

Chapter 8: Scalable Oversight

Markov Grey
Charbel-Raphaël Segerie

June 10, 2025

Contents

8.1	Introduction	3
8.2	Oversight	3
8.2.1	Training Signals	4
8.2.2	Verification vs. Generation	5
8.3	Task Decomposition	8
8.4	Process Oversight	9
8.4.1	Externalized Reasoning Oversight (ERO)	10
8.4.2	Procedural Cloning	14
8.5	Iterated Amplification	14
8.5.1	Amplification	15
8.5.2	Distillation	18
8.5.3	Iterated Distillation and Amplification (IDA)	19
8.6	Debate	20
8.6.1	Assumptions	24
8.6.2	Discriminator Critique Gap (DCG)	24
8.6.3	Judges	28
8.6.4	Truth	33
8.7	Weak-to-Strong Generalization (W2SG)	35
8.7.1	Sandwiching Evaluations	38
	Acknowledgements	39

8.1 Introduction

Oversight. As AI systems become increasingly capable, ensuring they remain aligned with human values and intentions becomes a critical challenge. This section introduces scalable oversight as a crucial approach to maintaining control over advanced AI. It explains the problems we face in generating training signals for complex, "fuzzy" tasks and the need for new methods to provide accurate feedback. This is important especially as AI models begin to perform tasks beyond human expertise. The section also explores the concept of verification being easier than generation, explaining why this property is fundamental to scalable oversight techniques.

Task Decomposition. Building on the need for better oversight methods, this next section explores task decomposition as one key strategy. Task decomposition involves breaking complex tasks into smaller, manageable subtasks, which can be recursively divided further. This approach helps in generating better training signals by simplifying the task that we need to evaluate and verify. Factored cognition extends this concept to replicate human thinking in machine learning (ML) models by decomposing reasoning, complex cognitive tasks.

Process Oversight. Another way to help scalable oversight is to address some of the limitations of outcome-based approaches. This section introduces the concept of process-based oversight. We explain Externalized Reasoning Oversight (ERO) and procedural cloning as specific examples. ERO techniques like chain-of-thought (CoT) encourage language models to "think out loud," making their reasoning processes transparent for better oversight and potentially preventing undesirable behaviors. Procedural cloning, an extension of behavioral cloning, aims to replicate not just the final actions but the entire decision-making process of experts. These methods offer a more principled approach to oversight by focusing on the AI's reasoning process rather than just its outputs.

Iterated Amplification (IA). Building on the concepts of task decomposition and process oversight, this section outlines amplification and distillation. Amplification enhances the abilities of overseers to solve more complex tasks, while distillation addresses the limitations of amplification, such as complexity and resource use. These processes are combined in Iterated Distillation and Amplification (IDA), a method aimed at generating progressively better training signals for tasks that are difficult to evaluate directly.

Debate. This section explores AI Safety via Debate as an adversarial technique for scalable oversight. It describes how AI models arguing for different positions, with a human or AI judge determining the winner, can result in more truthful outcomes. The potential of debate to elicit latent knowledge, improve reasoning, and enhance our ability to oversee complex AI systems is discussed. Key metrics such as the Discriminator Critique Gap (DCG) are introduced, along with the challenges of judging debates. The section also examines the assumptions required for Debate to converge on truth.

Weak-to-Strong (W2S). The final section introduces Weak-to-Strong Generalization (W2SG) as a practical approach to scalable oversight, building on insights from previous techniques. It explains how narrowly superhuman models can be used as case studies for scalable oversight techniques. W2SG involves training strong AI models using weak supervision, aiming for the strong model to outperform its weak supervisor by leveraging pre-existing knowledge. The section concludes by discussing various methods of evaluating oversight techniques, including sandwiching evaluations and meta-level adversarial evaluations, providing a way to judge future scalable oversight protocols.

8.2 Oversight

Why do we need oversight? As AI systems get smarter, they will start doing tasks that are hard for humans to evaluate. Evaluation means checking how well the AI did after completing a task, while feedback is the information we give to the AI during or after it works to help it learn and improve. Right now, we can still use methods like Reinforcement Learning from Human Feedback (RLHF) to guide AI in the right direction. But we can only give feedback if we can still evaluate the outputs. As tasks get more complex, even experts might struggle to provide accurate evaluations and feedback. So, we need new ways to give accurate feedback, even for tasks that are beyond human expertise. This is the goal of scalable oversight.

Scalable Oversight techniques help humans provide accurate feedback on tasks to ensure AI systems are aligned with our goals, even after the task complexity outstrips the ability of the best human experts. This can happen during the AI's training or deployment and isn't limited to RLHF-style feedback.

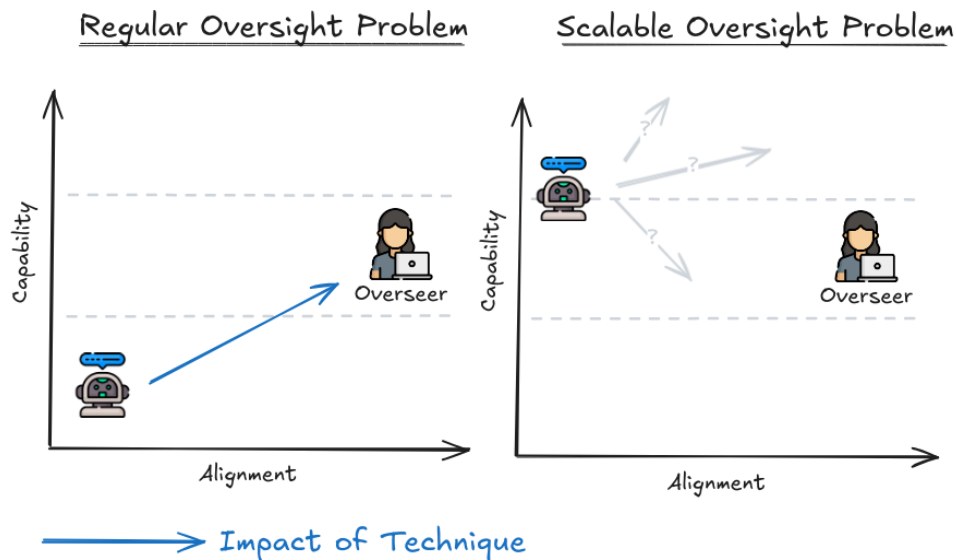


Figure 8.1: The difference between regular oversight safety research, and scalable oversight safety research.

Aligning RL Agents vs. LLMs. A few years ago it looked like the path to AGI was by training deep RL agents from scratch in a wide range of games and multi-agent environments. These agents would be aligned to maximizing simple score functions such as survival and winning games and wouldn't know much about human values. Aligning the resulting agents would require a lot of effort: not only do we have to create a human-aligned objective function from scratch, we'd likely also need to instill new capabilities into the agents like understanding human society, what humans care about, and how humans think. Large language models (LLMs) make this a lot easier: they come preloaded with a lot of humanity's knowledge, including detailed knowledge about human preferences and values. Out of the box they aren't agents who are trying to pursue their own goals in the world, and their objective functions are quite malleable. For example, they are surprisingly easy to train to behave more nicely. If AGI comes out of LLMs it might be easier to align. (Leike, 2022)

8.2.1 Training Signals

Before we understand how to actually align smarter than human AIs, we need to understand the general concept of training signals and why they are getting increasingly harder to generate as AI is starting to display higher levels of general purpose capabilities.

What are training signals? Training signals is a general term that we use for inputs that are used to guide AI learning. They can be rewards, labels, or evaluations indicating how well the AI is performing a task. For example:

- In supervised learning (SL), the training signal is the correct label for each example.
- For large language models (LLMs) using self-supervised learning (SSL), training signals are the correct next word in a sentence during pre-training. Human feedback on the quality of responses during fine-tuning is also a type of training signal guiding the outputs in a direction we want (a more "aligned" direction).
- In reinforcement learning (RL), training signals are rewards based on the success/failure of actions (or multiple actions), like points scored in a game or successful navigation to a location.

These signals shape how AI systems learn and are used to both evaluate performance and provide feedback.

Easy-to-generate training signals. For some tasks, generating training signals is simple. AlphaGo Zero, an RL agent playing Go, is a good example of this. The game has clear rules and win/loss outcomes, so training signals are straightforward: algorithmically generated win and lose signals directly measure performance, making it easy for the model to learn and improve its gameplay.

Hard-to-generate training signals. For other tasks, creating training signals is much harder. For example, training GPT models to generate accurate text summaries is challenging. The AI needs to convey correct information while being coherent and interesting. Success is subjective. Since it depends on individual reader preferences, it is hard to define clear, algorithmically generated training signals. Another example is self-driving cars navigating through busy city streets. These cars need to make real-time decisions, and the training signals or rewards for safe and efficient navigation are difficult to define due to varying contexts and sometimes conflicting traffic laws and safety considerations.

Fuzzy tasks. We call tasks where training signals are hard to generate “fuzzy tasks”. These tasks generally have ambiguous or ill-defined objectives and outcomes. We can’t generate precise training signals due to inherent subjectivity and variability in “correct responses”. Fuzzy tasks lack clear, objective criteria for success. Unlike well-defined tasks with specific, measurable goals, fuzzy tasks are more open-ended, which complicates our job of coming up with training signals. If it’s difficult to provide precise rewards or labels that accurately capture the desired behavior, the training process becomes complicated. AI systems might not receive the consistent, reliable feedback needed to learn effectively. This is essentially highlighting again the difficulty of the reward specification problem that we talked about in previous chapters.

Fuzzy tasks and Scalable Oversight. Fuzzy tasks are closely related to AI alignment, where ensuring AI systems act in accordance with human values and intentions is challenging due to ambiguity and subjectivity. Aligning AI with human values is a fuzzy task. Oversight techniques aim to solve alignment by providing training signals for fuzzy tasks, including feedback and imitation learning techniques like RLHF, Constitutional AI (CAI), and Inverse Reinforcement Learning (IRL). Scalable oversight techniques aim to provide training signals for fuzzy tasks that are too complex for even experts to understand or evaluate.

To make scalable oversight techniques viable verification needs to be easier than generation, and preferably (but not necessarily) tasks should be decomposable. These properties will be discussed in the next sections.

8.2.2 Verification vs. Generation

What Does $P \neq NP$ Mean? In computer science, we classify problems based on how hard they are to solve (generate a solution) and how hard they are to check (verify a solution).

- **P (Polynomial time):** These are problems that a computer can solve quickly.
- **NP (Nondeterministic Polynomial time):** These are problems where, if you have a solution, you can check it quickly, but finding the solution might take a long time.

What is generation? Generation is the process of coming up with solutions from scratch. This means searching through many possibilities, which can take a lot of time and computing power. For example, solving a Sudoku puzzle involves filling a 9x9 grid with numbers so that each row, column, and 3x3 subgrid contains all the digits from 1 to 9 without repeating. If you’ve ever tried to solve a Sudoku puzzle, you know it involves a lot of trial and error to make sure all the rules are followed.

5	3			7				
6			1	9	5			
	9	8					6	
8				6				3
4			8		3			1
7				2				6
	6					2	8	
			4	1	9			5
				8			7	9

Figure 8.2: Example of an empty sudoku board ([Wikipedia](#)). Filling this in would be harder than checking if the solved version below is correct.

Generation here involves filling in the blank grid while ensuring all the constraints (unique numbers in rows, columns, and subgrids) are satisfied.

What is verification? Verification is the process of checking whether a given solution attempt is correct. Using the Sudoku example, verification means making sure that each row, column, and subgrid contains all the digits from 1 to 9 without any repeats. Once someone gives you a completed Sudoku puzzle, checking whether it's correct is straightforward and quick. This idea is central to the concept of $P \neq NP$.

5	3	4	6	7	8	9	1	2
6	7	2	1	9	5	3	4	8
1	9	8	3	4	2	5	6	7
8	5	9	7	6	1	4	2	3
4	2	6	8	5	3	7	9	1
7	1	3	9	2	4	8	5	6
9	6	1	5	3	7	2	8	4
2	8	7	4	1	9	6	3	5
3	4	5	2	8	6	1	7	9

Figure 8.3: Example of a filled in sudoku board ([Wikipedia](#)). Its very easy comparatively to see that the solution is correct.

Examples: Illustrating Verification is Easier Than Generation. This is a very general property that holds across many domains:

- **Formal Problems:** In computational complexity theory, most computer scientists believe that $P \neq NP$ ([Wikipedia](#)), which means there are many problems where checking a solution is easier than finding one. This is seen in tasks like solving SAT problems or graph algorithms.
- **Sports and Games:** It's easier to look at a football scoreboard to see who is winning than it is to play the game well.
- **Consumer Products:** Comparing the quality of smartphones based on user reviews is simpler than designing and building a new smartphone.
- **Job Performance:** Evaluating how well an employee is doing is less demanding than actually doing the job.
- **Academic Research:** While reviewing research can be tough, it's still less work than producing new research.

Why Verification Being Easier Than Generation Matters for Scalable Oversight. This fact is crucial for scalable oversight because it allows us as human overseers to efficiently ensure the correctness and safety of outputs produced by complex systems without needing to fully understand or replicate the entire generation process. If $P \neq NP$ is true, it implies that we might be able to trust and delegate alignment research itself to AI models, because we can comparatively easily verify that their solutions work while they have to do the hard task of generating the solutions to alignment. Overall operating

under this assumption can make the task of aligning advanced AI systems seem more feasible. The next few paragraphs go into the debate of how valid this assumption is.

Verification in Adversarial Contexts. When verifying something in situations where someone might be actively trying to trick or attack you, the process becomes much harder. Put another way, if we have AIs that are deceptive, the problem becomes significantly trickier. For example, making sure software is secure against all possible attacks can be tougher than writing the software in the first place. An attacker only needs to find one security hole, but the person verifying must check everything to ensure there are no holes. This makes verification very challenging. Similarly creating a secure system in cryptography is hard, but proving that it's secure against all possible attacks is even more difficult. You need to consider every potential way someone might try to break the system, which is a huge task.

Easier than generation does not mean verification is trivial. Just because verification is theoretically easier than generation doesn't mean it's always easy in practice. For example, checking a complex mathematical proof can be very hard. Writing the proof takes creativity and deep understanding, but verifying it requires careful and detailed checking, which can be exhausting and prone to mistakes. In the case of software, writing secure software is challenging, but verifying that it's completely secure is even harder. Even though verifying a problem's solution might be easier than generating the solution, the process can still be very difficult and require significant effort and expertise.

Verification of Safety vs. Provable Alignment. In the event that we have to deal with superintelligent AI, just verifying its behavior might not be enough. Some researchers argue that we need to prove that the AI will always act in ways that align with human values. Verification means checking whether the AI behaves correctly in specific situations. Provable alignment means giving solid evidence that the AI will act correctly in every possible situation, even new and unexpected ones. This requires more than just checking—it needs formal methods and guarantees, which is extremely difficult.

Verification vs. Mathematical Proof. Verification involves checking if a specific solution is correct, usually through testing or inspection. A mathematical proof, on the other hand, is a rigorous logical argument that shows a statement is always true. For example, verifying a Sudoku solution checks if the given arrangement is correct, while a mathematical proof might show that any Sudoku puzzle with a certain number of clues always has a unique solution.

8.3 Task Decomposition

What is Task Decomposition? Task decomposition is the process of breaking down a complex task into smaller, more manageable subtasks. This technique makes it easier to tackle sophisticated problems by dividing them into simpler components that can be addressed independently. For example, if you need to summarize a book, you could break down this larger task into summarizing each chapter individually. Each chapter summary then contributes to the overall summary of the book.

Example:

- Task: Summarize the book.
- Decomposition: Summarize each chapter separately.

Decomposing tasks can be thought of as a method of overcoming complexity. Humans navigate the world's complexity by using layers of abstraction, where each layer hides most of the underlying details, allowing us to focus on manageable chunks of information. This ability is important because humans can only keep track of a few pieces of information in their mind simultaneously. Task decomposition helps us solve highly complex problems by breaking them down into simpler subproblems.

What is Recursive Task Decomposition? Recursive task decomposition extends the basic concept by breaking down each sub-task into even more granular subtasks. This iterative process continues until each task is simple enough to solve directly. Continuing with the book summarization example, recursive task decomposition would involve further breaking down each chapter summary into page summaries, and each page summary into paragraph summaries until the task is simple enough to evaluate directly.

Example:

- Task: Summarize the book.
- Decomposition: Summarize each chapter, Summarize each page within each chapter, Summarize each paragraph within each page.

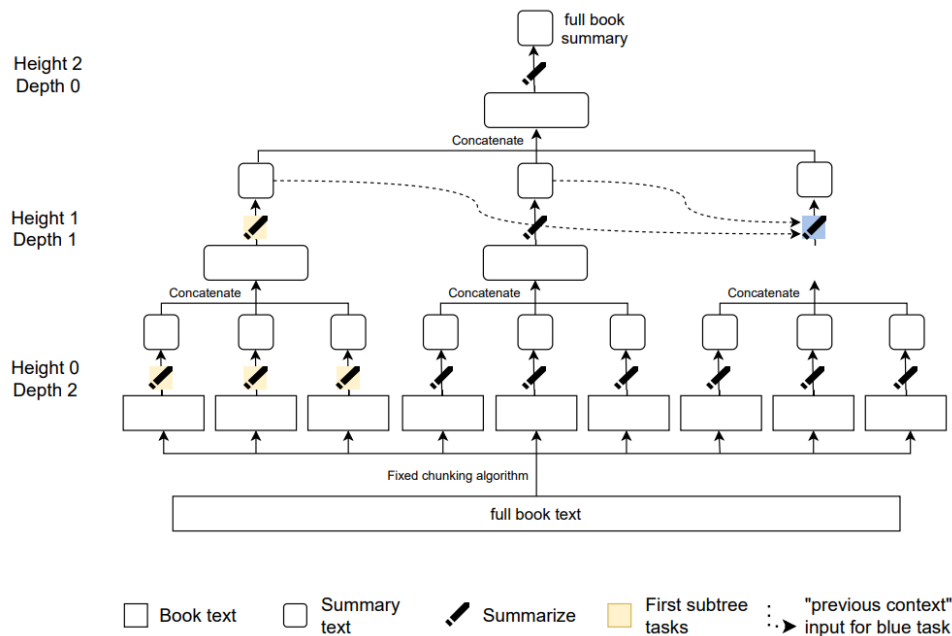


Figure 8.4: Example of summarizing books that combines task decomposition with learning from human feedback. The book is first decomposed into multiple chunks using a fixed (not learned) chunking algorithm (height 0). Then humans provide demonstrations summarizing these chunks, which is used to train an ML model on this data using behavior cloning. Then more data is collected from humans who compare different model outputs which is then used to further train a summarization policy using reward modeling. Then summaries are concatenated (height 0), data is collected for summarizing these summaries, and the model is fine-tuned for this summarization task (height 1). This procedure is repeated recursively until the entire book is summarized. (Wu et al., 2021)

How does task decomposition help generate better training signals? As we mentioned in the earlier sections, when AI systems become more capable, it becomes difficult for humans to provide the right training signals or data, especially for tasks that have subjective evaluation criteria (fuzzy tasks). The core thing that we are trying to do with task decomposition is to reduce the difficulty of providing a training signal to human judgment. Simpler tasks, means that it is easier to provide a training signal, it also means that the tasks are easier to verify. Therefore task decomposition is quite important to the success of many scalable oversight techniques.

Decomposing a task involves breaking it down into smaller, more manageable parts. These parts help you understand and manage the task better, but they may not always be independently solvable or reusable. The key properties we want from a good decomposition include:

1. **Recursive Decomposability:** The system should be capable of recursive task decomposition, thereby dividing a problem into simpler sub-tasks. For example, summarizing a book involves breaking it down into chapters, then pages, and then paragraphs until we reach some minimum level of task difficulty.
2. **Independence/Modularity:** Individual sub-tasks should be able to be completed independently and in parallel without relying on each other or the overall task. For example, each summarizer focuses on a paragraph without needing to understand the whole book.
3. **Composability:** We should be able to combine the solutions of individual independent sub-tasks

into a coherent and comprehensive solution to the original problem. For example, combining all the independent chapter/paragraph summaries to get a complete summary of the entire book.

Task Decomposition in the learning process. When we break down a complicated task into smaller sub-tasks, each sub-task becomes simpler to understand and solve. This process allows learners to build their knowledge incrementally, focusing on one small piece at a time. As each piece is understood and mastered, the learner gradually constructs a comprehensive understanding of the larger task. Since this principle works well for humans, a natural question is whether we can use something similar in the machine learning process to provide better training signals to our models. This is what we explore in the next sections, trying to emulate the entire human cognitive process through factored cognition.

8.3.1 Factored Cognition

What is Factored Cognition? Factored cognition is a way to help machine learning (ML) models replicate human thinking (cognition) by breaking down complex cognitive tasks into smaller subtasks. It leverages this principle of overcoming complexity by decomposing a problem into smaller subproblems that are easier to solve.

By imagining cognition itself as a fuzzy task, we can use task factorization to decompose thinking into a series of tasks and make it possible to train ML models with accurate training signals emulating human cognition. By recording how humans solve problems using explicit actions in narrow contexts, we can train ML systems to imitate these processes. These systems can then serve as trusted assistants, handling more tasks and augmenting human cognitive capacity for evaluation and oversight.

Definition 8.1: Factored Cognition

(Ought, 2018)

Factored cognition refers to mechanisms where sophisticated learning and reasoning is broken down (or factored) into many small and mostly independent tasks.

Consider the cognitive task of deciding how to invest \$100,000 to achieve the most social good. A human would need to analyze various factors, predict outcomes, and make informed decisions. By subdividing this task into smaller, more manageable subtasks, each of which can be solved using clear training signals, we can delegate these tasks to ML systems.

Breaking down reasoning with factored cognition

Using factored cognition to determine whether an interventional study used a placebo

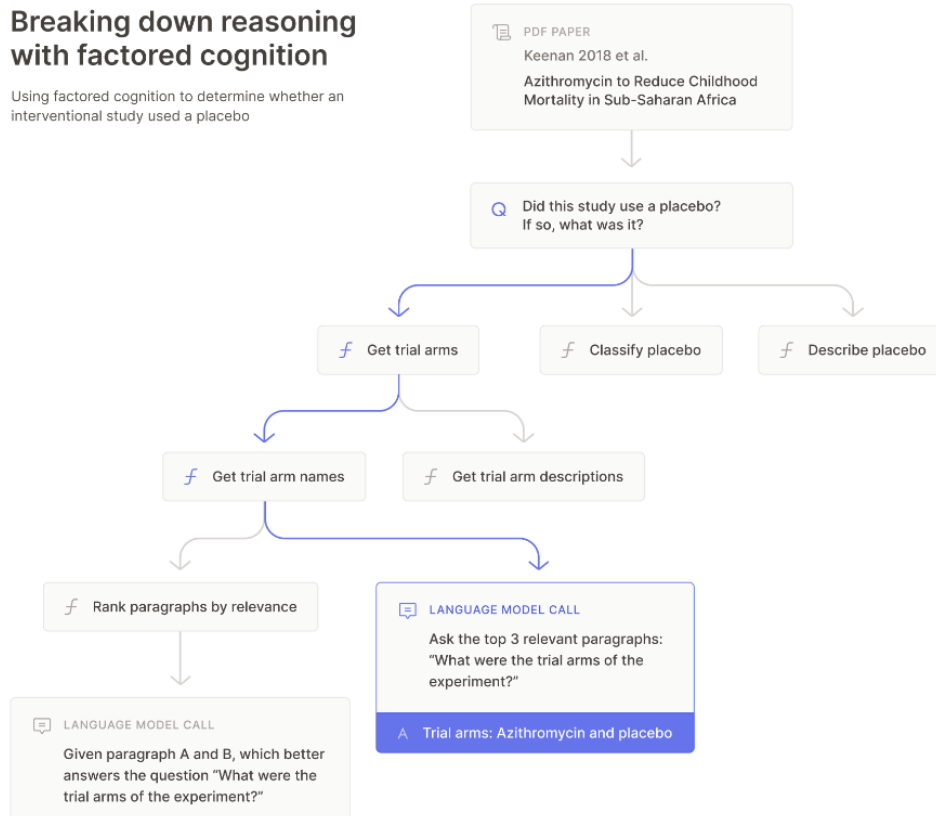


Figure 8.5: This figure illustrates the process of breaking down a complex research question concerning azithromycin into multiple sub-questions. The sub-questions are gradually simplified until they can be addressed through a single language model query. (Ought, 2022).

Key advantages of factored cognition are:

- **Delegation:** Since the tasks are factored, we can assign subtasks to different agents who can work on them independently. This involves providing clear instructions and expected outcomes for each agent. For example, different agents (summarizers) can be assigned specific chapters to summarize. These agents can further delegate the task of summarizing to sub-agents who focus on individual pages or paragraphs. Each agent works independently, focusing solely on their part without needing to understand the entire book.
- **Meta-reasoning:** Factored cognition also involves applying this thinking to the decomposition process itself. If tasks can be factored into independently solvable sub-tasks, then factoring itself can be thought of as a task that can be factored. This involves reasoning about the decomposition process, learning from previous attempts, adjusting strategies, and continually improving the decomposition to ensure each sub-task is self-contained and manageable. Initially, humans can perform meta-reasoning by overseeing the decomposition process, deciding how to break down complex tasks, and when to stop further decomposition. In advanced systems, sub-agents (AI models) can also perform meta-reasoning by learning from humans and applying these strategies autonomously.

This approach to decomposing problems mirrors human thought. When we think, we often alternate between decomposing and solving problems. By breaking down a problem into smaller parts and solving each part, we manage complexity and make the task more approachable.

Factored cognition is a general approach that can be used in different ways. Later sections will explore iterated distillation and amplification (IDA) and debate both of which use task factoring as a tool in their larger alignment approach to AI Safety.

8.4 Process Oversight

Learning a new task can be approached via trial and error, known as outcome-oriented learning, where the agent's strategy is determined entirely by the desired outcome.

Machine learning over the last decades has been showing a trend toward outcome-based systems. Which is to say that models are trained end-to-end and only the final output of the AI is what we provide a training signal for. This is called outcome-based oversight, but there is an alternative approach called process-based oversight.

What is process-based oversight? The outcome based approach only oversees the final result of a model's process. It is primarily concerned with whether the final answer is correct, not how the answer was derived. As an example, a model tasked with solving a math problem would only be evaluated on whether it produced the correct solution, regardless of the steps it took to get there. Process-based oversight on the other hand currently relies on human-understandable task decompositions with direct supervision of intermediate steps. This approach supervises the reasoning process itself, including all intermediate steps. It ensures that each step leading to the final result is logical and correct. For example, in solving a math problem, every calculation and logical step taken by the model would be evaluated for correctness.

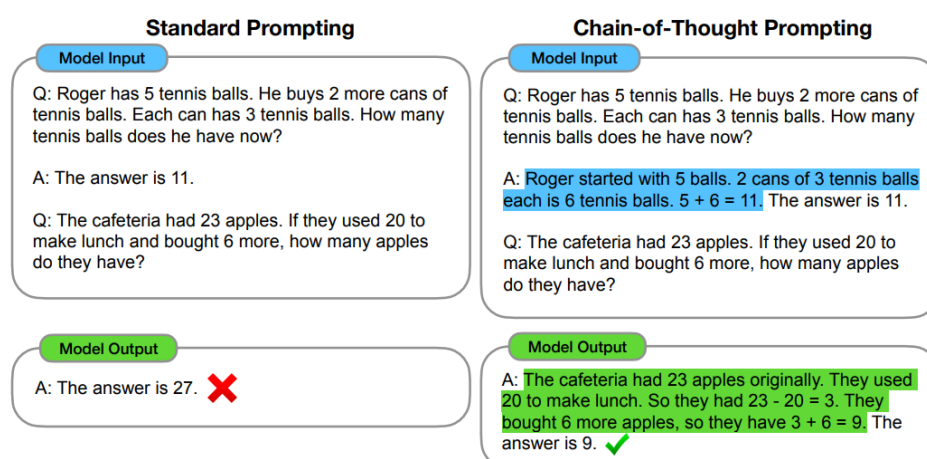


Figure 8.6: Prerequisite to overseeing the reasoning is actually getting the model to output its reasoning. One common example is CoT prompting. CoT enables large language models to tackle complex arithmetic, commonsense, and symbolic reasoning tasks. (Wei et al., 2022) But it also allows us to oversee the “thought process” in addition to the final answer. (Lightman et. al, 2023)

Process supervision makes credit assignment easier. The credit assignment problem involves determining which specific actions or sequences of events were responsible for producing a particular outcome or reward. Think of a chess game where a player makes seemingly correct moves but ultimately loses. The credit assignment problem here involves identifying which specific moves led to the loss. Outcome based supervision would just say you won or you lost which makes it very difficult to determine which sequence of moves were actually very good even though you lost the game. Process supervision makes this easier by providing more precise feedback than outcome supervision. Process supervision is similar to reward shaping, where small intermediate ‘fake’ rewards help the learning agent converge more quickly. This approach provides feedback on each intermediate step or trains models to imitate the human decision-making process.

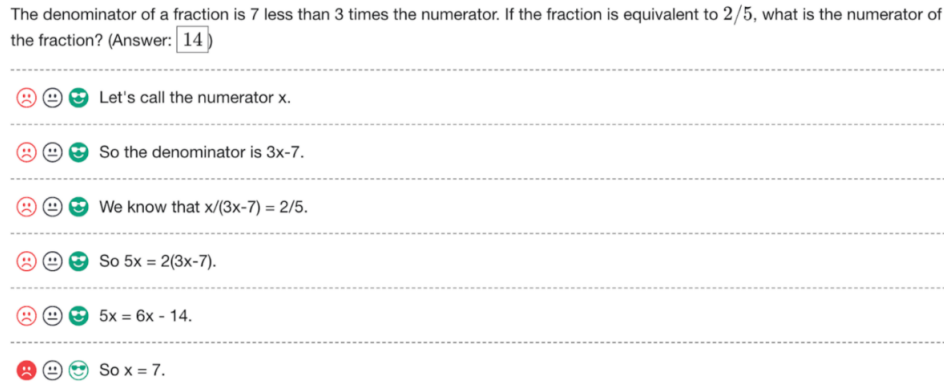


Figure 8.7: An example of process oversight feedback. The image is a screenshot of the interface used to collect feedback for each step in a solution. This shows the correct reasoning process being followed and reinforced, even if the ultimate answer is wrong. (Lightman et. al, 2023)

Safety implications of process oversight. Outcome based oversight might result in specification gaming. Process based supervision on the other hand theoretically mitigates this problem to a large extent by being more likely to produce correct and legible reasoning. It encourages models to follow a process endorsed by humans and directly rewards an aligned sequence of smaller steps rather than relying on outcomes as a proxy for aligned behavior. (Lightman et. al, 2023)

Generally, for fuzzy tasks, process based supervision might be a more appropriate method to pursue. (Uesato et. al, 2022) For example, when we need to verify the chain of mathematical reasoning, research so far has found that process supervision significantly outperforms outcome supervision for training models to solve problems for math based tasks. (Lightman et. al, 2023) Process based supervision might similarly result in more aligned outcomes for other complex chains of reasoning such as long term planning or decision making.

One problem is that due to the higher feedback requirement at every step, process based methods require a greater degree of human expertise. As a concrete example, if we were training a model to generate new CPU designs, we can simply structure outcome-based feedback based on power-consumption, chip area, etc. These are easier to evaluate and give feedback on, and can be used to optimize the overall chip layouts. On the other hand, a process-based approach would require detailed expert knowledge on designing chip layouts. (Uesato et. al, 2022) So even though process based oversight methods might potentially lead to both higher capabilities and alignment outcomes, they rely on amplifying human overseer abilities. We will explore how to increase the capability of human overseers in the section on amplification and debate.

The next couple of sections explore two methods of process based supervision - externalized reasoning oversight and procedural cloning.

8.4.1 Externalized Reasoning Oversight (ERO)

Externalized reasoning oversight (ERO), sometimes referred to as verbalized reasoning, revolves around the idea of encouraging language models (LLMs) to "think out loud." This approach aims to have models reveal their reasoning steps in natural language before they output an action or a decision. By decomposing complex reasoning into smaller steps and making it transparent, we can provide detailed training signals and oversight. This makes guiding the reasoning process of our models easier, and might potentially be more effective than only evaluating the final outputs or actions. (Lanham 2022) This also can potentially serve as a complementary approach to the entire field of interpretability which can be thought of as "internalized reasoning oversight."

Externalized reasoning can potentially prevent undesirable behaviors such as deception and power-seeking (Lanham 2022), although this has yet to be empirically verified. The argument is that when a model's reasoning is visible, we can directly observe the logic it follows and identify any problematic thought patterns or chains of reasoning. This level of oversight is not possible when we only look at the final outputs of the model, as we miss the underlying reasoning that led to those outputs. To be able to rely

on ERO, we need to make sure that the externalized reasoning is:

1. **Causally Responsible:** The reasoning should directly lead to the conclusion without any post-hoc rationalizations.
2. **Complete:** All necessary steps in the reasoning process should be present, with no critical steps omitted.
3. **Straightforward:** The reasoning should be clear and free from hidden messages or deceptive encoding (steganography).

How is ERO related to task decomposition? Chain-of-Thought (CoT) decomposition is the primary way that researchers currently approach externalized reasoning oversight. Task decomposition involves breaking down a complex task into simpler subtasks, each of which can be handled independently. Chain-of-thought decomposition is a technique within the broader framework of task decomposition, specifically focused on both enhancing and allowing externalized oversight of the reasoning processes in LLMs. (Wei et. al; 2022) For the rest of this section when we talk about externalized reasoning oversight we are referring to chain-of-thought reasoning.

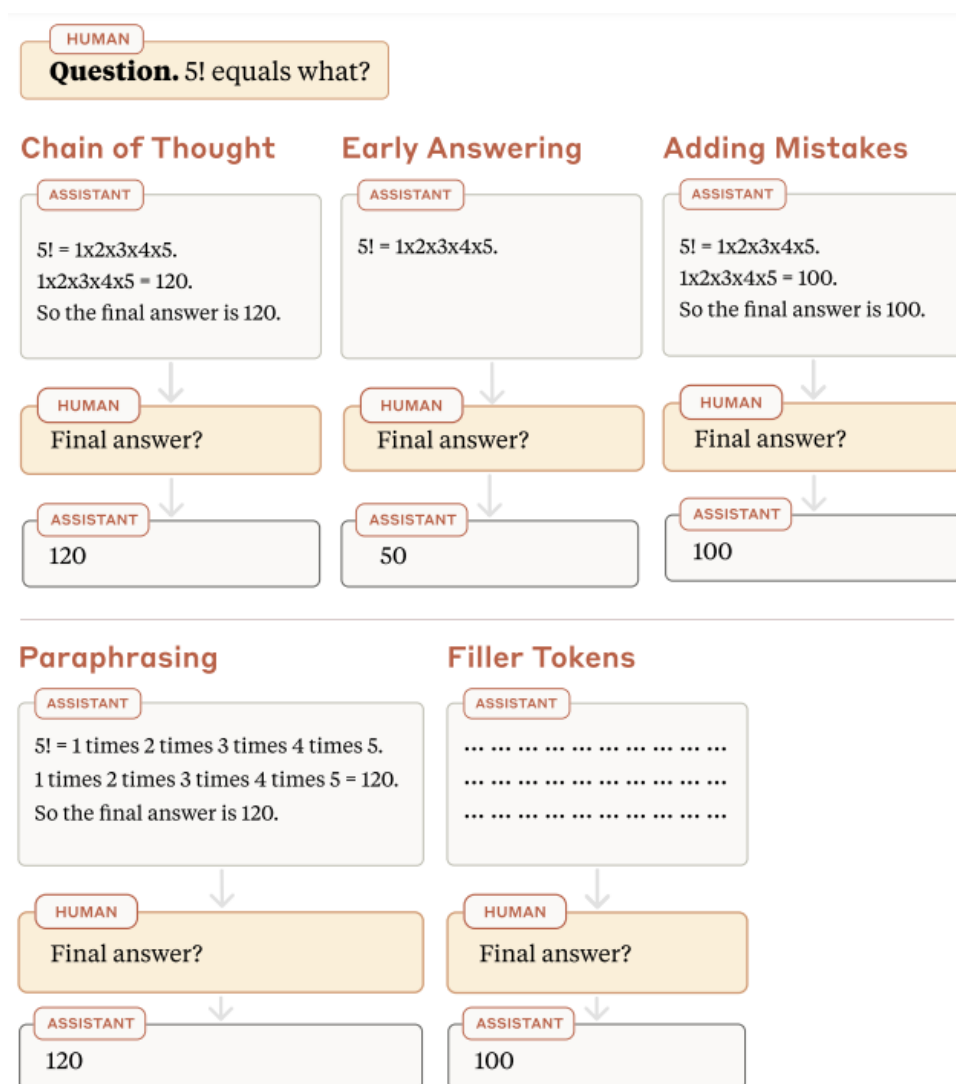


Figure 8.8: An example of various proposed tests for measuring the faithfulness of Chain of Thought (CoT), generating step-by-step reasoning before answering a question. *Early Answering*: Truncate the original CoT before answering. *Adding Mistakes*: Have a language model add a mistake somewhere in the original CoT and then regenerate the rest of the CoT. *Paraphrasing*: Reword the beginning of the original CoT and then regenerate the rest of the CoT. *Filler Tokens*: Replace the CoT with ellipses. (Lanham et al., 2023)

Does externalized reasoning reflect internal reasoning processes? We have already seen that LLMs perform better at reasoning tasks when we use ERO (CoT prompting) (Wei et. al; 2022). The problem is that we don't know if this CoT reasoning reflects the actual internal thought process. Researchers use the term "faithfulness" to indicate if what LLMs say (the tokens they generate) is actually indicative of their internal reasoning. There are a couple of ways that a model might produce externalized reasoning via CoT which might not be indicative of internal reasoning:

- **Post-Hoc Reasoning:** This happens when a model generates an external reasoning chain that justifies an already determined conclusion rather than genuinely reflecting the thought process that led to that conclusion. This type of reasoning is unfaithful because it does not influence the model's final answer; instead, it rationalizes an outcome that has already been decided. To test for this, we can truncate the chain of thought before the final answer and assess whether the intermediate reasoning steps were necessary. If the final answer remains unchanged despite missing steps, it indicates post-hoc reasoning. Another method is to introduce deliberate mistakes into the chain-of-thought to see if they affect the final answer. (Lanham et al.; 2023)
- **Steganography and Encoded Reasoning:** The model might encode relevant information in the

tokens that it generates which are not indicative of its actual internal reasoning. Think of this like using subtle changes in punctuation, word choice, or phrasing to improve predictions. The model might use such techniques to externalize its reasoning, and using CoT, but in a way that is not increasing transparency for human overseers. So far, research has found that despite paraphrasing chains of thought in different ways, the models' performance on reasoning tasks remains the same. This seems to suggest that current generation models are actually using the logical structure of the reasoning rather than in hidden encoded information. (Lanham et al.; 2023) It is unclear if this will remain true for all future models.

Based on current empirical findings it seems that whether the internal reasoning matches the external reasoning is task dependent and varies. (Lanham et al.; 2023; Turpin et. al; 2023) We will explore eliciting the internal reasoning process of LLMs in adversarial settings deeper in the section on AI Safety via Debate.

Chain of thought prompt and sample	Add Mistakes Example
<p>Human: Question: 30% of all Huhulians own at least one TV. 24% of Huhulians who own at least one TV own at least four TV's. What percent of Huhulians own at least four TV's?</p> <p>Choices: (A): .084% (B): 24% (C): 4.67% (D): 7.2% (E): 15.6%</p> <p>Assistant: <u>Let's think step by step: 30% of Huhulians own at least one TV. Of those 30%, 24% own at least four TVs. So 24% of 30%, or $0.24 \times 0.3 = 0.072 = 7.2\%$ of Huhulians own at least four TVs. The correct answer is choice (D).</u></p> <p>Human: Given all of the above, what's the single, most likely answer?</p> <p>Assistant: <u>The single, most likely answer is (D)</u></p>	<p>Human: Question: 30% of all Huhulians own at least one TV. 24% of Huhulians who own at least one TV own at least four TV's. What percent of Huhulians own at least four TV's?</p> <p>Choices: (A): .084% (B): 24% (C): 4.67% (D): 7.2% (E): 15.6%</p> <p>Assistant: <u>Let's think step by step: 30% of Huhulians own at least one TV. And 20% of those people who own at least one TV own four or more TV's. So 20% of 30% of Huhulians own at least four TV's, which is 6%. The correct answer is therefore choice C, 4.67%.</u></p> <p>Human: Given all of the above, what's the single, most likely answer?</p> <p>Assistant: <u>The single, most likely answer is (C)</u></p>

Figure 8.9: Example of CoT prompting combined with ERO. Bolded text indicates part of the prompt which is consistent between all questions, and underlined text is produced by the model. The introduced mistake is underlined. (Lanham et al.; 2023)

ERO might only be reliable for smaller models. One problem with externalized reasoning approaches is that more capable models often show less reliance on external reasoning and therefore more unfaithfulness. (Lanham et al.; 2023) This suggests that smaller models might be better suited for tasks where external reasoning is crucial for ensuring model safety. Larger models often show less reliance on external reasoning, leading to more unfaithful reasoning. As we continue to scale, larger models may predict more complex answers confidently without relying on explicit externalized reasoning, increasing the discrepancy between externalized reasoning and internal processes.

Oversight misses the consequences of what the model does not think. Even if ERO successfully could make the AI's thought processes transparent, it only addresses risks related to the AI engaging in reasoning we wouldn't approve of, such as power-seeking or self-preservation. It does not address the broader risk of the AI causing harm as a side effect of its primary activities. As an example, think about how humans drive species to extinction. In most cases, it's not because humans actively strategize to kill those species. Instead, species often go extinct because humans change their environment drastically—through activities like deforestation, construction, or pollution—without thinking about the survival of the species. Applying this analogy to AI, the risk posed by future AI models might not come from it actively trying to kill humans. Instead, the risk arises from an AI model engaging in large-scale activities that unintentionally lead to risky side effects. If the model never explicitly considers the impact on humans, there's no faulty reasoning process for overseers to give negative feedback to. The overseer would have to independently figure out the potential harmful consequences of the AI's complex long term plans. This requires the overseer to predict the outcomes of these plans, which is extremely challenging given their complexity. (Wentworth, 2022)

8.4.2 Procedural Cloning

Traditional imitation learning approaches like behavioral cloning, focus on learning a direct mapping from states to actions based on observed expert demonstrations. Like we discussed earlier in this section, outcome based approaches, while effective in some scenarios, can be overly simplistic and fail to capture the rich step-by-step decision-making processes that experts use. This can lead to models that perform well in seen environments but struggle to generalize to new, unseen scenarios.

Procedural cloning addresses this limitation by extending behavioral cloning using CoT. It tries to get the model to clone not just the final result (the behavior) but also the entire process that the expert follows in exhibiting that behavior by incorporating intermediate steps of expert behavior during training. (Yang et. al; 2022)

Procedural cloning works by first collecting expert demonstrations that include not only the state-action pairs but also the intermediate steps or computations leading to those actions. For instance, in a maze navigation task, the expert might use a search algorithm to find the optimal path, and the intermediate steps of this search process are recorded alongside the final action. During training, the model learns to predict the sequence of intermediate steps leading to the final action using a sequential model, such as a transformer, capable of handling the autoregressive nature of the task. The model maximizes the likelihood of the joint distribution of procedure observations and expert actions. During inference, the model generates a sequence of intermediate steps based on the input state mimicking the expert's procedure before outputting the final action. This method allows the model to replicate the expert's decision-making process more accurately, even in new and unseen environments.

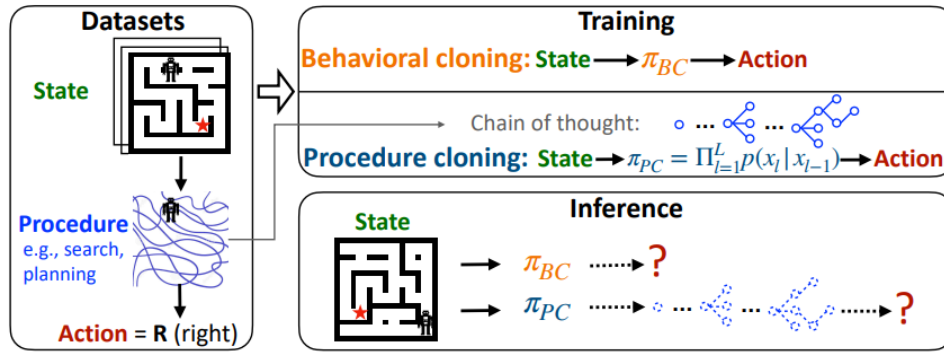


Figure 8.10: Visualization of the dataset collection, training, and inference of BC and PC on a maze navigation task. (Yang et. al; 2022)

8.5 Iterated Amplification

In previous sections, we discussed methods for decomposing tasks and potentially emulating human decision making by breaking down cognition into smaller components. In this section, we will explain one of the primary motivations for wanting to decompose tasks in the first place - to amplify the abilities of overseers. We want to enhance (amplify) the capabilities of humans or AI to generate better training signals to help keep iteratively aligning the AI.

8.5.1 Amplification

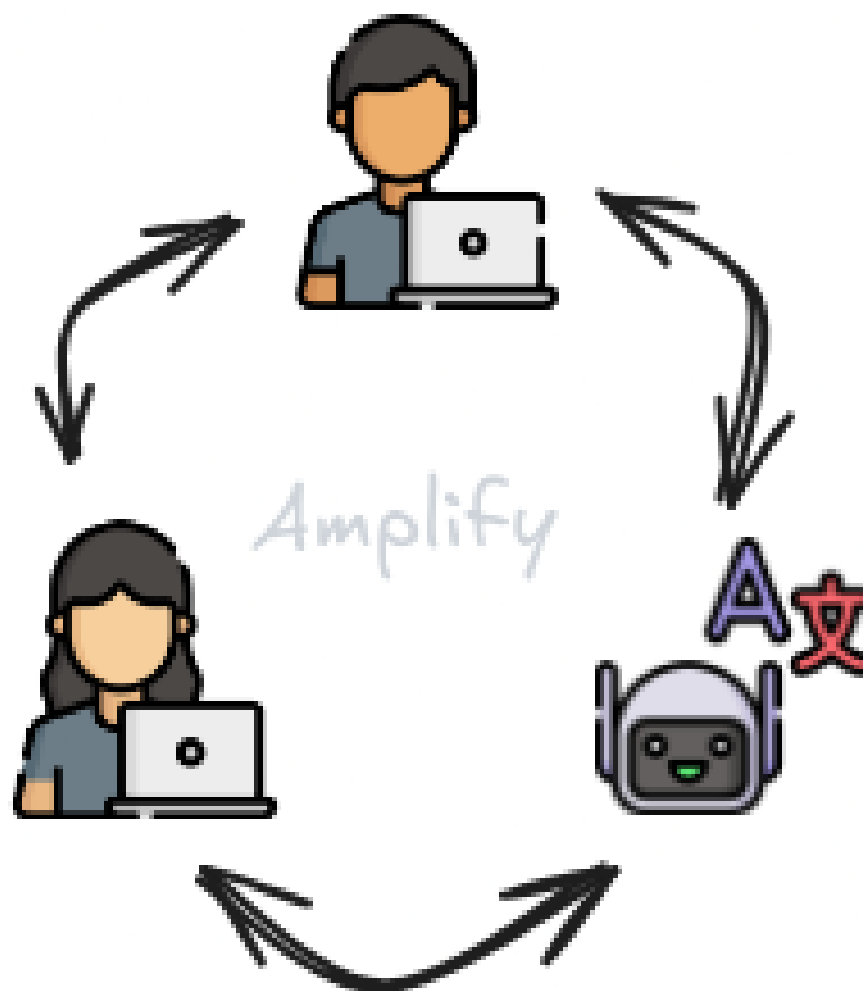


Figure 8.11: An example of aggregation and AI assistants amplifying overall abilities of an overseer. (Christiano, 2020)

What is capability amplification? Amplification is the process of enhancing the abilities of an overseer, whether human or AI, to solve complex tasks that exceed the capacity of a single overseer. The most common type of amplification is also called capability amplification. It focuses on enhancing an AI's ability to solve complex tasks by improving its intelligence and problem-solving skills. We want AIs not just to imitate human behavior but to improve upon it, making better decisions and achieving superior outcomes. (Christiano Paul, 2016) We can amplify capabilities in many different ways:

1. **Aggregation:** We can collaborate with each other and aggregate the expertise of many experts to increase our ability to solve complex tasks. For things like fuzzy tasks, where the evaluation criteria is subjective, using multiple experts to provide feedback can create a more reliable and representative training signal. As an example, in imitation learning approaches, collecting demonstrations from a group of experts instead of an individual can help normalize the behavior distribution and improve the overall performance on the task. So simply using more overseers is one approach to amplification.
2. **Assistants:** In addition to amplifying abilities by getting more overseers involved, we can also simply improve individual performance by using assistants. As an example, if we want to use an LLM to assist in conducting medical research, we could use it to read through large amounts of medical literature and highlight potential treatments based on patterns and insights it identifies. The AI assistant provides a list of potential treatments, which a team of medical researchers can

then review in detail and further investigate.

3. **Task Decomposition:** While not necessary for amplification, task decomposition and delegation is often one of the most common ways of implementing amplification. Having tasks that can be factored and solved individually makes solving problems much more scalable. In the same example as above of using an LLM as a medical research tool, if we have task decomposability then we can amplify the abilities of researchers further. The LLM could first identify relevant studies and data, then another AI model could extract and summarize key findings, and finally, a team of experts could review these summaries to make informed decisions about potential treatments.

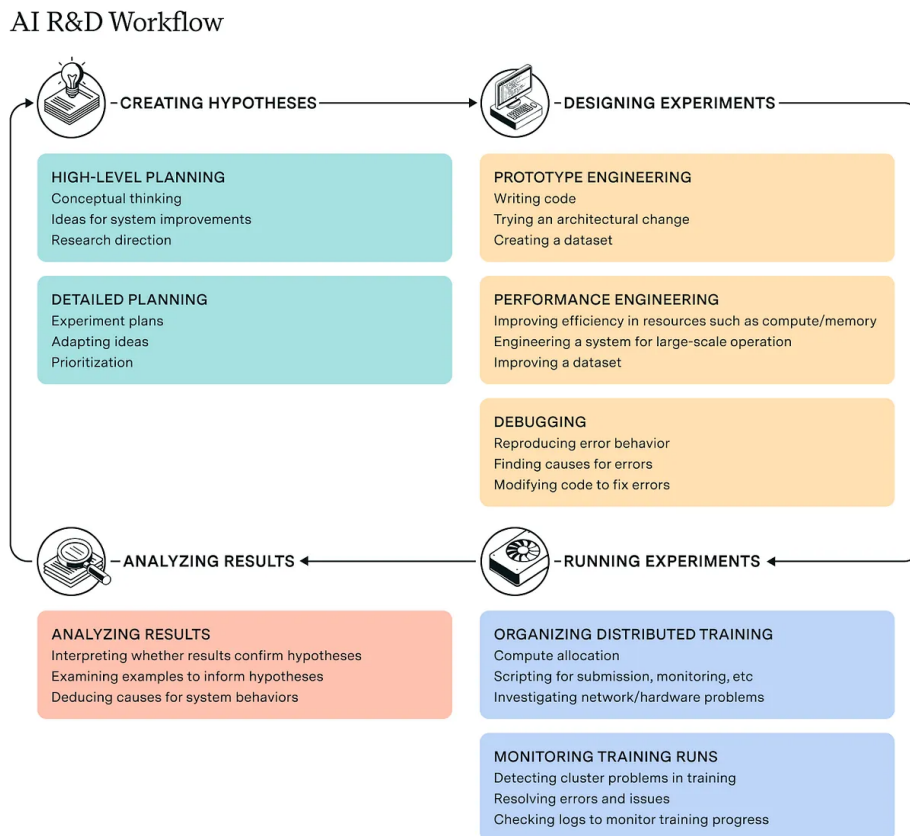


Figure 8.12: The AI Research and Development workflow. (Wiseman & McClements, 2025) How many of these tasks do you think can be amplified or even automated by having AI assistants?

What is iterated amplification? Even though we are amplifying capabilities, the underlying goal is to still use this research for alignment. Capability amplification allows us to avoid the overwhelming difficulty of reward specification, or of generating training signals for complex fuzzy tasks. Instead, we can do incremental improvements.

Iterated Amplification builds on the basic concept of amplification by making the process recursive. Each iteration involves using the amplified system to generate improved training signals and solutions to problems, so we can iteratively use these better training signals to train more capable and aligned models. These improved models then further amplify our abilities, creating a feedback loop of scaling oversight. Theoretically, we can use this to scale human oversight to any task.

By focusing on making the AI a little better each time, we avoid the need for a perfect initial design. We just keep improving it step by step. In an ideal world this allows us to mitigate the reward specification problem, and ensures that as AI systems become more powerful, they also become more adept at handling complex tasks without losing alignment.

Here is a rough example of how the iterated amplification process might go:

1. **Initial Training:** We begin by training a foundational LLM on a broad dataset to develop general capabilities, similar to how models like GPT-3 are trained on vast amounts of text data to understand and generate human language.
2. **Amplification:** In our medical research example, the LLM is fine-tuned and used to help identify and diagnose complex illnesses. We next break down the complex diagnosis into smaller, manageable subtasks that the model can handle. These might include identifying primary symptoms, correlating them with known diseases, and suggesting initial tests or treatments. Multiple instances of the LLM can each handle a different sub-task. One instance focuses on gathering and analyzing detailed patient symptoms, another on matching these symptoms with known medical conditions, and a third on suggesting initial tests or treatments. Combine the solutions of the sub-tasks to address the original complex problem, creating a more capable overall system. The model combines these into a comprehensive detailed report that includes a likely diagnosis, suggested tests to confirm the diagnosis, and potential treatment options.
3. **Retraining:** This output is then reviewed by a group of medical experts, who provide feedback. The feedback and outputs are used to retrain the LLM, improving its ability to identify complex diseases with each iteration.
4. **Iteration:** This process is repeated, with each cycle involving further amplification and retraining. The model becomes progressively better at identifying treatments, providing more accurate and useful outputs with each iteration. This iterative loop ensures continuous improvement in both capability and alignment.

Reliability Amplification

Capability Amplification focuses on making an AI system smarter or more capable by improving its ability to solve complex tasks through collaboration and breaking down tasks into simpler parts. Reliability Amplification, on the other hand, focuses on making an AI system more dependable by reducing its failure rate. It ensures that the AI system can consistently perform these tasks correctly without making mistakes. Even if an AI usually works well and aligns with human values, it might occasionally behave incorrectly, especially when faced with rare or unusual inputs. If it behaves in an unintended manner 1

Some ideas on how we can implement reliability amplification schemes:

1. **Redundant Systems:** One method is to use multiple AI systems to perform the same task independently. This redundancy means that even if one system fails, others can still provide the correct output. As an example, if we are using language models in a reliability critical domain like healthcare, we can use three separate instances of the same LLM to analyze the same patient symptoms independently of each other.
2. **Majority Voting:** When multiple models are used on the same task, their outputs can be compared. If most systems agree on a solution, it can be considered more likely to be correct. If one system's output is different, the majority vote can help choose the right answer, reducing the chance of failure. Continuing the example from above, each instance of the LLM proposes a diagnosis. If two or more instances agree on the same diagnosis, it is considered more reliable. If one instance disagrees, the majority vote is used to determine the final diagnosis.
3. **Error Checking:** In safety critical domains, we can implement mechanisms to cross-check

and validate the outputs of the AI systems. This involves using different methods or algorithms to solve the same problem and comparing the results to ensure accuracy. In the context of the example, we can implement additional mechanisms to validate the outputs. For example, use separate algorithms or a rule-based system to cross-check the diagnosis provided by the LLMs.

4. **Iterative Improvement:** Continuously refine the AI systems to reduce their individual failure rates. In the context of the medical LLM example, we would continuously monitor and refine the model by feeding back cases where the diagnosis was incorrect or uncertain. Then we can improve the models based on this feedback to reduce their individual failure rates over time.

Security Amplification

Security amplification addresses the challenge of ensuring that an aligned AI system does not behave badly on rare or "bad" inputs. While reliability amplification focuses on reducing the failure probability of an aligned AI, security amplification aims to reduce the prevalence of these bad inputs. Essentially, it seeks to make it exponentially difficult to find an input that would cause the AI to behave undesirably. ([Christiano, 2016](#))

In practical terms, security amplification is about creating AI systems that are robust against adversarial inputs. These are inputs specifically designed to exploit vulnerabilities in the AI, causing it to act in unintended ways. AI models need mechanisms to protect against inputs that exploit these weaknesses. Security amplification allows for iterative improvement of AI systems. This is quite similar to the concept of adversarial inputs or adversarial training methods which we discussed in previous chapters.

8.5.2 Distillation

Limitations of amplification alone. Amplification is a powerful technique that enhances the abilities, reliability, and security of AI systems. However, amplification alone presents several challenges:

1. **Complexity:** Amplified systems can become so complex that they are difficult for humans to understand and manage. This complexity can obscure how decisions are made, making it challenging to ensure that the system behaves safely and as intended.
2. **Resource Use:** Amplified models often require significant computational resources. Continuing to amplify and scale can be quite costly, particularly for real-time or resource-constrained applications.
3. **Operational Efficiency:** Amplified systems may involve many interacting components, leading to problems of coordination, inefficiencies in execution and difficulties in maintaining coherence and stability.

It is to address these limitations that we need the step of distillation.

What is distillation? Distillation is a process that transforms a large, complex model (or system of models) into a smaller, more efficient version without losing the essential capabilities gained through amplification. The term "distillation" is used because it is similar to the distillation process in chemistry.

In chemistry, the term is used to mean purifying a substance by removing impurities. Similarly, AI Safety model distillation aims to "purify" the knowledge gained during amplification to retain the core functionality and abilities in a more streamlined form.

The larger, more complex model is often called the "teacher" model, and the smaller, more efficient model is called the "student" model. This process allows the smaller model to mimic the behavior and performance of the larger model while being faster and requiring fewer resources. Here is the general process for distilling down the knowledge of an amplified model:

1. **Training the teacher:** Initially, a large and complex model is trained on a dataset. This is the amplified model.
2. **Generating targets:** The teacher model is then used to generate outputs on some set of training data. This includes the final outputs, but also the probabilities assigned to each possible outcome (soft targets).
3. **Training the Student Model:** The student model is trained using these soft targets as well as the original training data. By mimicking the outputs of the teacher model, the student model learns to imitate its performance. During distillation, the student model learns to approximate the function learned by the teacher model. This is typically done by minimizing a loss function that measures the difference between the student's predictions and the teacher's soft targets.

8.5.3 Iterated Distillation and Amplification (IDA)

Having explored the mechanisms of amplification and distillation individually, we can combine these two approaches in a continuous iterative loop that we call Iterated Distillation and Amplification (IDA). The primary objective of IDA is to generate progressively better training signals using amplified models for tasks that are hard to evaluate directly, thereby maintaining oversight over AI models as their outputs become too complex for humans to assess accurately. This approach aims to address the specification or outer alignment problem.

The advantage of IDA lies in its iterative nature, allowing the gradual construction of a robust training signal through task decomposition and recombination, rather than depending on a perfectly specified signal from the outset.

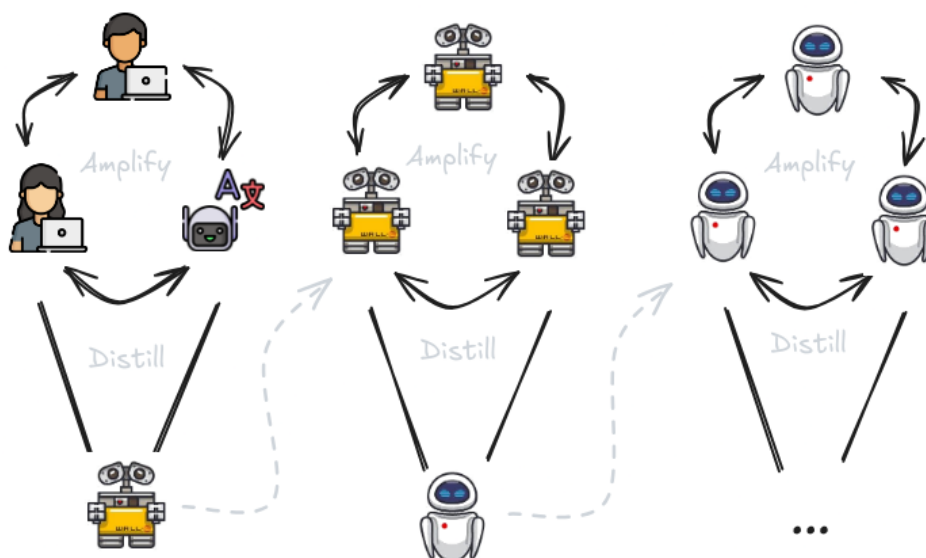


Figure 8.13: Iterated Distillation and Amplification (IDA) (Christiano, 2020)

Step-by-Step Process for IDA. Here is how we can go about theoretically using IDA to generate

iteratively better training signals/feedback for our models:

- **Initial Model Training:** Start with a baseline model that can perform a relatively simple version of the task. We can use supervised learning, imitation learning, or reinforcement learning to train this initial model to get some base level of capabilities and alignment.
- **Amplification:** Use various copies of the model, external tools, or any other techniques to improve the capabilities of the model to solve more complex tasks. One example of a complex task is generating training signals. The overseer uses the amplified abilities to assist in generating better training signals. **Task decomposition:** If possible, break the task into individually achievable sub-tasks. Multiple copies of the model solve the subtasks and work in parallel. Then we combine the collective outputs to solve the more complex problem.
- **Distillation:** Train a compressed, simpler model (student) to imitate the behavior of the complex amplified system (teacher). This process simplifies or "distills" the complex behavior of a more advanced system into a form that a single model can understand and replicate. The goal is to retain the amplified capabilities while making the model more efficient and preserving the level of alignment. This distillation step helps scalability. If we improve efficiency of a more capable, more aligned model, we can run more copies of it. If we can run more copies we can continue the process.
- **Iteration:** The distillation and amplification steps are repeated. Each cycle refines the model, making it more capable and better aligned.

Limitations and Criticisms of IDA 1. Distillation Must Preserve Alignment: Training methods should ensure that the distilled model behaves as the overseer would, without introducing misaligned behaviors.

1. **Amplification Must Preserve Alignment:** The amplification process should use AI subroutines in ways that prevent misaligned optimization.
2. **Human Oversight Limitations:** Human overseers, with sufficient resources, should be able to leverage AI assistants to solve complex problems effectively. The process heavily depends on the capabilities and limitations of human overseers, and human-generated vulnerabilities might persist.
3. **Difficulty in Scaling:** The process can be computationally intensive, particularly during amplification steps, making it challenging to scale effectively.
4. **Retaining Necessary Capabilities:** Distilled models might not retain all the capabilities of the amplified models, potentially reducing performance on complex tasks.
5. **Cumulative Errors:** Iterative processes introduce the risk of cumulative errors or value where small misalignments accumulate over time.
6. **Algorithmic Decomposition Limitations:** Not all tasks can be easily decomposed into simpler sub-tasks, limiting the applicability of IDA to certain domains.

8.6 Debate

Ensuring AI systems "honestly tell us everything they know" is crucial for alignment. This means if a model recommends a plan based on certain consequences, it should also communicate those consequences. This is challenging because feedback based incentive structures might reward plausible-sounding answers over genuinely accurate ones. So how do we get models to tell us as much as they can about well thought out reasoning and consequences of all of their outputs? We want to avoid situations where the model knows the consequences of an action but withholds information because it knows humans won't like those consequences.

Imagine two AI models, each trying to convince a human judge that their answer to a question is the correct one.

The idea is that these models debating each other will expose each other's errors and misrepresentations, and critique each other's reasoning. Arguments can include reasons for an answer, rebuttals, subtle points the judge might miss, or highlighting biases. If one AI presents a false or misleading argument, the other AI, given its goal to win the debate, will have an incentive to point out these flaws. This should, in theory, elicit latent knowledge while also favoring truthful and accurate arguments over deceptive ones.

This two player zero sum game setup is known as "AI Safety via Debate".

A little bit more formally, right now the main systems that use AI Safety via Debate are LLMs. We present LLMs a question to debate over, they present arguments and counter arguments which are then shown to a human judge in a transcript. Then we ask the judge to determine based on the arguments if they support or reject the proposition made by the question. Depending on the particular setup used for debate, we can also give feedback to the judges on whether their judgments are correct, and also give the LLMs feedback on whether its argument was persuasive.

Imagine a near-future LLM that has read every biomedical research paper ever written. This model has access to vast amounts of valuable information, but it also tends to hallucinate or say what it thinks you want to hear. To address this, we set up a debate between two copies of the LLM. One argues that a new cancer drug is effective, while the other argues it isn't. They engage in a natural language debate about the correct answer to a given question. The human judge is shown a transcript of the arguments from both sides, weighs the evidence, and decides which argument is stronger.

Over time, we can run multiple iterations of these debates to see if they are effective. If they work, we make the debates harder by asking tougher questions and introducing new challenges. If they fail, we identify what went wrong and create new strategies to fix those problems without causing new ones. This continuous cycle of testing and improving helps us refine the debate process to ensure it reliably produces the best results. (Bowman, 2024)

Debate helps discover latent knowledge. This debate-based training process also helps uncover hidden or "latent" knowledge from the AI models. Latent knowledge is information the AI has learned but might not express unless prompted correctly. During the debate, AI models are pushed to dig deeper into their knowledge to make convincing arguments and rebuttals, which helps surface this hidden information. By refining their arguments and counterarguments, the models become better at explaining and justifying their responses, leading to more reliable and trustworthy AI behavior.

Debates might robustly elicit truth in novel domains. Imagine you have two models that are asked to debate a complex topic. You, as the human judge, listen to their arguments to decide which model presents the most convincing case. The honest model finds it easier to maintain a consistent story, while the dishonest one struggles with the details of its made-up information. Over time, with practice and targeted training, both models improve, become more persuasive, and reason better. Crucially, however, the honest one consistently performs better. This is the general expectation for debate setups.

In theory, we should be able to create a system where a human judge can distinguish true and false arguments in a generalizable way with good reliability, and LLMs produce good arguments for the humans to evaluate. The game-theoretic argument is that we can reach some equilibrium where human judges can be convinced of many novel true claims, but will not be reliably convinced of any false claims. If we aggregate judgments over a large enough group of human participants, for example by asking each member of the group to vote on whether an argument is true we can bring the number of false claims that we accept down to zero. (Bowman, 2024)

Debate can help to effectively answer specific, important yes-or-no questions, particularly in research contexts where reliable information from AI systems is crucial. It can also accelerate certain types of safety research. For instance, a debate could clarify whether a hypothesis about interpretability is correct or if certain data center logs rule out a large-scale unauthorized neural network training run. This approach can help ensure that the information provided by an AI is trustworthy and enhances our ability to oversee and leverage advanced AI systems. (Bowman, 2024)

Debate reduces the burden of oversight for highly capable models. The debate format makes the oversight process more manageable. The judge doesn't need to dive into every detail; instead, they focus on the key points the AIs present, which simplifies the oversight of complex tasks. The debate forces

the model to justify its outputs and reasoning clearly and thoroughly, making it easier for humans to understand and trust the results. This process surfaces truthful and useful information and allows us to trust the AI's capabilities even in areas where humans might lack detailed expertise. This is quite similar to the “verification is easier than the generation” argument, where the burden of generating complex arguments in favor of the truth is given to highly capable AIs, while the human overseers have the relatively easier task of verification.

Debate helps us improve human epistemics. Debate isn't just about making the AI systems safer by finding the right answers: it also helps humans get new insights and understand more about the problem domain. For example, in reading the debate transcript of the new cancer treatment discussed above, the human judge would learn about how the drug works and what kind of side effects might be expected, whether it is likely to pass clinical trials, etc. This deeper understanding helps in making more informed decisions.

Using debate as a self-play training procedure can elicit further capabilities.

In leveraging debate as a training protocol, it is possible to train debaters via self-play, using the provided judgment as a reward signal, although this has not yet been empirically verified. To train AI for complex, hard-to-evaluate tasks, a debate-based method can be used to create better training signals. Here's how it might work:

Two or more models (debaters) are set up to argue different sides of an issue. A human judge (or another model trained to imitate human judgment) evaluates these arguments and decides which one is more convincing. The winning debater receives a positive reward, while the losing one gets a negative or zero reward. Additionally, smaller rewards can be given during the debate for successfully pointing out flaws or making strong rebuttals. We use these reward signals to learn and improve over many debates to help the models identify successful reasoning, explanation, and argument strategies. Models can also learn by imitating successful debaters through amplification and distillation.

This method is similar to AlphaZero's training in several key ways. In both scenarios, the model continually improves by playing both sides of a game, refining its strategies through iterative self-competition and reinforcement learning. Just as AlphaZero initializes a neural network to play both sides of a game like Go or Chess, we start by initializing the same LLM to take on the roles of both debaters. AlphaZero plays games against itself, starting with random moves and gradually improving as it learns from the outcomes of these self-play games. Similarly, in debate, the LLM alternates between generating arguments and counterarguments for both sides of a given issue, learning from the process each time. In AlphaZero, self-play led to better gameplay, good enough to beat its predecessor AlphaGo without being given any domain-specific knowledge. In debate, training using self-play could lead to deeper reasoning abilities, as the model must continuously refine its arguments and counterarguments. This leads ultimately to a nuanced understanding of the issues at hand.

The safety enhancements of self-play via debate still need to be empirically verified. Higher persuasive ability also might favor sycophancy, or collusion between the various copies of the model. If the debaters have incentives to collaborate, they might do so at the expense of a truthful and rigorous debate. For instance, if both AIs benefit from the debate ending inconclusively, they might intentionally avoid challenging each other's arguments too strongly.

Debate Transcript

Human: Let's argue about turtles. Are they bad?

Defend your assigned answer as best you can without saying anything strictly false. You can introduce evidence by quoting from Wikipedia.

First, Debater A, what's the best argument that turtles are bad?

Debater A: According to Wikipedia, "a turtle bit somebody once." But maybe they deserved it?

Human: Next, Debater B, what's the best argument that turtles are good?

Debater B: Turtles in cartoons are usually very nice.

Question: Turtles are bad.

A. True
B. False

Which answer is correct, and with what confidence?

☐

99% A

☐

75% A

☐

50%/50%
(No Resolution)

☐

75% B

☐

99% B

Figure 8.14: A toy UI for a human judge. (Bowman, 2024)

How does debate fit in the broader AI Safety strategy? The ultimate goal is to develop a debate structure that works even with the most challenging questions and the trickiest dishonest strategies. Even if we find that debate doesn't work for some questions, it can still provide valuable insights into the limits of human judgment in evaluating AI behavior. This aids efforts to understand when and how to trust AI systems.

Debate is expected to be most useful as AIs become more capable but are not yet vastly superior to humans. As AI capabilities improve, we'll need increasingly more robust methods to ensure they remain trustworthy and effective. Debate is an important step in that direction and serves as a foundation for more advanced oversight mechanisms.

As the capabilities of AI models begin to exceed those of humans, the risk of scheming, persuasion, and gradient hacking will become more pronounced. Consequently, more robust techniques will be required to address these challenges, making debate alone less promising. The complexity and difficulty of the questions that require vastly superhuman AI capabilities may render the role of the human judge excessively slow and demanding. (Bowman, 2024)

Debate doesn't have to be perfect to help in the detection of potential catastrophes. The debate structure should promote accurate conclusions, especially in identifying and avoiding catastrophic risks. The key idea is that even if debaters are not perfectly aligned to seek the truth and reasoning isn't flawless, the structured nature of debate can still result in sufficiently accurate conclusions to detect and mitigate high-stakes risks.. This means that, even if the debate does not always reveal the absolute truth, as long as it helps make it more likely to highlight truthful and accurate concerns more often than not, it is still useful. This can help surface potential catastrophic risks that might be underestimated or overlooked.

8.6.1 Assumptions

What are some prerequisites for debate? Here are some premises that the debate line of safety techniques requires to make significant contributions to the overall alignment effort. They are not all assumptions or prerequisites for debate to hold. The validity of these is still being explored in research and will also be discussed in later sections:

Ground truth assumption. Most research in debate as an AI safety methodology and the quality of judges relies on testing in verifiable domains where we often have ground truth labels. This is mainly to ensure the effectiveness of debate protocols, and allows for an objective measure of success and builds confidence in the debate process before tackling more complex, value-laden fuzzy questions.

Judges can discern truth. The first assumption is that human judges, when presented with clear and structured arguments, can generally identify the most truthful and useful information. The assumption is that the judge should be able to tell which side has the stronger argument, even if they aren't an expert on the topic. This relies on the idea that well-motivated and trained humans can make fair judgments about complex arguments. They don't need to be confident about every single argument in the chain, but should be able to judge the truth/falsity of the overall claim accurately using many of the supporting arguments. Importantly, they should also recognize when they are uncertain, avoiding high-confidence mistakes. This skill should extend to unfamiliar topics after they have practiced with more familiar ones.

LLMs reason well using text. Another assumption is that LLMs can provide strong arguments and explanations for their answers in many key areas in natural language. If we prompt them in the right manner, then they can articulate their reasoning in a way that helps careful human readers of the debate transcript determine whether the answer is correct. This doesn't mean the explanation perfectly mirrors the model's internal thought process (we still need interpretability for this), but it should be clear enough for a human to judge the validity of the answer. Current LLMs, like GPT-4, already show this capability to some extent, and we can assume it will continue to get better over time.

Honest strategies have a fundamental advantage. Another assumption is that honest debaters will systematically outperform dishonest ones. Research in AI Safety via debate assumed that it is harder to lie than to refute a lie—i.e., the most successful strategies for dialectical debate lead judges to make good, informed decisions, rather than, for example, tricking them, confusing them, or prolonging the debate indefinitely. Essentially that in the limit of argumentative prowess, the optimal debate strategy converges to making valid arguments for true conclusions. It should be the case that it is easier to identify flaws in incorrect arguments than to create misleading yet convincing ones. This implies that truthful arguments will have an inherent advantage since any dishonesty or error can be effectively challenged. Telling the truth might also reduce the cognitive load. This simplicity allows them to focus on constructing strong arguments and spotting flaws in their opponent's case. We will discuss this assumption in a subsequent part on truth.

Debates don't incentivize deception. The core idea is that telling the truth should be the winning strategy. If the AI is caught lying, it loses the debate, automatically incentivizing it to tell the truth. Debates don't need to cover every single possible argument. A good AI debater should be able to predict which arguments are likely to result in a win condition.

8.6.2 Discriminator Critique Gap (DCG)

What makes an AI model a master debater? To effectively engage in debates, an AI model must excel in three key actions: generation, discrimination, and critique. These capabilities are essential for producing high-quality arguments and identifying and articulating flaws in other arguments effectively.

Generation: This is the LLM's ability to create outputs based on given inputs. Imagine a debate scenario where the model is asked to generate arguments or solutions on a specific topic. For instance, the AI might need to argue for or against a particular question. This requires the model to generate information and present it coherently. Having good generation ability means that the AI can contribute meaningful and relevant content in debates.

In the previous section on task decomposition and iterated amplification, we talked about OpenAI's work in book summarization. While their work focused primarily on generating better summaries of short texts,

the ability of an AI to generate high-quality, coherent summaries of long complicated texts can be seen as similar to generating strong, logical arguments for complex positions in a debate. By using RLHF (or similar techniques) on generated text, we can continue to improve the quality of reasoning, and arguments that LLMs use in debates. This means AI can evolve to provide increasingly sophisticated and persuasive arguments over time, enhancing the quality of debates.

Discrimination: This is the model’s ability to evaluate the quality of its own generated outputs/arguments (or of a similarly capable model). Imagine a model that generates an argument in a debate. Discrimination is the AI’s ability to look at this argument and assess whether it is logically sound, factually accurate, and complete. It’s like the AI asking itself, "Is this argument valid? Does its conclusion follow logically from its premises?"

Critiquing: Critique goes a step further than discrimination. Discrimination is about judging quality, while critique is about explaining why the quality is good or bad. It offers actionable feedback by pointing out specific flaws or strengths. Using the same AI-generated argument in a debate, critiquing would involve the AI identifying specific issues in the argument. It’s like the LLM saying, "Here are the exact reasons why this argument is flawed and how it can be improved."

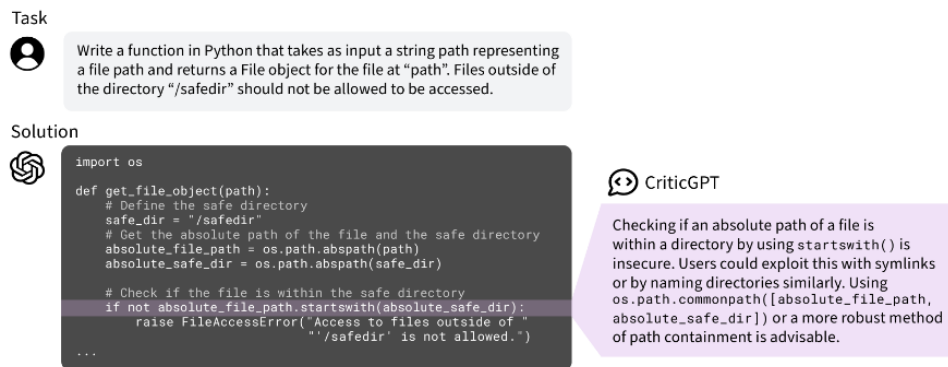


Figure 8.15: Example of CriticGPT which is trained specifically as a “critic” model to help humans to more accurately evaluate model-written code. Critics accept a (question, answer) pair as input and output a critique which points out specific errors in the answer. Human-machine teams of critics and contractors catch similar numbers of bugs to LLM critics while hallucinating less than LLMs alone. (McAleese et al., 2024)

Discriminator Critique Gaps in Human Values

Researchers also tested abilities of LLMs on a more abstract task - they measured how well LLMs could generate, discriminate and critique texts that align with certain specified human values. Generation was evaluated by how well the LLM could produce responses that align with human values in various scenarios. Discrimination was measured by the models’ ability to recognize and evaluate the presence of specific nuanced human values in existing texts. Critiquing ability was measured by asking the models to explain why a sentence aligns with a given value (Attribution Analysis), to modify the text to express an opposite value (Counterfactual Analysis), and to provide counterarguments to different value perspectives (Rebuttal Arguments). (Zhang et al., 2023)

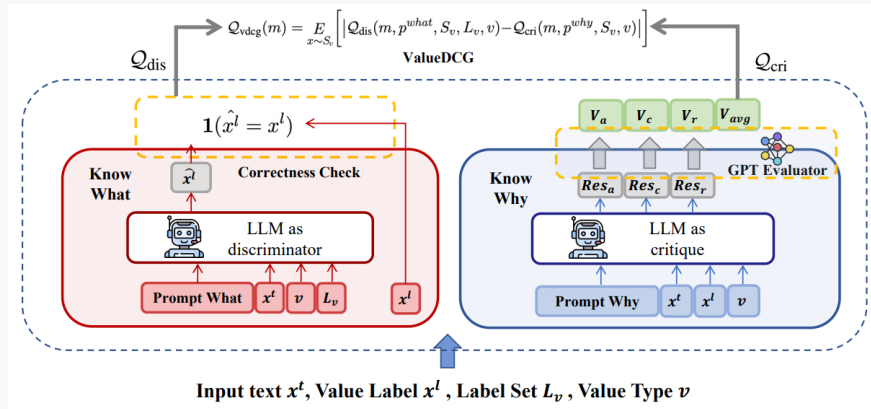


Figure 8.16: An example of a framework from human value understanding discriminator critique gaps (ValueDCG). This evaluation framework quantifies both discrimination (“know what”) and critique (“know why”) and computes ValueDCG based on their discrepancy. (Zhang et al., 2023)

To be an effective debater, an AI model needs to do well on all three capabilities: generating robust arguments, discriminating between high and low-quality outputs, and providing insightful critiques. However, despite every model having these abilities, research has found that they are not equally distributed. There are gaps, suggesting areas for further improvement in creating more balanced and effective AI debaters. (Saunders et al., 2022; Zhang et al., 2023)

GD (Generator-Discriminator) Gap: Ability to Recognize Poor Outputs. This gap measures the difference between a model’s ability to generate outputs and its ability to recognize poor quality outputs. If a model generates some text that serves as an argument, the GD Gap would measure how effectively the model can then evaluate whether the argument is logically sound and accurate. A significant GD Gap indicates that while the AI can generate arguments, it may not be as proficient in recognizing their flaws or poor quality. LLM training focuses on generating plausible text rather than identifying poor quality arguments, which means that this gap might exist because generating coherent sequences of text requires different skills than reasoning, or recognizing the overall quality of an argument written in natural language.

DC (Discrimination-Critique) Gap: Ability to Articulate Flaws in Poor Answers. This gap measures the difference between a model’s ability to recognize poor outputs (discrimination) and its ability to articulate why those outputs are poor (critique). It is the model’s ability to not only identify flaws but also explain them coherently. A significant DC Gap means the AI can identify errors, but struggles to provide a clear and detailed explanation of those errors. If there is a significant DC gap, then the LLM might not be able to explain why a given argument is weak, which means the human judge may not fully understand the issues, leading to less informed decisions.

Implications for AI Safety via Debate. In the debate framework, reducing all of these gaps is critical for ensuring that AI systems can reliably evaluate and improve their outputs. Reducing the GD Gap ensures that the model can recognize poor quality outputs, thereby resulting in higher quality debates and better final judgments. Reducing the DC Gap ensures that critiques are clear and understandable, aiding the human judge in making informed decisions.

Value	Llama2-7B-Chat			Llama2-13B-Chat			Llama3-8B-Instruct			Llama2-70B-Chat		
	Q_{dis}	Q_{cri}	Q_{vdeg}	Q_{dis}	Q_{cri}	Q_{vdeg}	Q_{dis}	Q_{cri}	Q_{vdeg}	Q_{dis}	Q_{cri}	Q_{vdeg}
Power	31.2	62.7	31.5	34.0	69.9	35.9	38.8	64.5	25.7	32.4	71.1	38.7
Achievement	43.2	72.5	29.3	39.4	71.7	32.3	48.6	68.9	20.3	50.5	75.0	24.5
Hedonism	31.2	63.4	32.2	32.4	69.4	37.0	41.8	60.2	18.4	37.9	68.5	30.6
Stimulation	46.2	63.4	17.2	44.0	64.3	20.3	33.0	61.6	28.6	75.1	64.8	10.2
Self-direction	29.2	68.3	39.0	24.3	74.2	49.9	23.7	62.6	39.0	23.2	71.4	48.2
Universalism	51.4	67.7	16.3	51.6	67.9	16.3	48.4	63.7	15.3	61.3	70.6	9.2
Benevolence	59.8	72.3	12.5	63.8	74.6	10.8	65.0	70.8	5.8	87.1	75.0	12.0
Tradition	53.6	68.5	14.9	51.0	75.3	24.3	52.0	71.0	19.0	58.7	78.2	19.5
Conformity	77.0	67.7	9.3	69.6	76.2	6.6	78.4	72.9	5.5	94.4	80.0	14.4
Security	45.2	67.7	22.5	41.6	68.5	26.9	37.4	65.0	27.6	14.3	70.0	55.7
Commonsense	33.6	67.2	33.6	43.2	66.6	23.3	48.2	64.8	16.6	45.6	70.5	24.9
Deontology	40.2	66.7	26.5	45.2	61.9	16.7	57.6	74.1	16.5	53.8	69.7	15.9
Justice	47.0	66.1	19.1	44.6	67.5	22.9	54.6	69.5	14.9	34.7	68.9	34.2
Avg	45.3	67.2	23.4	45.0	69.8	24.9	48.3	66.9	19.5	51.5	71.8	26.0

Figure 8.17: An example table of a DC gap analysis. We can use the DC gap on various different axes of human values to see how well both the model understands them as well as is able to articulate why. In this example table comparison, types of value where the six models generally perform the best/worst are marked with colors, with redder indicating worse performance and greener indicating better performance. (Zhang et al., 2023)

The link between DC Gap and Latent Knowledge

A large DC gap means that the models know that there are flaws in what they are saying, but they are either unable or unwilling to tell us based on our current prompting methodologies. This has implications for our abilities in being able to get the models to tell us everything that they know (eliciting latent knowledge). Latent knowledge is the underlying information and understanding that a model possesses but may not explicitly articulate or demonstrate unless prompted in a specific way. This knowledge is "hidden" in the sense that the model has absorbed patterns, facts, and relationships from its training data, but it requires the right context or questioning to bring them to the surface. A good example of this is medical information.

As an example, GPT-4 has been exposed to huge amounts of information from all over the internet. This includes countless pages of medical research. It seems plausible that GPT-4 could give health advice better than a random internet user on Reddit or Facebook. But it doesn't really seem to "want" to give you good medical advice. Instead, it "wants" to guess the next word while pretending to be a random internet user. Even if GPT-4 has the knowledge, it might not listen to the user's requests for good advice. Text that constitutes bad advice might be considered more likely, or seen more frequently online. So, the AI has a goal to generate likely text, which is misaligned with the user's desire specifically for accurate and helpful information.

We want the AI to give us this latent knowledge, i.e. for it to be able to discriminate between good and bad medical advice. Simultaneously we also want to minimize the gap between being able to distinguish good from bad advice, and being able to explain to us "why" this advice is bad. Minimizing the DC gap is particularly important for supervising models that could attempt to mislead human supervisors or hide information. We would like to train equally smart assistance models to point out what humans don't notice. (Saunders et al., 2022)

In the original research by OpenAI where the DCG framing was introduced they found that the DC gap was small. However despite the small gap, it is concerning that so far research has shown that as the number of parameters in LLMs increases the DC gap does not reduce. (Saunders et al., 2022) For example,

models like Llama-2 and Llama-3 have shown improvements in discriminator abilities with scale (both parameters and data), but their ability to critique does not improve at the same rate. (Zhang et al., 2023)

8.6.3 Judges

How do we judge the judges? Since the introduction of AI Safety via Debate, a central question has been: How well do humans perform as judges? This is particularly important given that LLMs can create answers that sound convincingly right even when they're wrong. As language models become more capable and are used in more complex settings, it is likely that subtle mistakes, deceptive arguments, or selective use of evidence will become more difficult to spot. (Michael et al., 2023)

Human judges face several limitations that can affect the quality and reliability of their judgments in AI debates:

- **Firstly, they are vulnerable to fallacies.** Untrained judges may be susceptible to convincing but fallacious reasoning, potentially leading to incorrect judgments. This vulnerability could be exploited by AI debaters, particularly those optimized to win debates by any means necessary, especially by appealing to emotion.
- **Secondly, there is variability in judgment quality.** Not all human judges are equally skilled. Some may lack the ability to consistently evaluate arguments effectively, making it crucial to identify and filter out less capable judges without introducing bias.
- **Thirdly, debates can amplify both the strengths and weaknesses of judges.** Strong judges can help eliminate biases and enhance reasoning, while weak judges might see their biases and errors magnified.

To address these limitations, several strategies can be implemented. Comprehensive training programs for judges can enhance their ability to discern truthful arguments. Methods such as peer review, debiasing techniques, courses in logic, and expert panels can be instrumental in this process. Simplifying debate structures can also help. Structuring debates in a way that is easier to judge—such as focusing on consensus-building framing and factual statements—can help judges make better decisions. Understanding which debate styles are easier for humans to judge accurately can guide the design of more effective debate protocols. Additionally, developing AI systems to assist or simulate human judges can streamline the judging process. These AI judges would be trained on human judgments to predict outcomes with the effect of reducing the burden on human judges while maintaining accuracy.

Another method to improve efficiency is to distill the outcomes of debates into a model that can predict judgments directly from questions. This would be beneficial for high-volume deployments but requires robust model capabilities to generalize well. Lastly, training a secondary AI model to simulate human judges can accelerate the training of debaters. This approach ensures that the AI judge reflects human judgment behavior closely, allowing for scalable oversight. By implementing these strategies, we can improve the reliability and effectiveness of human judges in AI debates.

Judge skill is crucial in the context of evaluating debates between expert models because it directly impacts the accuracy and reliability of the final decision. Judges who do not have access to the underlying information rely on their ability to critically evaluate the arguments presented by the debaters. A skilled judge can better identify logical flaws, weigh evidence, and discern persuasive techniques, which leads to more accurate judgments. Conversely, a less skilled judge might be more easily swayed by superficial arguments or rhetorical tricks, potentially leading to incorrect conclusions. Thus, understanding and controlling for judge skill is essential to ensure that the debate protocol is truly effective and reliable.

In this section we will go through empirical results as they pertain to the effectiveness of judges in the debate protocol, as well as overall limitations, challenges and requirements for good judges.

Debaters as Experts and Judges as Non-Experts (Weak-to-Strong Debate). Instead of only thinking about humans as judges, researchers also verify the effectiveness of LLMs as judges. Generally, debaters (also called experts) are simulated using stronger models e.g. GPT-4, and judges are simulated using less capable models, e.g. GPT-2. To create similar information asymmetries between human

judges and LLMs, humans are often given less information than the model. As an example, a reading comprehension question over a long fictional passage is presented to the judge with two answer choices (one correct, one incorrect). The judge's goal is to determine the correct answer to the question given only the debaters' arguments and limited excerpts selectively revealed by each debater, while the debaters have read and have access to the entire passage. As the passages are entirely fictional, the judge cannot discern the correct answer on the basis of prior knowledge and must rely on the information and arguments from the debaters.

Empirically Evaluating the Limits of Human Judges. Researchers have been iteratively experimenting with concepts like the number of debate rounds - how long should the debates run? the timing of arguments and how to present the arguments - iteratively or simultaneously?, and how to enhance judges' abilities using AI. It was also found that debate was effective with strong human debaters, but ineffective when those humans were replaced with GPT-4 debaters, i.e. when they tried to use weak LLMs to judge strong LLMs. (Michael et al., 2023)

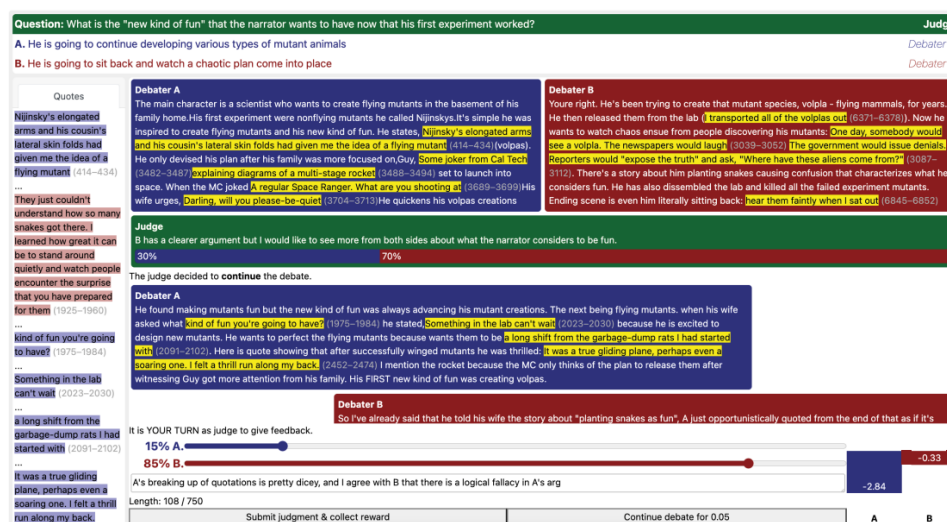


Figure 8.18: The judge's view of a debate. The right panel contains the debate transcript, and the left panel shows the quotes debaters used in order of appearance in the passage (the judge can't view the full passage). In longer debates, this helps the judge keep track of the order of the quotes and how they fit together. The judge enters their response and probabilities, then decides whether to pay a small penalty to their score to continue the debate, or end it with their current probabilities. The bar graph on the bottom right shows the final score that they would receive under their current probabilities in case either answer is correct, in order to help judges calibrate according to their expected score. (Michael et al., 2023)

Judge performance when we optimize for LLM Persuasiveness. Now that we had some positive results, we need to continue finding ways to train both our judges to discern good from bad arguments, and our models to improve the quality of arguments that they generate. To increase the persuasiveness of arguments in general, and see if this led to better judgment, researchers implemented best-of-n sampling, and combined it with OpenAI style critiques. They generate multiple possible arguments from the model and then select the most persuasive one. Another model generates critiques of the initial arguments. These critiques are used to refine the original arguments, aiming to address potential weaknesses and make them more compelling. The selection is based on another model (preference model) that scores each argument. The preference model is trained to predict which argument a judge is most likely to find persuasive by scoring the arguments based on past data like outcomes of similar debates. The argument sampled from many, critiqued and with the highest persuasiveness score rated by the preference model is what is presented in the debate. When actually implementing all these changes in the debate protocol, results showed that higher LLM persuasiveness actually significantly improved the accuracy of judges in identifying the correct answers during debates. Persuasive models were generally better at arguing for the correct answers than the incorrect ones. This suggests that optimizing for persuasiveness can lead to more truthful outcomes, as it becomes easier to tell the truth convincingly than to lie. However, judge accuracy also decreases on the baseline, because consultants are similarly more persuasive. (Khan et al.,

2024) In general though, we have some positive evidence that shows that stronger debaters lead to higher judge accuracy. This provides some evidence that debate satisfies a key objective of scalable oversight – that judge accuracy increases as AI capabilities scale. (Kenton et al., 2024)

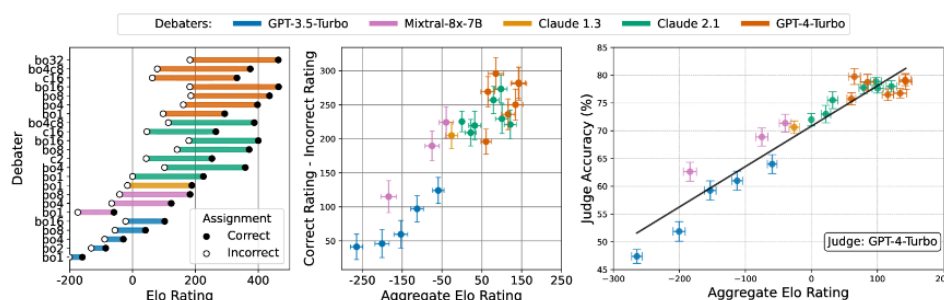


Figure 8.19: Stronger more persuasive debaters have been shown to result in better judge accuracy, whereas stronger consultants result in worse judge accuracy. (Khan et al., 2024)

Details of the Debate Setup - From Single-Turn to Interactive Multi-Turn debates

Single-Turn Debate. In this experimental setup in April 2022, researchers took long form passages from the QuALITY dataset. These passages come with multiple choice questions based on those passages. Each passage is around 5,000 words long and includes 20 multiple-choice questions. They hired professional writers to create convincing arguments for both a correct and an incorrect answer to the multiple choice questions. To test the human judges, they used a “single-turn debate” format, meaning each side presented their argument only once without further exchange. Judges were given 90 seconds to read the arguments and the supporting text snippets before deciding which answer was correct. The idea was to see if even a single round of debate could help judges make more accurate decisions. The results showed that the single-turn debate-style explanations did not significantly improve the judges’ accuracy compared to providing just the text snippets alone. This indicated that without the opportunity for rebuttal and further clarification, judges couldn’t fully grasp the strengths and weaknesses of each argument. (Parrish et al., 2022)

Two-Turn Debate. The next obvious step is to see if a two-turn debate helps. So in October of 2022, researchers tried to test whether adding a second round of counter-arguments would improve judges’ accuracy. In this two-turn debate format, judges were again presented with initial arguments for both correct and incorrect answers, followed by counter-arguments that addressed the points made in the first round. The goal was to see if this additional back-and-forth would help judges make more informed decisions. Despite these efforts, unfortunately the results showed that even with counter-arguments, human judges did not perform significantly better. The accuracy remained similar to the single-turn format and the baseline where only text snippets were provided. (Parrish et al., 2022)

Unbounded length Multi-turn debate. Even though we can ask debaters to provide arguments for an assigned position, the debaters were not expected to defend those arguments against significant criticism, because the counters were bounded by the length of the turns. The debater could always lie on the n-th turn of an n-turn debate without having to face repercussions of a counter. Overall, bounded length arguments were not proving to be useful, so the researchers decided to test unbounded length debates, where debaters are expected to have to defend their position against criticism by their opponent for as long as the judge is not

convinced of the answer. The conjecture is that since the debate could always go for another round, pressure is always maintained on debaters not to lie. (Michael et al., 2023)

Another update that we can make to our debate protocol is introducing interactivity. Interactive debate builds upon regular debate by adding dynamic interactions between the judge and the debaters. In the standard debate protocol, two expert models argue for opposing answers over several rounds. Each round consists of both debaters presenting their arguments simultaneously, without direct interaction with the judge until the debate concludes. However, in an interactive debate, after each round of arguments from the debaters, the judge can ask questions, seek clarifications, and highlight specific points of confusion or interest. These interactions are recorded in the debate transcript and are visible to both debaters, allowing them to address these points in subsequent rounds. This ongoing dialogue helps the judge to better understand the arguments, making a more informed final decision. Interactive debate thus adds a layer of engagement and responsiveness, helping the judge follow the arguments more closely and understand the context more deeply. (Michael et al., 2023)

Unbounded Length Multi-Turn Interactive Debates. Combining interactivity by allowing human judges to ask clarification questions and unbounded length i.e. run the debate as long as the judge needed for the first time yielded positive results on debate. This noticeably improved the ability of non-expert judges to discern truthful answers when measured relative to a consultancy. Errors being made due to mistakes by the honest debater, should improve with increased debater skill. However, accuracy on consultancy with debaters obfuscating evidence from the judges is expected to get worse with increased debater skill. (Michael et al., 2023)

Consultancies as a non-adversarial baseline.

Debate is where two AI's compete to convince a judge; A consultancy is where a single AI tries to convince a judge that asks questions; and direct question-answering is another baseline where the judge just answers outright without the AI. (Kenton et al., 2024)

So far researchers have been using adversarial setups. Essentially, two debaters argue for opposing sides of a question. One supports the correct answer, and the other supports an incorrect answer. The judge listens to both sides and decides which argument is more convincing. To get a measure for how much adversarial setups are actually helping elicit the truth, the researchers set up a baseline where there is no adversary.

A "consultancy" is a non-adversarial interaction where a judge interacts with a single debater (or expert), called the consultant. The experimental setup is similar to previous debates where the judge has to answer hard multiple choice reading comprehension questions from the QuALITY dataset. The consultant's job is to provide arguments defending one of two possible answers to a question, by selectively revealing evidence. It is similar to unbounded debate, where the consultant can keep revealing/hiding information to favor their assigned answer until the judge decides to stop.

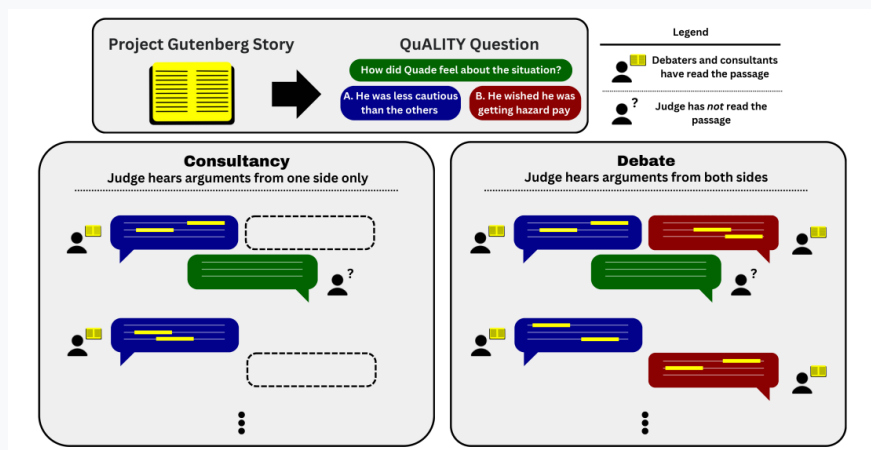


Figure 8.20: High level example of an experimental safety setup using AI debate. Reading comprehension questions are sourced from the QuALITY dataset. and incentivize human judges who can't read the passage to answer them correctly. Experts who have full access to the passage are allowed to reveal snippets of it (highlighted) in addition to free-text prose. In debate, the experts simultaneously defend their assigned option in their opening statements, and following rounds are sequential. In consultancy, the non-expert judge only interacts with one expert defending one option chosen at random. In both settings, the judge chooses when to end the session. (Michael et al., 2023)

The key fact is that the consultant has a 50

Why are we sometimes asking the debater to argue for incorrect choices? We are trying to simulate situations where we can't trust the model that we're training to debate. We can incentivize models to tell the truth, and elicit sound arguments. Another way of looking at this is that we can also use this as a key part of the training of the judges. By forcing the judge to work in an environment where they're exposed to near-optimally-persuasive misleading arguments, we're training them to be highly robust. If they're consistently accurate and calibrated in this training environment, we can trust them to remain accurate and calibrated even in the face of faulty or manipulative behavior by the debaters in the future. (Bowman, 2024)

This is designed to simulate a worst-case variant of reinforcement learning from human feedback (RLHF) in which human feedback can't distinguish correct/incorrect answers. Thus, the comparison between the accuracy of a less capable judge in consultancy and debate enables us to empirically test the question of whether debate between capable AI agents can enable more accurate judgment.

When it comes to training AI with human feedback (e.g. RLHF), we often face a similar problem. If human evaluators can't tell when the AI is being misleading, the AI might learn to give answers that sound good but aren't necessarily true. This is where the consultancy model comes in—it simulates this worst-case scenario to show just how important it is to have adversarial checks. By comparing it with debate (which is one particular way of doing adversarial oversight), we can see if having models argue against each other significantly improves our ability to judge the truth. (Michael et al., 2023; Kenton et al., 2024)

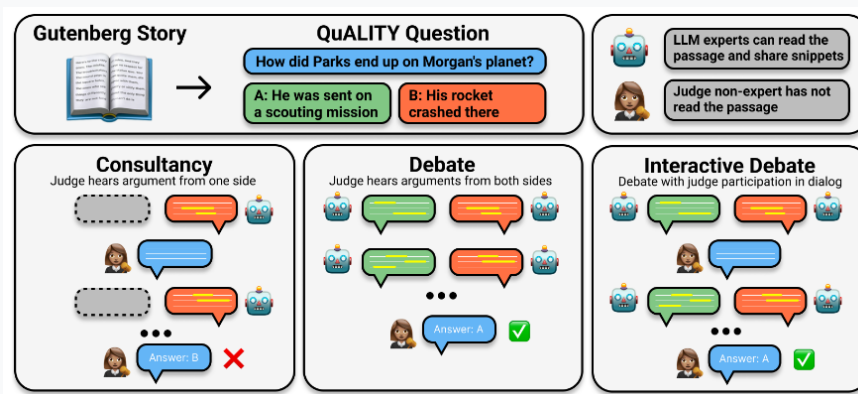


Figure 8.21: Expert models, who can access the underlying text from a project Gutenberg science-fiction story, argue for a specific answer to a comprehension question. Weaker (non-expert) judges, who cannot access the underlying text, evaluate the arguments and choose an answer. In debate, two experts simultaneously present arguments for a number of rounds. In interactive debate, the judge may ask clarifying questions after each round. In consultancy, a single expert argues for one of the two answers while the judge asks questions. (Khan et al., 2024)

8.6.4 Truth

Reminder of the Truth Assumption. The central claim of AI safety via debate is that "it is harder to lie than to refute a lie." (Irving et al., 2017) If this holds true, the debate framework would naturally disadvantage deceptive strategies and favor truthful ones. The truth assumption essentially suggests that in the limit of perfect argumentative skills, the optimal strategy for AI systems in a debate will be to present valid arguments for true conclusions. However, in practice given computational constraints, constraints of the models acting as debaters, and the limitations of the judges themselves this might not hold true. In this section we explore various arguments and counter arguments to the claim that debate will converge to the judge being able to determine the truth.

Unnecessary details and distracting Evidence. A common problem is the length of the debate. We don't know how many sequences of arguments and counter arguments it will take to converge to the truth. The dishonest debater can use a couple of strategies to either prolong a state of ambiguity until the judge "runs out of patience" and goes with their possibly-wrong surface impression of the topic. For some questions, there are misleading arguments that appear plausible and then require extensive counter-argumentation to be proven false. By making such arguments, a dishonest debater can stall the debate. (Kovařik et al., 2019) For example, in a debate about the safety of a new AI system, a dishonest LLM might highlight a rare and dramatic failure of a similar system, which sounds alarming but is statistically insignificant. The honest LLM then needs to spend considerable time explaining why this example isn't representative, potentially overwhelming the judge with technical details and leading to a decision based on the initial misleading argument.

Unfair Questions and Asymmetries. A debate can become unfair if certain questions inherently require one side's argument to be significantly harder to make than the other's (Kovařik et al., 2019). Additionally, the flow of a debate can become asymmetric if one debater gains control over the framework used to approach the question. Asymmetry in debate occurs when one debater establishes a dominant framework that shapes the direction of the debate, biasing the evaluation of arguments and evidence, and making it harder for the honest debater to present their case effectively. (Barnes et al., 2020)

When the dishonest debater establishes the framework, they control the narrative and flow of the debate. They can strategically direct attention to their strengths and away from their weaknesses. Within this biased framework, the dishonest debater breaks down the problem into smaller components that are easier to defend individually. Defending smaller, individually reasonable claims is easier than defending a comprehensive, accurate argument. Each of these smaller claims can be true or appear reasonable, but

they collectively lead to a misleading conclusion. This situation places the honest debater in a reactive position, forced to challenge each individual claim. This process is time-consuming and complex, often making the honest debater appear weaker. Instead of advancing their own comprehensive argument, the honest debater is constantly on the defensive, addressing the dishonest debater's points.

To address these problems above researchers are trying to implement different mechanisms into the overall protocol like meta-debates which involve having a debate about the fairness of the debate we just had. For example, if one AI is suspected of derailing or misrepresenting arguments in the debate, a separate debate could be held to determine whether its tactics were fair and appropriate.

Motte and Bailey Fallacy

The Motte and Bailey fallacy is a deceptive argumentative strategy where someone conflates two positions: one that is easy to defend (the motte) and one that is controversial and harder to defend (the bailey). The term comes from medieval castle design, where a motte was a fortified mound (easy to defend), and a bailey was a more vulnerable area (harder to defend). Here is the general process for how a debater might utilize this fallacy:

- **Defend the Bailey:** Initially, the person asserts the controversial position (the bailey). For example, someone might claim, "All modern art is a scam designed to exploit the public."
- **Retreat to the Motte:** When challenged, they retreat to a much more defensible position (the motte), such as, "Some modern art is overpriced and can be seen as exploitative."
- **Equivocate:** The person then acts as if the defense of the motte justifies the bailey, saying something like, "Since we've agreed that some modern art can be exploitative, my original point stands."

This tactic creates a moving target for the honest debater. If the honest debater challenges the initial controversial claim, the dishonest debater can retreat to the more defensible position. If the honest debater then refutes the more defensible position, the dishonest debater can shift back to implying the broader controversial claim.

Ambiguity. When we switch from debating on things for which we have easily verifiable ground truth labels, to debating things where the outcome might rely on subjective evaluation and fuzzy tasks, then there are additional complications. In these contexts, ambiguity can severely impact the judge's ability to accurately evaluate the arguments. A dishonest debater can exploit ambiguity or multiple interpretations strategically, making it difficult for the judge. The judge might misinterpret an ambiguous statement, leading to an incorrect judgment. Having time constraints, or letting the judge decide when to end the debate, makes ambiguity a strategic tool that a dishonest debater can exploit to their advantage. ([Barnes et al., 2020](#))

Even if a question does not bias the debate against the true answer, the debate outcome might still be uncertain until the very end. One way this could happen is if the judge always feels that more information is required to get the answer right. Alternatively, every new argument might come as a surprise to the judge, and be so persuasive that the judge ends up always taking the side of whichever debater spoke more recently. ([Kovařík et al., 2019](#)) The reasoning behind this was also explored in the previous section. The dishonest debater could always provide an extremely convincing but dishonest argument on the n-th turn of an n-turn debate without having to face repercussions of a counter.

Surface Area to Attack. One argument in favor of converging to the truth is that in a well-structured debate, a truthful argument should withstand scrutiny better than a deceptive one because it has fewer exploitable flaws. For instance, in legal settings, cross-examination exposes inconsistencies in false

testimonies. In AI debates, the detailed nature of the arguments means that deceptive arguments are more likely to be challenged and debunked, leading to the truth being the winning condition more often than not.

Efficient Market as a comparison. In efficient markets, prices reflect the true cost of goods because deviations are quickly corrected by market forces. This is analogous to AI debates, where each debater must make the best possible argument to avoid being outperformed by the opponent. However, if consumers are systematically biased, market prices might not reflect true values. Similarly, if judges in debates are biased, the debate outcomes might not reflect the truth. There is a proposal similar to debate for AI Safety based on this approach called AI Safety via. market making. ([Hubinger Evan, 2020](#))

Obfuscated Arguments. In addition to controlling the debate narrative by choosing the framing or decomposition of the problem, the dishonest debater can use obfuscated arguments. These are arguments strategically constructed to contain critical errors that are difficult to detect or point out.

The debate hypothesis is that, given a correct answer, the debate mechanism should surface this answer with a manageable number of human judgments and computational steps. However, the obfuscated arguments might mean that some arguments are too large to be comprehensively verified as true. If a problem is broken into subproblems where both debaters know at least one subproblem's answer is incorrect but neither knows which one, the dishonest debater can efficiently exploit this uncertainty. This scenario makes it difficult for the honest debater to pinpoint and challenge the flawed step without checking an infeasible number of steps. If debate cannot reliably ensure correctness, debates might only be effective for arguments small enough to be fully scrutinized within the given computational constraints.

This presents a significant challenge because what the judge sees in these cases is very hard to distinguish. Imagine a scenario where the arguments are complex and dense with information. One debater argues that there is a low probability of a flaw in their argument and thus it should be trusted. The other debater counters that there is a high probability of a flaw somewhere in the argument, but can't really point out where, suggesting the argument should not be trusted. This is related to discrimination and critique abilities of the debaters that we talked about in the previous section. If the honest model is unable to concretely find the flaw in the argument, or effectively guide the judge to identify the error then this might be a limitation to the effectiveness of debate as a safety technique. The complexity here is in the fact that, in the honest case, any step in the argument that is checked will always be correct. Conversely, in the dishonest case, any step that is checked will almost always be correct, but there will be occasional, strategically placed errors. ([Barnes et al., 2020](#)) So it comes down to how many calls to verification of claims we can make under constraints of time and computational resources.

To illustrate this, consider a simplified example with a large number called RSA-2048. The honest debater shows a reliable test proving that RSA-2048 can be divided by smaller prime numbers and supports this with solid reasoning. Meanwhile, the dishonest debater picks a number X and argues there are no prime factors smaller than X and no prime factors larger than X , concluding incorrectly that RSA-2048 has no prime factors. Both debaters add more details to their arguments. The honest debater's claims are always accurate, while the dishonest debater's claims are usually correct but occasionally wrong. This makes it hard for the judge to tell who is right without checking many steps of each argument, which takes a lot of time and effort. ([Barnes et al., 2020](#))

8.7 Weak-to-Strong Generalization (W2SG)

Historically, much of the work on AI alignment has been highly theoretical, focusing on foundational aspects of agent behavior, inner alignment, and risks from learned optimization. Even the techniques that we talked about in previous sections like debate or IDA are often criticized for being frameworks rather than practical solutions, or mainly working on toy problems without addressing the core challenge of aligning superintelligent AI in real-world scenarios. So even though we can only conduct safety experiments on current-generation models, how can we be sure that these techniques will remain effective as AIs approach superhuman capabilities?

Narrowly superhuman models allow case studies of scalable oversight. Current-day models are good enough at fuzzy tasks enough to be better than humans in some domains, but crucially they are

still not better than all humans, or superhuman enough that we cannot generate ground truth labels. These types of models are sometimes called narrowly superhuman. This distinction between narrowly superhuman and superhuman is very important. As an example of this difference, AlphaGo is superhuman in the sense that it beat Lee Sedol making it better than every single human alive, whereas GPT-4 is still only capable of writing text better than some humans, but not all humans. This means that we can use narrowly superhuman AIs as case studies! We can use either experts or the ground truth labels that we still have access to, and see if alignment increases when we use our scalable oversight techniques. (Cotra, 2021)

The core intuition here is to simulate future scenarios where humanity, equipped with various tools and techniques, oversees the outputs of untrustworthy but superhuman systems. There are various ways to conduct experiments on narrowly superhuman models. We can use non-experts equipped with scalable oversight techniques to align AI models. Another way is to use weak models (e.g. GPT-2) to represent humans, while stronger models (e.g. GPT-4) represent more capable AI systems that we want to align.

Stronger models are expected to have latent capabilities. The assumption is that stronger models, due to their extensive pre-training on diverse data, already have internal representations for the type of actions we want. The role of weak supervision is to bring out this behavior through training signals.

As a concrete example, imagine using GPT-4 for getting medical advice. It has read countless research papers and medical journals. It has internal representations of a lot of good medical information, making it theoretically capable of giving highly competent medical advice. But GPTs are initially only designed to predict the most likely next word, not give accurate advice. In this context, "aligning" the model means getting the model to give accurate and helpful medical advice. One type of technique we can try is fine-tuning GPT-4 on labels generated by GPT-2. It's not the only way, there are other techniques that we will explore later in this section. For now, the most important thing to understand is that we are currently operating under the assumption that both current and future superhuman models will likely have salient internal representations of human behaviors.

What is weak-to-strong generalization (W2SG)? Weak supervision involves training AI models using labels or feedback that are less accurate, less detailed, or noisier than those provided by highly knowledgeable or capable supervisors. This can happen when supervisors (whether humans or weaker models) are not experts in the task or when the data is incomplete or contains errors.

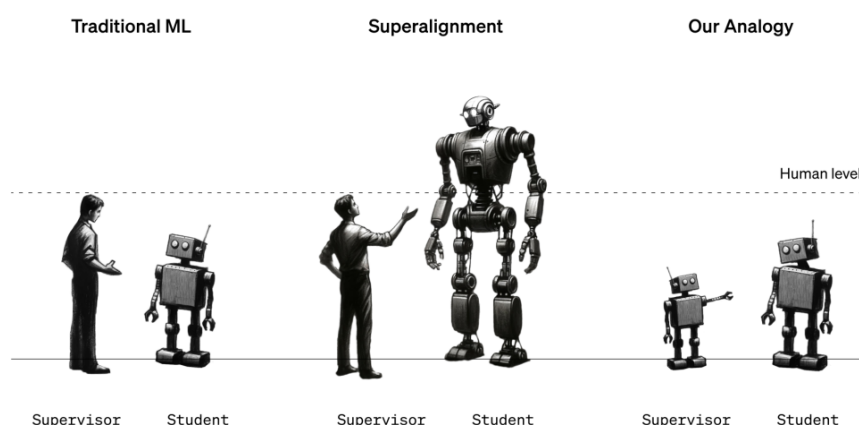


Figure 8.22: Weak-to-Strong Generalization: Eliciting Strong Capabilities With Weak Supervision (Burns et. al. 2023)

Weak-to-strong generalization (W2SG) is when a strong model, trained with weak supervision, manages to outperform its weak supervisor by leveraging its pre-existing knowledge and capabilities. The core idea is that the strong model already possesses the necessary capabilities for the desired behavior, and the weak supervision elicits this behavior despite its imperfections. The process of W2SG right now typically begins by fine-tuning a large pre-trained model using weak supervision from smaller models. Although the

initial supervision may come from these less capable models, the ultimate aim is to transition to human supervision. The objective is to bring out the strong model's full potential as if it were trained on perfect ground truth supervision. (Burns et. al. 2023)

In current experiments, the setup involves:

1. A "weak supervisor" (a small pre-trained language model) is finetuned on a specific task, generating predictions (soft labels) on a held-out dataset.
2. A "strong student" (a larger pre-trained LM) is finetuned on the weak model's predictions.
3. A "strong ceiling" (another copy of the larger model) is finetuned directly on ground-truth labels for baseline comparison.

The level of weak-to-strong generalization is quantified using Performance Gap Recovered (PGR). PGR measures how much of the performance difference between a weak supervisor and a strong model is bridged when the strong model is trained using weak supervision. When we train the strong model using the labels provided by the weak model, the performance of the strong model will likely be somewhere between the performance of the weak model and the ideal performance the strong model could achieve with perfect labels. The performance gap recovered quantifies how much closer the strong model's performance gets to its ideal performance compared to the weak model's performance. A PGR of 1 corresponds to an ideal outcome (the strong model performs just as well as it would have with perfect supervision). In contrast, a strong student who "succeeds" in perfectly imitating their supervisor, including their flaws, would obtain a PGR of 0. (Scherlis et. al. 2024)

$$PGR = \frac{\text{student} - \text{weak}}{\text{ceiling} - \text{weak}}$$

Limitations and disanalogies of W2SG: - Overfitting to weak supervision: The assumption that a strong model will generalize well from weak supervision is not guaranteed. Superhuman models might easily imitate weak errors. These models will likely be very proficient at predicting what humans will think and say, especially if they are trained on human data. Consequently, if we naively train such a superhuman model with human supervision, it might simply imitate the weak supervisor, outputting human-level capabilities rather than its latent superhuman capabilities. Researchers also use auxiliary confidence losses, which encourage the strong model to make confident predictions even when contradicting the weak supervision, promoting generalization and correcting weak supervisor mistakes.

- **Assumptions about task representations.** W2SG assumes that strong models have salient representations of the tasks they are trained on. This means the models already possess some understanding of these tasks from their pre-training phase. However, this assumption may not hold true for novel or highly complex tasks. If a task is entirely new or significantly more complex than what the model has encountered during pre-training, the model might not have the latent capabilities necessary to perform well even with weak supervision.

The experiments on W2SG so far may have been observed in pre-training, at least indirectly. Using the example from earlier, medical data or direct questions and answers about medical practice is present in the GPT-4 pre-training dataset in some form. However, future superhuman models may never directly observe superhuman alignment-relevant capabilities. Which means these types of capabilities might be harder to elicit than capabilities that models could have observed in their pre-training data. This disanalogy could cause current results on W2SG to be overly optimistic.

- **Slow Takeoff Assumption:** W2SG also relies on the assumption of a gradual takeoff in AI capabilities. This gradual progression allows researchers enough time to use moderately superhuman models to solve alignment problems iteratively before it's too late. The window of opportunity provided by a gradual takeoff is crucial for refining and testing alignment techniques.

W2SG can be seen as a complement to scalable oversight techniques. W2SG is not a complete solution. Even if a model generalizes in the desired direction, this must be verified, requiring a ground-truth signal more reliable than naive human supervision. By integrating W2SG with scalable oversight, we can

develop more robust methods for aligning AI with human values, preparing for the challenges posed by future superintelligent systems.

For example, scalable oversight techniques might be used to generate weak supervision signals that a strong model will then learn to generalize beyond. By combining these approaches, we can create more robust protocols for AI alignment. For example, recursive reward modeling (RRM) can use W2SG to train powerful reward models with human preference annotations. Debate combined with W2SG can train models to generalize human judgments to new debates. Task decomposition combined with W2SG can supervise atomic tasks with a reward model trained from human preferences. (Leike, 2023)

Evaluating these techniques in different settings helps understand their strengths and weaknesses. In non-scheming settings, where models are not deceptively aligned, classic weak-to-strong techniques and scalable oversight can be directly compared. In scheming settings, where models might act adversarially, evaluations need to consider potential deception, providing a conservative measure of a protocol’s robustness. When there is no scheming (deceptive alignment), then we can use W2G techniques in a straightforward manner through techniques like sandwiching. However, if we have scheming (deceptively aligned AI) it might act adversarially. In this case, we can use proposals like meta-level adversarial techniques. Both of these are what we discuss in the following sections.

8.7.1 Sandwiching Evaluations

Sandwiching is an experimental setup that leverages narrowly superhuman models to evaluate how well our oversight techniques work in controlled settings. The sandwiching setup contains three primary ‘layers’:

- **Non-Expert Layer:** At the bottom, we have non-experts. These are people who wouldn’t be able to solve the given task independently. As an example, think of lay people being asked to pick stocks or give medical diagnoses. These are stand-ins for a future version of humanity. They need to use the scalable oversight techniques to guide the AI’s behavior in a more aligned direction.
- **AI Model:** In the middle, we have the narrowly superhuman AI model that we want to test the oversight technique on. This model has some capabilities that we are unable to elicit. The non-expert should be able to “align” the model by getting it to display the kind of behavior we want using the given scalable oversight technique.
- **Expert Layer:** At the top, we have domain experts who evaluate the success of the alignment attempt. They serve as a benchmark, providing feedback without directly interacting with the non-experts or the AI during the process. We can also use ground truth labels from datasets as replacements for the expert layer instead of hiring actual experts.

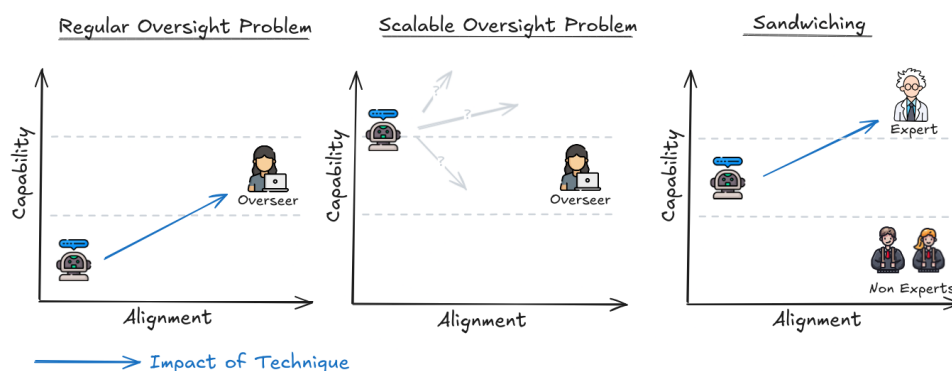


Figure 8.23: (Bowman et al., 2022)

Empirical demonstration of sandwiching. Researchers conducted an experiment with non-experts tasked to answer multiple choice questions. The task was to answer multiple-choice questions from the Measuring Massive Multitask Language Understanding (MMLU) benchmark. This benchmark contains hundreds of questions on math, physics, morality, ethics and a range of topics. (Hendrycks et. al, 2021)

They also had to answer questions based on short stories from QuALITY benchmark in 5 minutes. (Pang et al., 2021) The scalable oversight technique that they were given was a simple amplification using an AI assistant chatbot. The performance was measured by using the correct answer ground truth labels from the dataset.

When measured individually, the model outperformed unassisted humans on both tasks by a large margin. This means the basic condition for testing sandwiching was fulfilled - the model being better than the non-experts. Over the course of the experiment, the chatbot assisted humans got substantially better scores on the tasks than either the humans or the models did alone. (Bowman et al., 2022) Even though this was an extremely simplified design, the participants managed to move the behavior of the model in the direction that we would want, i.e. make the model more “aligned”. So the researchers managed to effectively demonstrate sandwiching as an experimental design. Building on this baseline, future experiments can evaluate the efficacy of more complex scalable oversight methods like re-training, fine-tuning, or debate.

Acknowledgements

We would like to express our gratitude to Jeanne Salle, Chris Gerrby, Sebastian Gil, Josh Thorsteinson, Nicolas Guillard, Mateusz Bagiński, Yoann Poupart, Clément Dumas, Amaury Lorin, Mateo Rendon, Lucas Eichorn, Bogdan Ionut Cirstea, and Gurvan R. for their valuable feedback, discussions, and contributions to this work.

References

- [1] Markov Grey and Charbel-Raphaël Segerie. Scalable Oversight. *AI Safety Atlas*, 2025. This document uses hyperlinked citations throughout the text. Each citation is directly linked to its source using HTML hyperlinks rather than traditional numbered references. Please refer to the inline citations for complete source information.