason about uncertain future events), ethics (weighting existential risks), and logic (argument structure). By formalizing arguments, we contribute to philosophical clarity and rigorous critical thinking.
 - Economics/policy: it's about making rational decisions under uncertainty – a classic economic problem (like cost-benefit analysis but for unquantifiable scenarios). Our model is essentially enabling a form of expected utility reasoning on existential risk, which is a very *economics meets ethics* kind of problem (like Pascal's wager but in practical terms).
 - This integrated approach showcases the *P&E program's values*: normative insight combined with analytical rigor. We should explicitly note that we're blending qualitative normative concerns (existential safety as an imperative) with quantitative decision tools (Bayesian networks, which are used in economics for decision analysis).
 - Possibly add that originality of this thesis lies in applying a tool (BNs) typically used in economics/engineering to a domain of philosophical significance (the fate of humanity, ethical risk). This cross-pollination is an innovative strength.
 - [] *In the discussion, ensure argumentation is strong (Rubric: Argumentation). We're not just listing results; we're building a case that the method is valuable despite limitations. Use the counter-objection structure to demonstrate critical engagement with potential criticisms (Rubric: Counterclaim & Rebuttal). Maintain clarity by structuring this section with subpoints or bolded objection/rebuttal labels, to guide the reader through the argument.*

5. Conclusion

- **Summary of contributions:** Recap the thesis in a synthesized way:
 - Restate the **core thesis**: that automating the modeling of AI risk arguments is both feasible and useful. Summarize how we showed this by actually implementing the pipeline and applying it to a significant case.
 - **Methodological innovation:** Highlight that we introduced or refined a novel methodology (ArgDown→BayesDown→BN) for knowledge representation in AI governance. This is a new contribution to the toolkit of both philosophers (argument mapping community) and policy modelers.
 - **Technical achievements:** List the concrete outputs: a working prototype software, a formal model of Carlsmith's analysis (which to our knowledge is the first of its kind), interactive visualizations, and analytical results (sensitivities, scenarios) that were previously not accessible.
 - **Insight generation:** Emphasize any key insights gained (e.g., identification of the most pivotal assumption in Carlsmith's model, or demonstration that certain interventions could dramatically lower risk in the model). These show that the approach doesn't just *work*, it also *added value* over the original purely narrative argument.

- **AI governance implications:** Summarize how these contributions matter for the broader goal: e.g., “By making expert assumptions explicit and computable, this approach can improve strategic planning for AI safety – a pressing need as the field advances.”
- Ensure the summary connects back to the Introduction’s promise, giving a satisfying sense that we addressed the initial problem (coordination failure via better modeling).
- **Limitations and open issues:** Acknowledge that this work is an early step and has limitations (some of which were discussed). Summarize the main ones:
- **Technical limitations:** The extraction is not fully autonomous – it may need supervision or fine-tuning for each new text. The BN models currently handle mostly binary or discretized variables; real scenarios might require continuous ranges or more complex state spaces. Also, our prototype might not scale well without optimization (e.g., parsing a 300-page book or integrating dozens of sources might be cumbersome now).
- **Conceptual limitations:** Our models are only as good as the input assumptions. If an important factor was omitted by the original author or by our interpretation, the model won’t include it. Moreover, not everything about AI risk is probabilistic or easily quantifiable (e.g., one might argue some outcomes are deeply unpredictable or qualitative); those aspects are not captured.
- **Scope limitations:** We focused on existential risk from misalignment. Other important facets of AI risk (like misuse, or systemic impacts) weren’t modeled. Also, we did not deeply engage with *ethical* dimensions like how to weigh future lives, etc., as our model was mostly factual/probabilistic. Those value-laden questions remain to be integrated with this framework.
- **Evaluation limitations:** We validated by reproducing Carlsmith’s result and logical structure, but we haven’t proven that an automated approach can handle *any* AI safety argument in the wild. A more comprehensive evaluation (taking multiple documents, getting experts to score the model quality) is future work.
- **Ethical considerations:** Note that increasing reliance on formal models could have downsides, such as over-reliance on numbers or potential misuse (imagine someone tweaking a model to justify a foregone conclusion). There’s a need for responsible use guidelines – which we haven’t fully fleshed out in this thesis.
- **Future research directions:**
- **Technical enhancements:** Outline how the pipeline could be improved. For example:
 - Use more advanced NLP or fine-tuned language models to improve ArgDown extraction accuracy, perhaps by training on a corpus of argument maps.
 - Extend BayesDown to handle **uncertainty ranges** or distributions rather than point probabilities, to incorporate confidence levels.
 - Integrate automated consistency checks, e.g., flag if the provided probabilities violate Bayes’ rule or if the network has unintended dependencies.
 - Scale up the visualization for very large graphs (maybe incorporating filtering or clustering so users can focus on parts of the model).
- **Integration with forecasting and data:** Suggest connecting these static expert models with dynamic data sources. For instance, linking nodes of the BN to real-time indicators or prediction market prices (where available) to keep the risk assessment up-to-date. This could move towards a live “AI risk dashboard” that updates as new information comes in.
- **Application to other arguments:** Propose applying AMTAIR to other substantive pieces of the AI risk literature or policy debates. E.g., model parts of *Nick Bostrom’s Superintelligence* arguments, or arguments about AI ethics (even if not existential risk). Each new application can reveal new requirements and enrich the methodology.
- **User study and refinement:** Future work could involve having actual policy-makers or researchers use the tool and gather feedback. How intuitive is it? Does it indeed improve their understanding or

decision-making? Their input would guide further development (for example, adding features or simplifying the interface).

- **Interdisciplinary research:** This project opens questions at the intersection of AI, philosophy, and governance. Future research might delve into *meta* questions like “How do formal models influence expert disagreement?” or “What is the epistemological status of a probability like 5% for an unprecedented event?” Such inquiry can deepen the theoretical grounding of AMTAIR.
- **Broader implications for AI governance:** End with a reflection on the big picture:
- **Improved strategic planning:** If widely adopted, approaches like AMTAIR could enable more **strategic foresight**. Governance discussions could shift from vague worries to model-guided exploration of scenarios. This might accelerate consensus on issues like which safety measures are most critical.
- **Institutionalizing modeling:** Imagine international bodies or think tanks maintaining open-source risk models that are continuously updated – similar to how climate models are used in climate policy. Our work could be a stepping stone toward that for AI. This would professionalize and systematize the analysis of AI risks.
- **Democratization of expertise:** If done right (interactive and transparent), these tools could allow a broader set of people to engage with complex risk arguments. For instance, students or journalists could play with a model to understand it rather than reading dense reports. This could lead to better public discourse on AI futures.
- **Ethical and normative considerations:** We should also recognize that formalizing decisions doesn't remove the need for value judgments. E.g., even if a model shows a policy reduces risk by X%, society must decide if the trade-offs (costs, restrictions) are worth it. Discuss how formal models fit into *normative decision-making*: they inform, but do not replace, ethical deliberation. It's important that the presence of a quantitative model doesn't overshadow qualitative values – instead, it should serve them by clarifying consequences.
- **Cautious optimism:** Conclude on a note that while transformative AI risks pose an unprecedented challenge, tools like AMTAIR provide a reason for optimism: we are not powerless to understand and shape these risks. By harnessing advanced AI (parsing arguments) to solve problems in AI governance, we turn the technology toward improving its own oversight. It's a meta approach that exemplifies the ingenuity needed to tackle emerging global risks.
- **Closing statement:** End the thesis by reasserting the thesis statement in past tense (“This thesis has shown that...”) and perhaps a call to action: encourage continued interdisciplinary work to refine these models and **embed them in real-world policy-making**. The final sentences might envision a future where AI risk assessments are as routine and trusted as financial models – a future where humanity is better equipped to navigate the transformative power of AI safely and coherently.
- [] *Ensure the conclusion is consistent with the introduction (Rubric: Conclusion). It should clearly refer back to the thesis question and claim victory in answering it, while also discussing broader implications and future work. Keep the tone optimistic but realistic, reinforcing the original contribution and inviting further inquiry. The conclusion should leave the reader with a clear understanding of what was accomplished and why it matters.*

</final_outline>

1 2 3 11 12 17 18 PY_Thesis_OutlineNDraft.md

https://github.com/VJMeyer/submission/blob/751d06f8e79384ebe2c1d699f3a25381e874e8c8/context/PY_Thesis_OutlineNDraft.md

4 5 6 9 10 **AMTAIR-FAQ_NotebookLM.md**

https://github.com/VJMeyer/submission/blob/751d06f8e79384ebe2c1d699f3a25381e874e8c8/context/AMTAIR-FAQ_NotebookLM.md

7 8 14 15 **AMTAIR_Prototype_example_carlsmithIPYNB.md**

https://github.com/VJMeyer/submission/blob/751d06f8e79384ebe2c1d699f3a25381e874e8c8/notebooks/AMTAIR_Prototype_example_carlsmithIPYNB.md

13 16 **AMTAIR_Prototype_example_rain-sprinkler-lawnIPYNB.md**

https://github.com/VJMeyer/submission/blob/751d06f8e79384ebe2c1d699f3a25381e874e8c8/AMTAIR_Prototype/data/example_rain-sprinkler-lawn/runtime_generated_data/AMTAIR_Prototype_example_rain-sprinkler-lawnIPYNB.md