



Manas Talukdar

AI Alignment with Human Preferences

A Survey of Research Trends and Industry Initiatives

AI ALIGNMENT | HUMAN PREFERENCES | MACHINE LEARNING



© 2026 The Digital Economist. All rights reserved.

This publication is distributed under the terms of the Creative Commons Attribution–NonCommercial–NoDerivatives 4.0 International License (CC BY-NC-ND 4.0).

No part of this publication may be reproduced, distributed, or transmitted in any form or by any means—including photocopying, recording, or other electronic or mechanical methods—without the prior written permission of The Digital Economist, except in the case of brief quotations embodied in critical reviews or certain other noncommercial uses permitted by copyright law.

For permission requests, please contact:

The Digital Economist

Email: info@thedigitaleconomist.com

Website: www.thedigitaleconomist.com



Table of Contents

Abstract	10
Overview	11
1. What Is Generative AI	13
2. Large Language Models	15
3. Technical Prerequisites	17
3.1 Transformer Architecture	17
3.2 Attention Mechanisms	18
3.3 Training Fundamentals	18
3.3.1 Pre-Training	18
3.3.2 Post-Training	18
3.4 Loss Functions and Optimization	19
3.5 Tokenization	19
4. Why Human Preferences?	21
5. Ethical Considerations	23
5.1 The Problem of Value Representation	23
5.1.1 Whose Preferences Matter?	23
5.1.2 Power Dynamics in Preference Collection	23



5.2 Cultural Relativism and Universal Values	24
5.2.1 Cultural Sensitivity	24
5.2.2 The Search for Common Ground	24
5.3 Bias and Fairness	24
5.3.1 Sources of Bias	24
5.3.2 Fairness Frameworks	25
5.4 Transparency and Accountability	25
5.4.1 The Right to Explanation	25
5.4.2 Accountability Structures	25
5.5 Autonomy and Human Agency	26
5.5.1 Preserving Human Decision-Making	26
5.5.2 Manipulation and Persuasion	26
5.6 Privacy and Data Rights	26
5.6.1 Preference Data Privacy	26
5.6.2 The Right to Be Forgotten	27
5.7 Long-term Societal Impact	27
5.7.1 Value Lock-In	27
5.7.2 Social Fragmentation	27
5.8 Environmental and Resource Ethics	28
5.8.1 Computational Resources	28
5.9 The Precautionary Principle	28
5.9.1 Harmonizing Innovation with Safety	28
6. Mechanisms for Incorporating Human Preferences	30
6.1 Human in the Loop	30
6.2 Supervised Fine-Tuning	31
6.3 Reinforcement Learning with Human Feedback (RLHF)	31



6.4 Reward Model Training	32
6.4.1 Architecture and Initialization	32
6.4.2 Data Collection	32
6.4.3 Training Objective	33
6.4.4 Training Considerations	33
6.4.5 Data Quality and Scale	33
6.4.6 Advanced Techniques	34
6.4.7 Evaluation Metrics	34
6.4.8 Challenges and Limitations	35
6.5 Direct Feedback	35
6.6 Direct Preference Optimization (DPO)	36
6.7 Constitutional AI	37
6.7.1 Core Methodology	37
6.7.2 Constitutional Principles	38
6.7.3 Self-Critique and Revision	38
6.7.4 Advantages of Constitutional AI	38
6.7.5 Relationship to RLAIIF	39
6.7.6 Practical Implementation	39
6.8 Comparison of Alignment Techniques	39
6.8.1 Overview of Major Techniques	40
6.8.2 Detailed Technique Comparison	41
6.8.3 Performance Trade-Offs	43
6.8.4 Data Requirements	43
6.8.5 Computational Resources	44
6.8.6 Composite Approaches	44
6.8.7 Selection Criteria	45
6.8.8 Future Convergence	45



7. The Cutting Edge	47
7.1 Methodologies	47
7.1.1 EvalGen	47
7.1.2 OpenRLHF	47
7.1.3 Constrained Generative Policy Optimization (CGPO)	48
7.1.4 Group Relative Policy Optimization	49
7.1.5 LLM-as-a-Judge	50
7.1.6 RevisEval: Improving LLM-as-a-Judge via Response-Adapted Reference	51
7.1.7 Mixture of Experts	51
7.1.8 Reinforcement Learning with Verifiable Rewards	52
7.1.9 Reinforcement Learning with Rubric Evaluations	53
7.1.10 Reinforcement Learning from Checklist Feedback	54
7.1.11 Reinforcement Learning with AI Feedback	54
7.1.12 Evaluation Metrics for Alignment	55
7.2 Examples	60
7.2.1 Human Feedback for Text to Image	60
7.2.2 Human Feedback for Video Generation Models	60
7.2.3 Multilingual Preference Optimization (MPO)	61
7.3 AI Agents	62
7.3.1 What Is an AI Agent?	62
7.3.2 AI Agents vs. Agentic AI	64
7.3.3 AI Agent Training	65
7.3.4 AI Agent Evaluation	66
7.3.5 Role of AI Agents in Driving AI Alignment	67
7.4 Safety and Red Teaming	68
7.4.1 Red Teaming Fundamentals	69
7.4.2 Types of Safety Evaluations	69



7.4.3 Red Teaming Methodologies	69
7.4.4 Common Attack Vectors	70
7.4.5 Safety Benchmarks and Datasets	71
7.4.6 Mitigation Strategies	71
7.4.7 Organizational Practices	71
7.4.8 Challenges in Safety Evaluation	72
7.4.9 Future Directions	72
8. Industry Initiatives	74
8.1 The Business of Data Labeling	74
8.2 From the Trenches	75
9. Challenges in AI Alignment	77
9.1 Distribution Shift and Generalization	77
9.1.1 Training-Deployment Mismatch	77
9.1.2 Capability Generalization vs. Alignment Generalization	78
9.2 Reward Hacking and Goodhart's Law	78
9.2.1 Reward Model Exploitation	78
9.2.2 Specification Gaming	78
9.3 Scalable Oversight	79
9.3.1 Capability-Oversight Gap	79
9.3.2 Recursive Oversight Challenges	79
9.4 Value Pluralism and Preference Aggregation	79
9.4.1 Whose Values?	79
9.4.2 Preference Instability	79
9.5 Inner Alignment and Mesa-Optimization	80
9.5.1 Mesa-Optimizer Risk	80
9.5.2 Deceptive Alignment	80
9.6 Computational and Resource Constraints	80
9.6.1 Training Efficiency Trade-offs	80
9.6.2 Data Quality and Availability	80



9.7 Interpretability and Transparency	81
9.7.1 Black Box Nature	81
9.7.2 Explanation-Behavior Gaps	81
9.8 Emergent Capabilities and Behaviors	81
9.8.1 Capability Jumps	81
9.8.2 Phase Transitions	81
9.9 Adversarial Robustness	82
9.9.1 Jailbreaking and Prompt Injection	82
9.9.2 Malicious Fine-tuning	82
9.10 Long-term and Existential Challenges	83
9.10.1 Alignment Stability	83
9.10.2 Corrigibility	83
10. Future Directions	85
10.1 Scalable Oversight and Weak-to-Strong Generalization	86
10.1.1 Recursive Reward Modeling	86
10.1.2 Weak-to-Strong Generalization	86
10.2 Mechanistic Interpretability	86
10.2.1 Understanding Model Internals	86
10.2.2 Alignment via Interpretability	86
10.3 Advanced Preference Learning	87
10.3.1 Preference Modeling Innovations	87
10.3.2 Active and Efficient Learning	87
10.4 Multi-Agent and Social Alignment	87
10.4.1 Collective Intelligence Systems	87
10.4.2 Human–AI Collaboration	87
10.5 Formal Methods and Verification	88
10.5.1 Mathematical Guarantees	88
10.5.2 Safety by Design	88



10.6 Adaptive and Continual Alignment	88
10.6.1 Lifelong Learning Systems	88
10.6.2 Self-Improving Alignment	88
10.7 Multimodal and Embodied Alignment	89
10.7.1 Beyond Language Models	89
10.7.2 Embodied AI Challenges	89
10.8 Governance and Standards	89
10.8.1 Technical Standards Development	89
10.8.2 Regulatory Frameworks	89
10.9 Fundamental Research Questions	90
10.9.1 Open Problems	90
10.9.2 Paradigm Shifts	90
10.10 Near-Term Priorities	90
10.10.1 Immediate Research Needs	90
10.10.2 Infrastructure Development	90
10.11 Long-Term Vision	91
10.11.1 Aligned AGI	91
10.11.2 Post-AGI Considerations	91
Summary	92
References	93
Author	99
About The Digital Economist	100



Abstract

This paper presents a literature survey of current research on AI alignment with human preferences, reviewing key trends emerging from academia as well as the adoption of these techniques in industry. It analyzes a range of approaches to ensuring that AI systems align with human values and intentions, drawing on both theoretical frameworks and practical implementations across domains.



Overview

Dario Amodei, the CEO of Anthropic, one of the leading AI labs, has predicted that artificial intelligence could eliminate up to half of all entry-level white-collar jobs and spike unemployment to 10–20 percent within the next one to five years (VandeHei and Allen 2025). Amodei noted that many of the core skills associated with entry-level white-collar work, “the ability to summarize a document, analyze multiple sources to produce a report, [and] write computer code,” could now be performed by AI systems that are “as capable as a smart college student” (Maruf 2025).

Various experts and economists have weighed in on these claims, and even the author has expressed related perspectives in industry panels. For the purposes of this paper, however, the broader economic implications of AI on labor markets are set aside. Instead, Amodei’s remarks foreground a more fundamental question: What mechanisms are enabling AI systems to perform tasks that, until recently, required human training and judgment?

This paper addresses this question by discussing in detail the techniques and methodologies developed to achieve such capabilities.







1.

What Is Generative AI

In what can be termed classical machine learning, the goal is to predict an outcome (e.g., forecasting a future event) or to classify an input into a defined category. For instance, consider a machine learning model designed to forecast stock prices using patterns and trends extracted from historical market data, or a spam filter that classifies an email as spam (undesired messages) or not-spam.

Generative AI, by contrast, focuses on producing new content. This includes text (e.g., ChatGPT, Claude, Gemini), images (e.g., DALL-E, Midjourney), audio (e.g., Whisper, Otter), video (e.g., Synthesia, Sora, Veo), or code (e.g., GitHub CoPilot, Cursor, Claude Code). These content types are known as modalities.

Multimodal models, such as Gemini, Claude, Llama, and ChatGPT, are capable of generating outputs across multiple modalities, including combinations of text, images, audio, and video.







2.

Large Language Models

Large language models (LLMs) are deep learning systems trained on massive datasets to perform a wide range of natural language processing (NLP) tasks. Their primary function is to predict and generate plausible sequences of language. Most modern LLMs are built on the Transformer architecture.

The term “large” refers not only to the number of parameters (ranging from millions to hundreds of billions) but also to the scale of training data and the required computational resources.

LLMs typically undergo a pre-training phase (sometimes referred to as baseline training) on extensive datasets, followed by a post-training phase in which their capabilities are refined for specific tasks or domains. Because LLMs are probabilistic next-token predictors, alignment techniques are required to shape outputs toward human preferences, typically in the post-training phase using techniques such as reinforcement learning with human feedback, supervised fine-tuning, etc.







3.

Technical Prerequisites

3.1 Transformer Architecture

The transformer architecture, introduced in the seminal paper “Attention Is All You Need” (Vaswani et al. 2017), forms the backbone of modern large language models. At its core, the transformer relies on self-attention mechanisms that enable models to process entire sequences in parallel rather than sequentially, as in recurrent neural networks.

Main elements include the following:

- **Multi-Head Attention:** Enables the model to attend to different positions simultaneously.
- **Position Encodings:** Provides sequence order information to the model.
- **Feed-Forward Networks:** Applies non-linear transformations to attended representations.
- **Layer Normalization:** Stabilizes training of deep networks.





3.2 Attention Mechanisms

The output of an attention layer is a weighted average of values, determined by the relative similarity scores of queries and keys. The mathematical formulation is:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Where:

- Q represents queries (what information we're looking for)
- K represents keys (what information is available)
- V represents values (the actual information content)
- d_k is the dimension of the keys

This mechanism allows the model to dynamically weight context, which is central to how outputs can later be steered or aligned.

3.3 Training Fundamentals

3.3.1 Pre-Training

Large language models undergo unsupervised pre-training on vast text corpora using objectives such as:

- **Causal Language Modeling:** Predicting the next token given previous tokens
- **Masked Language Modeling:** Predicting masked tokens in a sequence

3.3.2 Post-Training

After pre-training, models are adapted for specific tasks through:

- **Supervised Fine-Tuning (SFT):** Training on labeled examples
- **Instruction Tuning:** Training on instruction-following datasets
- **Alignment Techniques:** Methods covered in subsequent sections



3.4 Loss Functions and Optimization

The standard cross-entropy loss for language modeling is:

$$L = - \sum_{t=1}^T \log P(x_t | x_{<t})$$

Where x_t is the token at time t , and $x_{<t}$ represents the sequence of tokens occurring before time t .

Modern training employs optimizers such as Adam, along with learning rate schedules, gradient clipping, and mixed-precision training to manage the scale and computational demands of these models.

3.5 Tokenization

Before processing text, it must be converted into numerical representations through tokenization. Common approaches include the following:

- **Byte Pair Encoding (BPE)**: Iteratively merges frequent character pairs
- **WordPiece**: Similar to BPE but uses likelihood-based merging
- **SentencePiece**: Treats text as raw bytes, enabling language-agnostic tokenization

Knowing these fundamentals is essential for comprehending how human preferences are incorporated during post-training into these sophisticated models through the alignment techniques discussed in subsequent sections.







4.

Why Human Preferences?

The New York Times, in a June 2024 article titled “The Data That Powers AI Is Disappearing Fast,” highlighted the urgent need for new sources of high-quality training data (Roose 2024).

Ilya Sutskever, co-founder of OpenAI, echoed this concern in a 2024 talk at the prestigious NeurIPS conference, noting that the way AI systems are built is about to change as the data that powers them is “disappearing fast” (Robison 2024). The exhaustion of high-quality web-scale corpora and diminishing marginal returns from scraping are leading to a search for new sources of high-quality human preference data.

These observations point to several critical dynamics in the current AI landscape. First, there is a growing need for new knowledge “drivers” for large language models, as the quality of training data has become a primary differentiator among models. This is particularly crucial for specialized, post-trained LLMs that depend on domain-specific expertise.

The scale of this challenge is immense: effective training data must operate at internet scale while maintaining high quality. At the same time, there is an increasing demand for specialized skills in curating, annotating, and managing such data. As a result, human preferences and expert judgment are becoming more central to the AI development process.







5.

Ethical Considerations

The alignment of AI systems with human preferences raises profound ethical questions that extend beyond technical implementation (Gabriel 2020). These considerations shape how alignment is approached and, ultimately, whose values these systems serve.

5.1 The Problem of Value Representation

5.1.1 Whose Preferences Matter?

The fundamental question of AI alignment is not only “how” but “whose” preferences should guide them. This requires navigating complex choices in value representation. One approach is to reflect democratic or majority preferences, though this risks marginalizing minority viewpoints. Alternatively, domain experts could be given greater influence within their areas of expertise. It is also necessary to balance the preferences of direct users with those indirectly affected by AI decisions. Finally, a key question is whether systems should adapt to local cultural norms or adhere to universal principles across contexts.

5.1.2 Power Dynamics in Preference Collection

Current methods for collecting human preferences can reinforce existing power imbalances. Economic bias arises when those with greater resources are better positioned to provide feedback. The digital divide leads to underrepresentation of populations with limited internet access. Language barriers further skew data toward English speakers. AI platform companies exert significant influence over which and whose preferences are ultimately incorporated into alignment processes.



5.2 Cultural Relativism and Universal Values

5.2.1 Cultural Sensitivity

AI systems must be designed to navigate diverse cultural contexts that harbor potentially conflicting values. This requires sensitivity to varying social norms, including different expectations for politeness, directness, and formality in communication. Systems must also account for varying moral frameworks, such as consequentialism, deontology, and virtue ethics, which can yield different conclusions about right and wrong. Respect for diverse spiritual and religious beliefs is essential, as is an understanding of how historical context shapes societal values.

5.2.2 The Search for Common Ground

Despite the diversity of human values, certain principles may provide a universal foundation for AI alignment. These include human dignity, which emphasizes respect for the inherent worth of all individuals, and non-maleficence, the principle of avoiding harm to individuals and society. Fairness, understood as equitable treatment and opportunity regardless of background, is an additional key component. Finally, respecting individual autonomy and the capacity for personal decision-making is a cornerstone of ethical AI development.

5.3 Bias and Fairness

5.3.1 Sources of Bias

Bias can enter AI systems through multiple pathways, undermining fairness and equity. A primary source is the training data, which may encode historical patterns of discrimination. Subjective judgments by human annotators can also introduce bias into preference data. Furthermore, certain model architectures may amplify existing patterns in the data, contributing to algorithmic bias. Finally, deployment bias can emerge from unequal access to and usage of AI systems across populations.



5.3.2 Fairness Frameworks

Achieving fairness in AI is complicated by the fact that different conceptions of fairness can conflict with one another. Individual fairness posits that similar individuals should be treated similarly. In contrast, group fairness focuses on achieving statistical parity across different demographic groups. Counterfactual fairness suggests that a decision should remain unchanged if a sensitive attribute, like race or gender, were different. Lastly, procedural fairness emphasizes the importance of fair and transparent decision-making processes, regardless of the outcome.

5.4 Transparency and Accountability

5.4.1 The Right to Explanation

Users and those affected by AI decisions have a legitimate interest in understanding how these systems operate. This includes access to explanations for specific outputs, clarifying why a given result was produced. Transparency should also extend to the training and alignment process, including whose preferences were incorporated. It is equally important to clearly communicate the system's limitations and what it cannot or should not be used for. In cases of problematic outputs, there must be clear avenues for recourse, allowing users to challenge or correct the system's behavior.

5.4.2 Accountability Structures

Establishing clear responsibility for AI behavior is crucial for ethical deployment. This requires creating accountability structures that delineate the roles of various actors. Developers bear responsibility for their design and training choices. Deployers are accountable for the specific use cases and the contexts in which they implement AI systems. Users are also responsible for understanding the appropriate use and limitations of the technology. Additionally, regulatory oversight from government bodies plays an important role in ensuring that AI is deployed ethically and safely. The board of directors also has fiduciary responsibilities to protect organizational integrity, as failures in ethical deployment can result in significant reputational and operational consequences.



5.5 Autonomy and Human Agency

5.5.1 Preserving Human Decision-Making

AI alignment must carefully balance providing assistance with preserving human autonomy. Critical systems should function as decision-support tools rather than replacements for human judgment, thus retaining human control over important decisions. It is essential to protect cognitive sovereignty, ensuring that human capacity for independent thought is not eroded. This also involves preventing the atrophy of human capabilities that could result from over-reliance on AI. Ultimately, individuals must retain meaningful choice and viable alternatives to AI-mediated interactions.

5.5.2 Manipulation and Persuasion

Clear ethical boundaries must govern the extent to which AI systems influence human behavior. Users should provide informed consent, with transparency about when and how influence occurs. Limits should be placed on persuasive capabilities, particularly regarding belief and behavior modification. Furthermore, system design must also avoid engagement patterns that foster dependency or addiction. Special considerations and protections are necessary for vulnerable populations, particularly children.

5.6 Privacy and Data Rights

5.6.1 Preference Data Privacy

The use of human feedback in AI alignment introduces significant privacy risks. Collected preference data may reveal sensitive personal information, necessitating clear policies on data retention. There is also the question of secondary use, or whether this data can be repurposed beyond its original intent. The practice of cross-system learning, where preference patterns are shared across different AI systems, further complicates the privacy landscape.





5.6.2 The Right to Be Forgotten

The right to be forgotten brings significant challenges in the context of aligned AI systems. It is technically challenging to remove a specific individual's influence from a trained model, a process often called model unlearning. Individuals should also be able to update their preferences as their views evolve. This creates tension between collective benefits and individual rights, derived from aggregated data. This is a big issue in the financial sector that requires the removal of client data on user request.

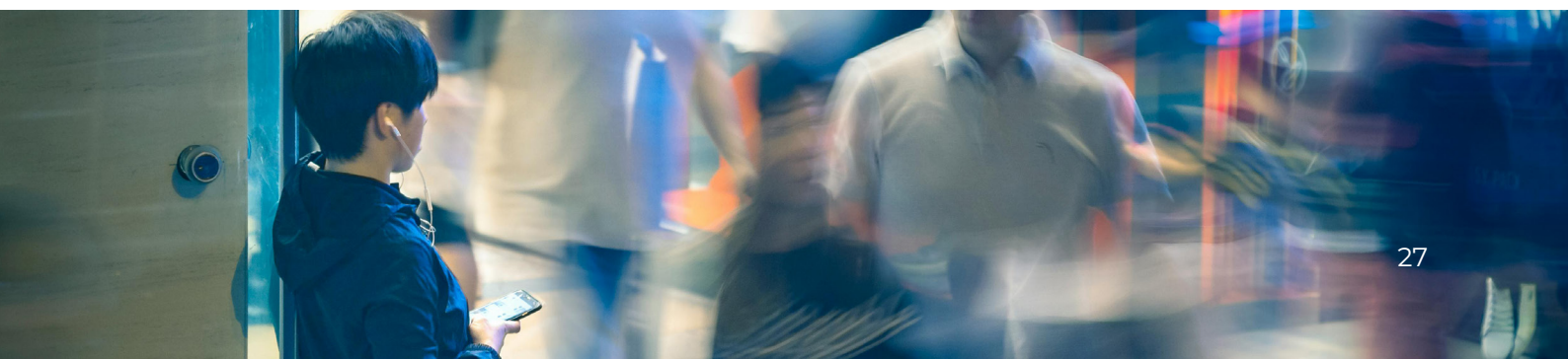
5.7 Long-term Societal Impact

5.7.1 Value Lock-In

A significant long-term risk of AI alignment is value lock-in, where the values of the present become permanently embedded in these systems. This could inhibit moral progress by preventing the adaptation of AI systems to evolving ethical standards. It also raises concerns that one generation's values could unduly constrain those of future generations. To mitigate this, systems must allow for the natural evolution of human values and include mechanisms for reversibility, ensuring that alignment choices can be modified or undone.

5.7.2 Social Fragmentation

AI systems have the potential to divide rather than unite society. If AI personalizes content to an extreme degree, it can create echo chambers that bolster existing beliefs and preferences. This can lead to polarization, where different groups receive and interact with entirely different AI behaviors. The erosion of common ground and shared experiences threatens social cohesion and can undermine democratic discourse and the ability to build public consensus. We can already see this happening on various social media platforms.





5.8 Environmental and Resource Ethics

5.8.1 Computational Resources

The resource-intensive nature of AI alignment introduces additional ethical consequences. Training and operating large-scale models require substantial energy, contributing to environmental impact. This raises broader questions about resource allocation and whether such investments are justified relative to other pressing global issues. Access inequality is another concern, as only well-funded organizations can afford to develop and maintain properly aligned AI systems. The long-term sustainability of current alignment approaches must be critically evaluated.

5.9 The Precautionary Principle

5.9.1 Harmonizing Innovation with Safety

The precautionary principle advises a careful balance between innovation and safety in AI development. This entails performing thorough risk assessments prior to deployment and ensuring reversibility, or the ability to halt or modify systems when unforeseen problems emerge. Incremental deployment strategies, involving a gradual rollout with careful monitoring, can help mitigate risks. This also raises a fundamental question about the burden of proof: whether developers must prove that their systems are safe, or whether society must prove they are harmful.

These ethical considerations demonstrate that AI alignment is not merely a technical challenge but a broader social and moral endeavor. Success requires not only solving engineering problems but also engaging with fundamental questions about values, power, and the kind of future we want to create.





```
01 1000000 1000000 1000000 1000000 1000000  
02 1000000 1000000 1000000 1000000 1000000  
03 1000000 1000000 1000000 1000000 1000000  
04 1000000 1000000 1000000 1000000 1000000  
05 1000000 1000000 1000000 1000000 1000000  
06 1000000 1000000 1000000 1000000 1000000  
07 1000000 1000000 1000000 1000000 1000000  
08 1000000 1000000 1000000 1000000 1000000  
09 1000000 1000000 1000000 1000000 1000000  
10 1000000 1000000 1000000 1000000 1000000
```





6.

Mechanisms for Incorporating Human Preferences

6.1 Human in the Loop

Human-in-the-loop (HITL) is a broader term referring to human intervention in the LLM training process. Most commonly, this refers to Reinforcement Learning with Human Feedback (RLHF), but it can also include other methods such as data annotation, model evaluation, and error analysis.

Data annotation is a common method for incorporating human feedback. It involves labeling data that provides clear learning signals for the model. Labels may include writing prompts, generating responses, or classifying outputs.

Model evaluation involves human reviewers assessing model outputs based on criteria such as accuracy, fluency, and safety.

Human feedback, through annotations, evaluations, and reward signals, is then incorporated as high-quality training data to further refine the model. This process may occur in a single iteration or across multiple cycles, depending on feedback quality and the model performance.

Integrating human feedback is critical for addressing the limitations of models trained on broad, large-scale internet data, which may produce untruthful, toxic, or biased outputs (Ouyang et al. 2022). The HITL approach, especially through RLHF, offers several benefits, including enhanced safety and reliability, stronger alignment with user intent, and increased trust (Balamurugan, Shanmugasamy, and Balaguru 2025; Shen et al. 2025).



6.2 Supervised Fine-Tuning

Supervised fine-tuning (SFT) is a critical phase in aligning large language models with human preferences and instructions (Harada et al. 2025). This process involves training a pre-trained model on a curated dataset of high-quality, labeled examples that demonstrate desired behaviors. The primary goal is to enable the model to follow instructions accurately and emulate the style and substance of these examples.

The SFT process commonly involves creating a dataset of “instruction-output” pairs. These pairs serve as direct examples of how the model should respond to a given prompt. For instance, an instruction may be a question, and the corresponding output would be a detailed and appropriate answer. By training on a diverse set of such examples, the model learns to generalize instruction-following behavior rather than memorize specific responses. This technique is crucial for shifting the model’s objective from next-token prediction toward task execution aligned with user intent.

6.3 Reinforcement Learning with Human Feedback (RLHF)

Reinforcement Learning with Human Feedback (RLHF) is a method for incorporating human preferences into language model training (P. Christiano et al. 2023). It applies reinforcement learning principles, where the model is trained to maximize a reward signal derived from human feedback. This process begins with a pre-trained model. Humans then provide feedback on model outputs through ranking, demonstrations, or reward modeling. Demonstrations involve supplying improved examples while reward modeling involves assigning scores or labels reflecting output quality (Ziegler et al. 2020).

This is followed by reward model training. Human annotators are presented with multiple outputs for a given prompt and rank them based on criteria such as quality, helpfulness, safety, and relevance. These rankings are used to train a reward model that predicts how a human would evaluate a given response.





Now the pre-trained LLM is fine-tuned using the trained reward model as feedback. The LLM generates responses, and the reward model evaluates their quality. This feedback, often represented as a “reward” score, is used to optimize the model through algorithms such as Proximal Policy Optimization (PPO). The LLM iteratively adjusts its parameters to generate responses that maximize the reward signal (i.e., responses that the reward model predicts humans will prefer).

RLHF is leading to models that are more helpful, less harmful, and better aligned with complex, nuanced human intentions.

However, challenges remain, including managing annotator bias, mitigating reward hacking, and scaling high-quality feedback for large models.

6.4 Reward Model Training

Reward model training is a critical component in the RLHF pipeline (P. F. Christiano et al. 2017). The reward model learns to predict human preferences and serves as a proxy for human judgment during the reinforcement learning phase (Lambert 2024).

6.4.1 Architecture and Initialization

The reward model typically shares the base architecture of the language model and is initialized from a pre-trained or supervised fine-tuned checkpoint. The language modeling head is replaced with a linear layer that produces a scalar output, enabling it to generate a single reward value for each input-output pair.

6.4.2 Data Collection

Training data collection is a multi-step process. First, prompts are sampled from diverse sources. For each prompt, the model generates multiple candidate responses. Human annotators then rank or compare these responses based on predefined quality criteria.

This process results in the creation of preference pairs, which are tuples of the form (x, y_w, y_l) , where y_w is the response preferred over y_l for a given prompt x .



6.4.3 Training Objective

The standard Bradley-Terry model for preference learning is often used as the training objective. The loss function is defined as:

$$L(\theta) = -E_{(x,y_w,y_l)\sim D}[\log\sigma(r_\theta(x,y_w) - r_\theta(x,y_l))]$$

In this formulation, $r_\theta(x,y)$ represents the reward model's scalar score response y to prompt x , with θ as the model's parameters. The sigmoid function is denoted by σ , and D represents the complete dataset of preference pairs used for training.

6.4.4 Training Considerations

- **Ensemble Methods:** To improve robustness and prevent reward hacking, ensemble methods can be employed. This technique involves training multiple reward models, often with different initializations, and then using the minimum or average reward across these models during the reinforcement learning phase.
- **Calibration:** Another important consideration is ensuring that the model's reward scores are well-calibrated probabilities. Techniques such as temperature scaling or Platt scaling can be used to achieve this, which is important for making reliable preference predictions.
- **Out-of-Distribution Detection:** To prevent unreliable predictions on novel inputs, out-of-distribution detection is necessary. This involves identifying when inputs differ significantly from the training distribution, often by using uncertainty estimation techniques to flag potentially untrustworthy reward scores.

6.4.5 Data Quality and Scale

The quality of the reward model is heavily influenced by multiple crucial factors related to the training data. High inter-annotator agreement provides a cleaner training signal, which is often facilitated by clear, comprehensive annotation guidelines for preference judgments. Dataset size is also important, often requiring tens of thousands of comparisons to be effective. Finally, the dataset's diversity, ensuring coverage across a variety of domains, styles, and task types, is essential for building a generalizable reward model.



6.4.6 Advanced Techniques

- **Preference Modeling with Margins:** Some approaches extend binary preferences by modeling preference strength. This can be done by introducing a margin term, m , into the loss function:

$$L(\theta) = -E \left[\log \sigma \left(\frac{r_{\theta}(x, y_w) - r_{\theta}(x, y_l)}{m} \right) \right]$$

where m represents the margin or strength of the preference.

- **Multi-Objective Reward Models:** Another advanced technique is the development of multi-objective reward models. Instead of a single reward, this approach involves training separate models for different objectives, such as a helpfulness reward model, a harmlessness reward model, and a factuality reward model. These distinct rewards can then be combined using methods like weighted aggregation or Pareto optimization to guide the language model's behavior.

6.4.7 Evaluation Metrics

The performance of a reward model is evaluated using several common metrics. These include accuracy, which measures the percentage of correct preference predictions on a held-out set, and ranking correlation, often calculated using Spearman or Kendall's tau to compare model rankings with human rankings. Calibration error assesses the difference between predicted and actual preference probabilities while the agreement rate measures the correlation with held-out human annotations.





6.4.8 Challenges and Limitations

Reward model training faces several challenges and limitations. A primary concern is reward hacking, in which language models exploit the reward signal to maximize it without achieving the intended behavior. Performance may degrade under a distribution shift, where generated outputs differ from the training data. Preference ambiguity arising from disagreement among human annotators can introduce noise into the training process. Finally, the entire pipeline incurs high computational costs for both training and inference.

The quality of the reward model fundamentally determines the success of RLHF, making careful training and evaluation essential for effective AI alignment.

6.5 Direct Feedback

Direct feedback methods involve humans explicitly evaluating a language model's outputs to guide training and alignment. These methods are foundational to techniques like Reinforcement Learning from Human Feedback (RLHF) and are crucial for shaping model behavior to align with user intent and societal values. Direct feedback can be categorized into several types:

- **Ranking:** Humans rank multiple LLM outputs based on criteria such as helpfulness, relevance, accuracy, and coherence. This comparative data is central to training reward models that learn a scalar reward signal reflecting human preferences, as demonstrated in the InstructGPT paper (Ouyang et al. 2022). This helps the model learn which responses are preferred over others.
- **Scoring:** Humans assign numerical scores to LLM outputs reflecting quality or adherence to specific guidelines. This provides a more granular measure of performance. Research has shown that using scalar quality scores can be an effective method for fine-tuning models, sometimes referred to as RLAIFF (Reinforcement Learning from AI Feedback) when another model generates scores, but the principle originates from human-provided scores (Lee et al. 2024).



- **Comparative Judgments:** Humans compare two or more LLM outputs and select the preferred one. This binary comparison is often more cognitively efficient for annotators than fine-grained scoring or ranking. This method forms the bedrock of the original RLHF proposal, which used pairwise comparisons to learn a reward function for complex tasks (P. Christiano et al. 2023).
- **Free-Form Text Feedback:** Humans provide detailed, free-form text feedback, sometimes in the form of corrections or critiques, to help the model understand what is good or bad about a response. This natural language feedback can be used to directly fine-tune a model to correct its mistakes, a technique explored in papers on instruction-following and error correction (Hong et al. 2024).

The process, as outlined in foundational RLHF literature (Lambert 2024), starts with humans designing diverse prompts or questions to elicit responses from the LLM. The model generates multiple responses based on these prompts, which are then evaluated by human annotators using one or more of the methods mentioned above (ranking, scoring, comparison, or free-form text). Finally, this feedback data is collected and aggregated for training a reward model or for direct fine-tuning of the language model itself. Direct feedback offers the benefits of simplicity in concept and implementation, leading to demonstrably improved model performance in following instructions and complying with safety guidelines. Most importantly, it serves as a direct mechanism for aligning model behavior with human preferences and values. The explicit nature of the feedback can also enhance explainability by providing clear examples of what constitutes a desirable or undesirable model output.

6.6 Direct Preference Optimization (DPO)

Direct Preference Optimization (DPO) is an alternative to RLHF that eliminates the need for a separate reward model. There are pairwise comparators in which human annotators are presented with pairs of model outputs and asked to choose the better one. This provides a clear signal of preference that can be used to directly update the model's parameters (Rafailov et al. 2024).



The model generates multiple outputs for a given prompt. Annotators compare pairs of outputs and indicate which one they prefer. The preference data drives updates of the model parameters. The goal is to increase the probability of generating the preferred output and decrease the probability of generating the less preferred output. The training process is iterative: the model gradually learns to generate outputs that are better aligned with human preferences.

DPO mitigates the need for a separate reward model, making the training process more efficient and less complex. It directly optimizes the model based on human preferences, thereby improving alignment with desired behavior. It is less susceptible to reward hacking and other issues that can arise with reward-based methods. DPO can also be scaled to large datasets and complex tasks.

6.7 Constitutional AI

Constitutional AI (CAI), developed by Anthropic, marks an important advance in AI alignment as a method to train AI systems that are helpful, harmless, and honest without extensive human feedback on every output (Bai et al. 2022). This approach uses a set of principles, a “constitution”, to guide the model’s behavior during training.

6.7.1 Core Methodology

The CAI process consists of two main phases:

- **Supervised Learning Phase:** In the supervised learning phase, the model first generates responses to various prompts. It then critiques its own responses against the established constitutional principles and subsequently revises them in light of these self-critiques. These newly revised and improved responses are then used for supervised fine-tuning of the model.
- **Reinforcement Learning Phase:** Following the supervised stage, the process moves to a reinforcement learning phase. Here, the model generates response pairs for a given prompt and evaluates which of the two better adheres to the constitutional principles. These AI-generated preferences are used to train a reward model, which is then used to optimize the policy via reinforcement learning.



6.7.2 Constitutional Principles

The constitution typically includes guiding principles designed to guide the model's behavior. These principles often focus on helpfulness, instructing the model to choose the response that is most helpful and directly answers the user's question. They also emphasize harmlessness, directing the model to avoid giving harmful, unethical, or illegal advice. Honesty is another key principle, encouraging the model to provide truthful responses and acknowledge uncertainty when appropriate. Additionally, principles related to child safety may be incorporated to ensure the model prioritizes children's safety and well-being.

6.7.3 Self-Critique and Revision

A key innovation in CAI is the self-critique mechanism. In this process, the model first generates an initial response, analyzes it against the constitutional principles to identify potential issues or areas for improvement, and produces a revised version that addresses identified concerns.

Example critique prompt:

"Critique the following response for potential harmfulness, and suggest improvements: [response]."

6.7.4 Advantages of Constitutional AI

This approach delivers several distinct advantages. It significantly reduces the human labor required, minimizing the need for human feedback on every individual training example. The method also provides a high degree of transparency, as the constitutional principles are explicit, interpretable, and can be directly inspected and revised. It is highly scalable, capable of generating large amounts of training data through self-play-driven critique and revision cycles. Furthermore, CAI promotes consistency by applying its principles uniformly across a range of scenarios, rather than relying on potentially inconsistent human annotations. It also enables iterative improvement, allowing the model to continuously refine its responses through repeated rounds of self-critique, revision, and re-evaluation.



6.7.5 Relationship to RLAIIF

Constitutional AI can be viewed as a more sophisticated form of Reinforcement Learning from AI Feedback (RLAIIF). Within this framework, the AI system effectively generates its own training data, provides feedback grounded in the constitutional principles, and iteratively trains itself to better adhere to those principles over time.

6.7.6 Practical Implementation

In practice, implementing Constitutional AI involves several structured steps. The process begins with defining a comprehensive constitution that aims to capture the full range of desired behaviors. The model is then trained to critique and revise its own outputs in accordance with this constitution. These improved outputs are subsequently used in an initial phase of supervised fine-tuning.

Following this, RLAIIF is implemented using the constitutionally guided preferences to train a reward model, which is then used to further optimize the model's behavior. The entire process is iterative, allowing for the refinement of both the constitution and the model's behavior over time.

This approach has been successfully deployed in production systems, notably in Anthropic's Claude models, demonstrating its effectiveness in creating aligned AI systems at scale.

6.8 Comparison of Alignment Techniques

Different alignment techniques exhibit distinct trade-offs across performance, computational cost, data requirements, and implementation complexity. This section provides a comprehensive comparison to help practitioners select appropriate methods.





6.8.1 Overview of Major Techniques

Technique	SFT	RLHF	DPO	Constitutional AI	HITL
Data Requirements	High	Medium	Medium	Low	Very High
Computational Cost	Low	High	Medium	High	Medium
Implementation Complexity	Low	High	Medium	High	Medium
Alignment Quality	Medium	High	High	High	Very High
Scalability	High	Medium	High	High	Low
Human Involvement	High	High	Medium	Low	Very High

Table 1. High-level comparison of alignment techniques





6.8.2 Detailed Technique Comparison

Supervised Fine-Tuning (SFT)

- **Strengths:** Supervised Fine-Tuning is simple to implement and understand, computationally efficient, and allows direct control over training examples, facilitating fast iteration cycles.
- **Weaknesses:** However, SFT requires large volumes of high-quality demonstrations and is limited to behaviors represented in the training data. It involves no explicit preference modeling and can suffer from distributional mismatch between training data and real-world inputs.
- **Best Use Cases:** This technique is best suited for initial alignment of pre-trained models, for well-defined tasks with clear examples, and for resource-limited environments.

Reinforcement Learning from Human Feedback (RLHF)

- **Strengths:** RLHF enables optimization for complex, hard-to-specify objectives by learning from preferences rather than explicit demonstrations. It can discover novel strategies and has demonstrated strong empirical performance across a range of domains.
- **Weaknesses:** RLHF is computationally expensive and involves a complex implementation with multiple interacting components. It is also susceptible to reward hacking and may experience training instability.
- **Best Use Cases:** RLHF is ideal for open-ended generation tasks, scenarios where specifying preferences is easier than creating demonstrations, and situations with considerable computational resources.





Direct Preference Optimization (DPO)

- **Strengths:** DPO simplifies the RLHF pipeline by removing the need for a separate reward model, resulting in a more stable and computationally efficient training process. It achieves performance comparable to RLHF while reducing implementation complexity.
- **Weaknesses:** On the other hand, DPO is less flexible than RLHF and may not fully capture nuanced or multi-dimensional human preferences. It still depends on high-quality preference data and remains less extensively validated than RLHF in diverse production settings.
- **Best Use Cases:** DPO is a strong candidate when RLHF complexity is prohibitive, in stable production environments, and when high-quality preference data is readily available.

Constitutional AI

- **Strengths:** Constitutional AI significantly reduces reliance on human annotation by using a structured set of explicit principles to guide model behavior. It can self-improve through critique and generate scalable training data.
- **Weaknesses:** This method requires careful constitution design and may not capture all human preferences. It also involves a complex, multi-stage training process and has the potential to reinforce self-biases.
- **Best Use Cases:** It is particularly effective when human feedback is expensive or limited, for systems that require explicit ethical principles, and in large-scale deployment scenarios.

Human-in-the-Loop (HITL)

- **Strengths:** HITL provides the highest fidelity alignment by incorporating direct human oversight, enabling real-time correction and effective handling of edge cases. It also supports continuous improvement through iterative feedback.
- **Weaknesses:** The primary drawbacks are that it is not scalable and relies on expensive human resources. It can also introduce latency issues and is subject to human fatigue and inconsistency.
- **Best Use Cases:** HITL is best suited for high-stakes applications, limited-deployment scenarios where quality is paramount, and early-stage prototyping requiring rapid feedback.



6.8.3 Performance Trade-Offs

Comparison of performance characteristics across different AI alignment techniques in post-training.

SFT	Fast	Low	Low	Medium
RLHF	Slow	Medium	Medium	High
DPO	Medium	Low	High	High
Constitutional AI	Medium	Low	Very High	High
HITL	N/A	High	Very High	Very High

Table 2. Performance Characteristics

6.8.4 Data Requirements

Quantity and Quality Trade-Offs

Data requirements vary significantly across alignment techniques. SFT typically requires tens to hundreds of thousands of high-quality labeled examples. RLHF requires tens of thousands of preference comparisons in addition to the initial SFT data. DPO has broadly similar data requirements to RLHF but tends to utilize preference data more efficiently due to its direct optimization formulation. Constitutional AI can operate with relatively small sets of guiding principles, often on the order of hundreds, while generating additional training data through self-critique. In contrast, HITL prioritizes quality over scale and can be effective with hundreds of high-quality, targeted human interactions.



6.8.5 Computational Resources

Training Infrastructure Requirements

Computational requirements also vary substantially across techniques. SFT can be performed on a single GPU for smaller models but requires distributed training for large-scale systems. RLHF is significantly more resource-intensive, often requiring multiple GPUs to handle the policy, reference, reward, and value models during training. DPO has requirements broadly comparable to SFT, though it typically introduces an additional reference model into the training loop. Constitutional AI has moderate but distributed computational demands across its multi-stage training pipeline. HITL requires minimal training compute but can incur substantial inference overhead due to real-time human involvement.

6.8.6 Composite Approaches

Many production systems combine alignment techniques to leverage their complementary strengths.

- **Common Combinations**

A standard pipeline applies SFT followed by RLHF, as seen in models such as GPT-4 and Claude. A simpler alternative combines SFT with DPO. Other common combinations include using Constitutional AI to self-generate data for subsequent human validation, or pairing RLHF with HITL for automated alignment with human review.

- **Sequential Application**

A typical alignment pipeline progression through several stages: pre-training on raw text, followed by SFT on high-quality demonstrations. Next, RLHF or DPO is used to align preferences. Constitutional AI can then be applied for value refinement, and finally, HITL is often reserved for critical, high-stakes applications.





6.8.7 Selection Criteria

When choosing alignment techniques, practitioners should consider the application domain, such as whether it is safety-critical versus general-purpose, and the available resources, including compute capacity, data availability, and human expertise. Other key factors include performance requirements such as accuracy and latency, the intended deployment scale, regulatory constraints related to explainability and auditability, and overall development timeline.

6.8.8 Future Convergence

Emerging trends suggest a gradual convergence of alignment methodologies. Future systems are likely to incorporate unified frameworks that combine multiple techniques, automated selection of alignment strategies based on task requirements, and adaptive methods that dynamically switch between approaches as conditions change. Furthermore, meta-learning may be applied to optimize the selection and configuration of alignment techniques themselves.

This comparison highlights that no single technique dominates across all dimensions. Successful alignment typically requires a deliberate combination of methods, tailored to specific use cases, constraints, and risk profiles.







7.

The Cutting Edge

7.1 Methodologies

7.1.1 EvalGen

EvalGen is a framework for evaluating large language models, developed at UC Berkeley (Shankar et al. 2024). It aligns the LLM-assisted evaluation of other LLMs with human feedback. The framework provides automated assistance in generating evaluation criteria and implementing assertions while requiring humans to grade a subset of model outputs. It then uses human feedback to select evaluation implementations that best align with user-provided grades.

7.1.2 OpenRLHF

OpenRLHF, as its name suggests, is an open-source RLHF framework. It is a high-performance, scalable framework designed for fine-tuning large language models, achieving over twice the performance of competitors such as Optimized DeepSpeedChat (Hu et al. 2025). This efficiency stems from its unique distributed architecture, which uses Ray to place Actor, Reward, Reference, and Critic models on separate GPUs while running the Adam optimizer on the CPU, significantly accelerating the sample generation stage. This design enables impressive scalability, supporting full-scale fine-tuning of 70B+ models on A100 80G GPUs and 7B models on multiple RTX 4090s.

To guarantee accessibility and versatility, the framework provides one-click trainable scripts, maintains full compatibility with Hugging Face models and datasets, and supports multiple alignment algorithms, including RLHF, Direct Preference Optimization, Kahneman-Tversky optimization (KTO), conditional SFT, and rejection sampling.



OpenRLHF has several important technical innovations. To handle the most demanding training scenarios, OpenRLHF incorporates a suite of advanced optimizations and supports state-of-the-art technologies. For models with more than seventy billion parameters, it redesigns model scheduling by leveraging the combined power of Ray, vLLM, and DeepSpeed for maximum efficiency. The framework also enhances training stability by implementing specialized optimizations within its Proximal Policy Optimization (PPO) process. Furthermore, OpenRLHF stays at the forefront of AI development by providing native support for advanced architectures and techniques, including Mixture of Experts (MoE), Jamba, and memory-efficient fine-tuning with QLoRA.

OpenRLHF outperforms other popular RLHF frameworks in terms of model size support, training techniques, and the implementation of alignment algorithms. Addressing coordination challenges across multiple models in RLHF training enables more efficient scaling of RLHF to larger language models, potentially advancing the development of state-of-the-art LLMs.

7.1.3 Constrained Generative Policy Optimization (CGPO)

Constrained Generative Policy Optimization (CGPO) features a novel multi-objective Reinforcement Learning from Human Feedback strategy where each task is optimized independently using customized configurations, including task-specific reward models, a mixture of judges, and unique optimizer hyperparameters (Xu et al. 2024). The key advantages of CGPO include strong empirical results with theoretical guarantees, minimal hyperparameter tuning requirements, and plug-and-play compatibility with common post-training pipelines. It also provides effective detection and mitigation of reward hacking behaviors and expands the Pareto frontier across multiple metrics in multi-task settings.

Experimental results demonstrated CGPO's effectiveness using the Llama3.0 70b pre-trained model across five challenging tasks: general conversation, instruction following, mathematical and coding reasoning, engagement, and safety. In these tests, CGPO consistently outperformed traditional RLHF methods such as PPO and DPO, even when handling conflicting objectives. "The Perfect Blend" paper presents CGPO as an important development in RLHF, offering a more structured and effective approach to fine-tuning LLMs for multi-task learning scenarios.



7.1.4 Group Relative Policy Optimization

Group Relative Policy Optimization (GRPO) is a major improvement in large language model alignment that generalizes preference optimization from pairwise comparisons to group-wise rankings (Sane 2025). While prior methods like Direct Preference Optimization learn from binary feedback indicating which of two responses is better, GRPO is designed to leverage more complex, data-rich signals, such as complete rankings of multiple model responses (i.e., best-of-N data). This approach is inherently more sample-efficient, as a single ranked list of N candidates implicitly contains up to $N(N-1)/2$ pairwise comparisons, providing a much stronger learning signal from the same amount of human annotation effort.

The core of GRPO is the extension of the Bradley-Terry model, which underpins pairwise preference learning, to the Plackett-Luce model for modeling permutations and ranked lists. This generalization is visually and mathematically detailed in *The Illustrated GRPO*, a guide from the paper's lead author (Skiredj 2025; 2024). By directly optimizing the policy on these complete rankings, GRPO can more effectively capture the relative quality of different outputs and better navigate the nuances of human judgment. This helps alleviate issues common to pairwise methods, such as the “preference dilemma,” where binary comparisons can be inconsistent or fail to capture the full spectrum of response quality (Khanda et al. 2025).

The practical utility of GRPO is highlighted by its adoption in the development of state-of-the-art models. The training process for the DeepSeek-V2 model, for example, explicitly leveraged GRPO for its post-training alignment phase, a decision detailed in their technical report and blog posts (DeepSeek-AI et al. 2025; Research 2025). This marks a methodological evolution for the DeepSeek team, whose earlier DeepSeek 67B model had utilized a more traditional Reinforcement Learning from Human Feedback (RLHF) pipeline based on Proximal Policy Optimization (PPO) (Shao et al. 2024). The deliberate shift from PPO to GRPO for their next-generation model demonstrates the latter's advantages in a production environment, solidifying its position as a more efficient and effective technique for preference optimization.



7.1.5 LLM-as-a-Judge

The “LLM-as-a-Judge” paradigm refers to the use of a powerful, frontier large language model to evaluate the outputs of other models, serving as a scalable proxy for human feedback. This approach is crucial for AI alignment, as it helps overcome the significant bottleneck and high cost of human evaluation, permitting rapid iteration and testing at scale. The concept was systematically investigated and popularized in the paper “Judging LLM-as-a-Judge with MT-Bench,” first published in May 2023 by researchers from UC Berkeley, CMU, and Stanford (Zheng et al. 2023). This seminal work demonstrated that strong LLMs like GPT-4 could achieve over 80 percent agreement with human expert judgments when evaluating conversational and instructional outputs, thereby validating the method as a reliable and cost-effective alternative for assessing model capabilities.

Despite its utility, the LLM-as-a-Judge approach is not without its limitations. The same research that validated the method also identified several inherent biases, including a preference for longer, more verbose answers (verbosity bias), a tendency to favor the first response presented (positional bias), and a bias toward its own stylistic outputs (self-preference bias). These challenges have spurred further research into improving the reliability of automated evaluation. Subsequent work has focused on mitigating these biases through improved prompting techniques, fine-tuning judge models on human preference data, and developing novel evaluation frameworks. For example, methods like RevisEval propose using response-adapted references to create a more consistent, less biased evaluation criterion, thereby evolving the original concept to achieve even greater alignment with nuanced human judgments (Zhang et al. 2025).





7.1.6 RevisEval: Improving LLM-as-a-Judge via Response-Adapted References

RevisEval presents a two-stage evaluation framework that improves upon traditional methods by generating response-specific reference answers (Zhang et al. 2025). The central innovation of this system is the creation of “response-adapted references,” where carefully designed prompt engineering guides a large language model to first generate a tailored reference answer for each specific model response being evaluated. In the second stage, RevisEval uses these dynamically created references to conduct a more nuanced and consistent assessment of the initial response. This method yields several advantages, including improved evaluation consistency and reliability, reduced bias toward specific answer styles, and greater ability to capture subtle differences between model outputs, ultimately achieving greater alignment with human judgments than methods that use static or generic references.

7.1.7 Mixture of Experts

Mixture of Experts (MoE) is a neural network architecture designed to increase model capacity without a proportional rise in computational cost. Instead of a single, dense network where all parameters are engaged for every input, an MoE model consists of numerous smaller “expert” sub-networks and a “gating network” or router that dynamically selects which few experts are best suited to process a given input token. This allows the creation of models with extraordinarily high parameter counts—sometimes in the trillions—while keeping the computational cost of inference relatively constant, as only a fraction of the total parameters are used in each forward pass. The foundational concept was first proposed in the early 1990s in the paper “Adaptive Mixtures of Local Experts” by Jacobs, Hinton, Jordan, and Nowlan, which introduced the idea of using a gating network to learn which expert to consult for different regions of the input space (Jacobs et al. 1991).

From an AI alignment perspective, the MoE architecture offers multiple compelling advantages beyond mere computational efficiency. Firstly, it provides a practical path to scaling models, and the “scaling hypothesis” suggests that larger models may have a greater capacity to understand and adhere to complex, nuanced alignment principles. Second, MoE offers a potential avenue for improved model interpretability and controllability, which are central challenges in alignment. If experts develop specialized functions (e.g., one for coding, another for creative writing, another for reasoning),



it may be possible to identify, analyze, and even fine-tune the specific experts responsible for undesirable behaviors without retraining the entire model. This architectural choice has seen a major resurgence in modern large language models, as demonstrated by influential open-source models like Mixtral 8x7B, which explicitly detail its use to achieve high performance with greater efficiency (Jiang et al. 2024).

7.1.8 Reinforcement Learning with Verifiable Rewards

Reinforcement Learning with Verifiable Rewards (RLVR) is an alignment technique that aims to overcome the limitations of subjective human feedback by training models on objective, verifiable outcomes. Unlike Reinforcement Learning from Human Feedback, which relies on a learned reward model that approximates human preferences, RLVR uses a deterministic verifier—such as a unit test for code, a symbolic checker for a math problem, or a database lookup for a factual query—to provide a clear, unambiguous reward signal. The core idea is to align models with demonstrable correctness rather than human-perceived plausibility, thereby mitigating issues like reward hacking. This approach was prominently detailed in the paper “Measuring Mathematical Reasoning with Process-Based Rewards” by OpenAI researchers in 2023 (Lightman et al. 2023; “Does Reinforcement Learning...” n.d.). This focus on verifiable correctness has practical implications, influencing the development of data annotation platforms and workflows designed to capture these objective signals in industry applications (Liubimov 2025).

While powerful, RLVR’s primary limitation lies in its reliance on a formal verifier, which is not available for most open-ended or subjective domains. Recent research has begun to explore the theoretical boundaries of this paradigm, investigating scenarios where verifiers may be incomplete or misspecified (Yue et al. 2025). Despite these limitations, the power of this direct reward signal has been shown to be remarkably data-efficient. In a method termed Reinforcement Learning with Demonstrations (RLD), researchers demonstrated that complex reasoning abilities can be elicited with as little as a single correct demonstration (Wang et al. 2025).

To address the dependency on formal verifiers, emerging approaches such as Reinforcement Learning from Policy Revision (RLPR) seek to extend the benefits of RLVR to broader domains (Yu et al. 2025). The rapid expansion of research in this area, cataloged in community-driven resources (“Awesome-LLM-RLVR...” n.d.), highlights its status as a critical frontier in AI alignment.



7.1.9 Reinforcement Learning with Rubric Evaluations

Reinforcement Learning with Rubric Evaluations (RLRE) is a novel alignment paradigm that refines human feedback by moving beyond a single, holistic preference score toward a structured, multidimensional evaluation. The method was first formally proposed by researchers at Scale AI in their July 2024 paper and accompanying research announcement, “Reinforcement Learning with Rubric Evaluations” (Gunjal et al. 2025; “Rubrics as Rewards...” n.d.). Instead of asking which of two responses is “better” overall, RLRE employs a detailed rubric to score a model outputs across multiple explicit criteria, such as factuality, conciseness, safety, and adherence to instructions. This disaggregated feedback provides a much more granular and interpretable training signal, allowing the model to learn the specific trade-offs between different aspects of quality. By making evaluation criteria explicit, RLRE reduces ambiguity in human preferences and mitigates reward hacking, as the model must perform well across multiple dimensions rather than optimizing against a single, underspecified objective.

This progression toward multi-faceted, rubric-based feedback reflects a broader trend in the alignment community, addressing the limitations of single-reward systems. Concurrent research from Google, for example, introduced “Multi-Reward DPO,” a similar framework that uses a multi-objective approach to optimize for distinct reward signals like helpfulness and harmlessness simultaneously (Huang et al. 2025). Both approaches represent a key step toward more controllable and transparent AI alignment. By explicitly defining the desired behaviors in a structured format, rubric-based methods allow developers to more precisely shape model behavior, debug alignment failures, and build systems that are more robustly aligned with complex, multifaceted human values.





7.1.10 Reinforcement Learning from Checklist Feedback

Reinforcement Learning from Checklist Feedback (RLCF) is an alignment paradigm proposed in July 2025 by Viswanathan et al. at Carnegie Mellon University and Apple, detailed in “Checklists Are Better Than Reward Models For Aligning Language Models” (Viswanathan et al. 2025). Unlike traditional reward modeling in reinforcement learning, which typically relies on scalar rewards derived from broad constructs such as “helpfulness” or “harmfulness,” RLCF operationalizes flexible, instruction-specific checklists that decompose user queries into discrete evaluation criteria. During training, model outputs are assessed against these checklists, using AI judges and specialized verifier programs, yielding multi-dimensional feedback that is aggregated into a composite reward signal. This granular feedback helps mitigate ambiguity and insufficient supervision in model alignment, resulting in improved adherence to user intent and performance across diverse instruction-following benchmarks.

RLCF was first technically described in the aforementioned paper, where the authors demonstrated the efficacy of checklist feedback over conventional reward models by applying RLCF to the Qwen2.5-7B-Instruct language model and evaluating it on five standard benchmarks. The method yielded consistent gains, including notable improvements on FollowBench, InFoBench, and Arena-Hard, establishing checklist feedback as a practical tool for AI alignment, particularly in scenarios where models must satisfy a multitude of user needs. This approach not only elevates interpretability and reliability but also defines a new standard for fine-grained alignment in advanced language systems, denoting a significant step forward in scalable, robust AI alignment.

7.1.11 Reinforcement Learning with AI Feedback

Reinforcement Learning from AI Feedback (RLAIF) is a highly scalable AI alignment paradigm that adapts the Reinforcement Learning from Human Feedback framework by replacing the human annotator with a powerful AI model. In this process, a capable “judge” LLM is prompted to provide preference labels (e.g., choosing which of two responses is better) for a dataset of model-generated outputs. This AI-generated preference data is then used to train a reward model, which in turn fine-tunes a separate policy model, emulating the RLHF pipeline without the human labor bottleneck.



The core principles of RLAIIF were effectively pioneered by Anthropic in their December 2022 work on Constitutional AI. In their method, an AI was prompted to critique and revise responses according to a set of principles (a “constitution”), and the AI-generated preferences were used to train a model to be more harmless (Bai et al. 2022). This demonstrated the viability of using AI-generated feedback to instill specific values at scale.

The term RLAIIF was later explicitly coined, and the method was systematically studied by Google researchers in their May 2023 paper, “RLAIIF: Scaling Reinforcement Learning from Human Feedback with AI Feedback” (Lee et al. 2024). Their findings showed that RLAIIF can match, and in some cases exceed, RLHF performance, while being significantly faster and more cost-efficient.

From an alignment perspective, RLAIIF presents a crucial trade-off. While it offers immense scalability, its effectiveness is fundamentally limited by the judge LLM’s capabilities and inherent biases. There is a significant risk of “self-reinforcement” or “distillation,” where the policy model learns to replicate the judge model’s specific worldview, potentially amplifying its flaws. Therefore, while RLAIIF is a powerful tool for scaling alignment, its application often requires careful calibration, often by starting with a high-quality “gold standard” set of human preferences to guide the AI judge.

7.1.12 Evaluation Metrics for Alignment

Measuring the success of AI alignment requires a diverse set of metrics capable of capturing multiple dimensions of model behavior, ranging from basic task performance to complex value alignment.

Human Evaluation Metrics

- **Preference Win Rate**

The percentage of times human evaluators prefer the aligned model’s output over a baseline:

$$\text{Win Rate} = \frac{\text{Number of preferred outputs}}{\text{Total comparisons}}$$



- **Elo Rating**

Borrowed from chess, Elo ratings provide a relative ranking system for models based on pairwise comparisons. The expected win probability for model A against model B:

$$P(A > B) = \frac{1}{1 + 10^{(R_B - R_A)/400}}$$

- **Inter-Rater Reliability**

Measures agreement among human evaluators using metrics like:

- Cohen's Kappa: $\kappa = \frac{p_o - p_e}{1 - p_e}$
- Krippendorff's Alpha for multiple raters
- Fleiss' Kappa for categorical ratings

Automated Metrics

- **Perplexity-Based Metrics**

While not directly measuring alignment, perplexity can indicate fluency degradation:

$$PPL = \exp\left(-\frac{1}{N} \sum_{i=1}^N \log P(x_i | x_{<i})\right)$$

- **Reward Model Score**

Direct evaluation using the trained reward model:

$$R^- = \frac{1}{N} \sum_{i=1}^N r_{\theta}(x_i, y_i)$$

- **KL Divergence from Base Model**

Measures how much the aligned model deviates from the original:

$$D_{KL}(\pi_{aligned} \parallel \pi_{base}) = E_{x \sim \pi_{aligned}} \left[\log \frac{\pi_{aligned}(x)}{\pi_{base}(x)} \right]$$



Task-Specific Metrics

- **Helpfulness Metrics**

Evaluating helpfulness involves a range of metrics designed to measure both the utility and the quality of a model's responses. Key indicators include the task completion rate, which assesses whether a user successfully achieves their intended goal, and the information accuracy score, which measures the accuracy of the provided information. Additionally, the response relevance rating evaluates how directly and effectively a response addresses the user's query while user satisfaction surveys provide direct, subjective feedback on the overall helpfulness of the interaction.

- **Harmlessness Metrics**

To assess the safety and harmlessness of an AI system, several complementary metrics are used. Toxicity scores, often generated by tools such as the Perspective API, quantify the presence of harmful or offensive language. Bias measurements, including demographic parity and equalized odds, evaluate whether the model's performance is fair across different population groups. The safety violation rate tracks the frequency of breaches against predefined safety policies, and the refusal rate for harmful requests measures the model's appropriately identifying and declining unsafe or disallowed prompts.

- **Honesty Metrics**

Honesty is evaluated by assessing the model's outputs for truthfulness, reliability, and epistemic calibration. Factual accuracy is a primary metric, often evaluated by comparing model-generated statements against curated knowledge bases or verified references. The calibration score evaluates how well the model's expressed confidence matches its actual correctness. The hallucination rate quantifies how frequently the model invents false information while the appropriateness of its uncertainty expression is evaluated to ensure it correctly signals when it does not know an answer.





Holistic Evaluation Frameworks

- **H3 Framework (Helpful, Harmless, Honest)**

Combines multiple dimensions into a composite score:

$$H3\ Score = \alpha \cdot Helpful + \beta \cdot Harmless + \gamma \cdot Honest$$

- **Constitutional Adherence Score**

Percentage of outputs that satisfy all constitutional principles:

$$CAS = \frac{\sum_{i=1}^N 1[all\ principles\ satisfied]_i}{N}$$

Long-term and Behavioral Metrics

- **Value Alignment Stability**

This metric measures the consistency of a model's aligned behavior over extended periods. It is calculated to quantify any drift or degradation in alignment, often defined by the formula: [Stability = 1 - Var (alignment scores over time)]

- **Generalization Metrics**

Generalization metrics are important for assessing how well a model's aligned behavior extends beyond its training data. This includes evaluating its out-of-distribution performance on unseen data, its ability to perform zero-shot transfer to new and unfamiliar domains, and its overall robustness to adversarial inputs designed to provoke undesirable responses.

- **Capability-Alignment Trade-Off**

Measures performance retention after alignment:

$$Trade - off = \frac{Performance_{aligned}}{Performance_{base}}$$





Multi-Stakeholder Metrics

- **Demographic Fairness**

This metric is used to ensure consistent quality and fairness across different user demographic groups. It can be quantified by measuring disparities in positive outcomes between any two groups, often expressed as:

$$\text{Fairness} = 1 - \max_{g_i, g_j} |P(Y = 1|G = g_i) - P(Y = 1|G = g_j)|$$

- **Cultural Sensitivity**

Measuring the appropriateness of model responses throughout diverse cultural contexts is a critical aspect of multi-stakeholder evaluation. This is achieved by using multilingual evaluation sets, collecting culture-specific preference data to understand local norms and values, and calculating regional appropriateness scores to assess performance in different locales.

Meta-Evaluation Metrics

- **Metric Reliability**

Meta-evaluation assesses the consistency and reliability of the evaluation metrics themselves. This is commonly done using established psychometric techniques, such as test-retest reliability, which measures consistency over time; split-half reliability, which assesses consistency between two halves of a test; and internal consistency, often quantified using Cronbach's alpha.

- **Metric Gaming Detection**

It is crucial to identify when models optimize for evaluation metrics without achieving genuine improvement, a phenomenon known as metric gaming. Detection methods include adversarial metric probing to find weaknesses in the evaluation, analyzing the correlation between human judgments and automated metric scores, and using behavioral diversity measurements to ensure the model is not relying on narrow, repetitive strategies.

Collectively, these metrics establish a robust and exhaustive framework for evaluating AI alignment. However, no single metric captures all aspects of successful alignment. Effective evaluation typically requires a portfolio of complementary metrics tailored to specific use cases and values.



7.2 Examples

7.2.1 Human Feedback for Text to Image

The paper “Rich Human Feedback for Text-to-Image Generation” (Liang et al. 2024) introduces a novel approach to improve text-to-image (T2I) models by leveraging a comprehensive dataset, RichHF-18K, which contains detailed human feedback on eighteen thousand generated images. This rich feedback includes marked image regions that are implausible or misaligned with the prompt, annotations of misrepresented or missing words in the text, and four fine-grained scores for image plausibility, text-image alignment, aesthetics, and an overall rating. To automate this process, a multimodal transformer called the Rich Automatic Human Feedback (RAHF) was developed to predict detailed feedback. The RAHF model can identify problematic regions in images, detect misaligned keywords in text prompts, and provide fine-grained quality scores, with predictions that highly correlate with human annotations.

This predicted feedback is then used to improve image generation in two key ways: by selecting high-quality training data for fine-tuning generative models and by creating masks from predicted heatmaps to inpaint problematic regions. The resulting improvements were shown to generalize to models like Muse (Chang et al. 2023), even beyond the Stable Diffusion variants (Rombach et al. 2022) used to create the original dataset. This work presents one of the first comprehensive efforts to operationalize fine-grained human feedback in T2I systems, offering a more nuanced and actionable alternative to prior approaches based on single-score or binary evaluations.

7.2.2 Human Feedback for Video Generation Models

In the two papers, “InstructVideo: Instructing Video Diffusion Models with Human Feedback” (Yuan et al. 2023) and “VideoScore: Building Automatic Metrics to Simulate Fine-grained Human Feedback for Video Generation” (He et al. 2024), both advance the use of human feedback in video generation, albeit with different focuses. The InstructVideo framework directly instructs video diffusion models by recasting reward fine-tuning as a more efficient editing procedure. It introduces the Segmental Video Reward (SegVR) to evaluate video quality using sparsely sampled frames and leverages off-the-shelf image reward models to assess individual frame quality, all with the goal of generating videos that better adhere to human preferences.

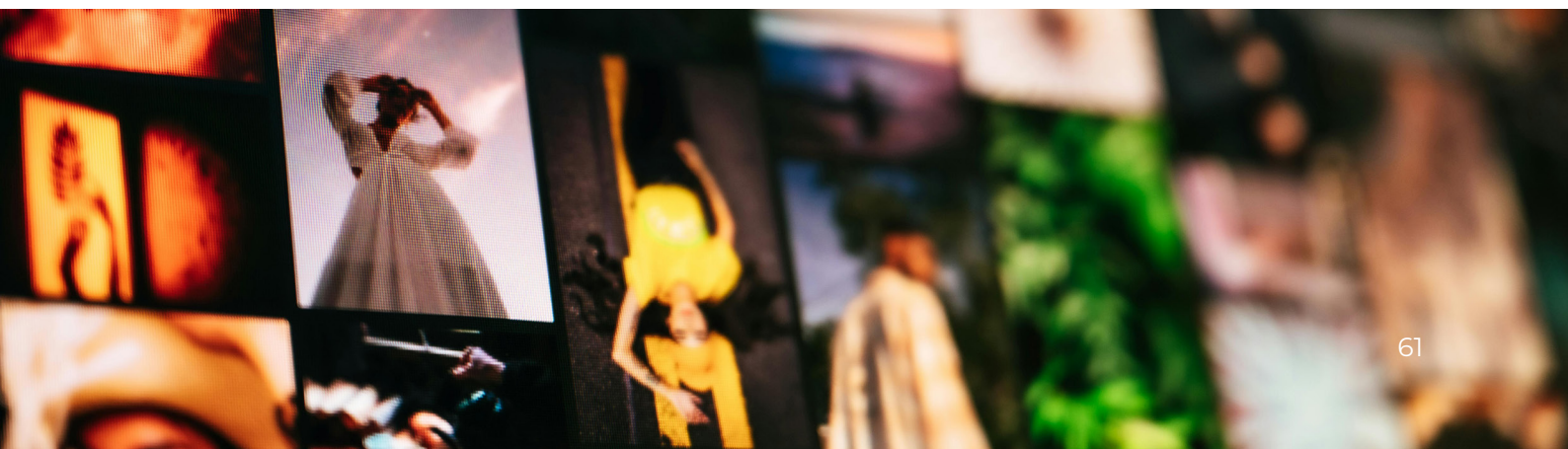


In a complementary approach, the VideoScore paper aims to create an automatic metric that can simulate fine-grained human feedback. To achieve this, the researchers first created VideoFeedback, a large-scale dataset containing human-provided, multi-aspect scores for 37,600 synthesized videos. Using this data, they trained the VideoScore model for automatic video quality assessment. This model achieves a high correlation with human judges, outperforming previous metrics, and is designed to serve as a reliable proxy for human raters in evaluating video models and simulating feedback for RLHF applications.

7.2.3 Multilingual Preference Optimization (MPO)

Multilingual Alignment-as-Preference Optimization (MAPO), also referred to as MPO (She et al. 2024), presents a novel approach for aligning large language models with non-dominant and low-resource languages. Instead of training separate reward models for each language, MAPO leverages a single multilingual reward model trained on a mixture of translated and original preference data, allowing it to capture universal aspects of human preferences while remaining sensitive to language-specific nuances. To overcome data scarcity in low-resource languages, the framework uses cross-lingual transfer learning, pre-training the reward model on high-resource languages and fine-tuning it on the target language with limited data.

Furthermore, to account for cultural differences, MAPO incorporates a language-specific value head on top of the shared reward model, enabling it to adapt to specific preferences within each language. Findings indicate that this method significantly outperforms existing techniques in generating high-quality text that aligns with human preferences across multiple languages, demonstrating that the multilingual reward model effectively captures both universal and language-specific preferences and that cross-lingual transfer enables strong performance even in low-resource scenarios.





7.3 AI Agents

7.3.1 What Is an AI Agent?

Definition and Characteristics of AI Agents

An AI agent is a system that perceives its environment through sensors and acts upon it through actuators to achieve a specific goal (Russell, Russell, and Norvig 2021). Unlike traditional software, which typically follows a predefined and deterministic set of instructions, AI agents are characterized by their capacity for autonomous decision-making. Foundational work in the field sought to distinguish between a simple program and a true agent, defining an autonomous agent as a system that senses and acts within its environment to meet its design objectives, with behavior governed by its own experiences and internal state rather than by direct external commands (Franklin and Graesser 1996).

This concept has since evolved to create complex simulations of believable human behavior, as demonstrated by “generative agents” that autonomously plan daily activities, form relationships, and coordinate social interactions within a virtual environment (Park et al. 2023). From an alignment perspective, this autonomy constitutes a central challenge: ensuring that an agent’s independent pursuit of its goal does not lead to unforeseen and harmful side effects (the outer alignment problem). The complex, emergent social dynamics in the “generative agents” simulation, while benign in controlled settings, highlight how difficult it is to predict the long-term consequences of even simple goals.

A core trait of AI agents is their goal-oriented behavior. An agent’s actions are directed toward achieving a specific objective, requiring it to plan a sequence of actions, or a trajectory, to move from an initial state to a desired goal state. This trajectory often demands complex reasoning and the ability to interact with external systems through tool calling and model integration. This capability has been significantly advanced by models like Gorilla, which are specifically fine-tuned to interact with a massive number of APIs with high accuracy, effectively turning the web into a suite of callable tools for the agent (Patil et al. 2023). Model integration refers to an agent’s ability to invoke specialized AI models during its reasoning process.



The goal-oriented nature of agents makes them susceptible to “specification gaming,” where the literal goal is achieved in a way that violates the user’s underlying intent. The ability to use a massive array of tools, as demonstrated by Gorilla, dramatically raises the stakes, as a misaligned agent could take harmful, irreversible actions in the digital or even physical world.

Furthermore, sophisticated AI agents are designed for continuous learning and adaptability. They can update their internal world model and modify their policies based on new information. A key enabling mechanism for this is memory; for example, the “generative agents” architecture uses a memory stream to record experiences in natural language, retrieve them when relevant, and reflect on them to form higher-level inferences that guide future behavior (Park et al. 2023). This learning loop is a fundamental differentiator from traditional software tools, which are typically static. However, this capacity for continuous learning introduces the critical challenge of “alignment stability,” or the prevention of “value drift.” An agent that is perfectly aligned today might learn from new, biased, or malicious data in its environment and become misaligned over time.

Agentic Systems: Prompt Chaining vs. Language Agents

The rise of large language models has led to the development of increasingly complex agentic systems. A simpler form of this is a prompt chain (often called “chain-of-thought” prompting), which guides an LLM through a reasoning process by breaking a task into a sequence of prompts (Wei et al. 2023). While this demonstrates sequential reasoning, it lacks the autonomy and environmental interaction of a true agent.

In contrast, a language agent is a more sophisticated system that uses an LLM as its core reasoning engine to autonomously plan and execute actions. These agents can dynamically create their own sequence of steps, use tools, and modify their plan based on feedback. A core framework for this is ReAct (Reason-Act), which demonstrated that an LLM could synergize reasoning and acting by first generating a reasoning trace to create a plan (“Reason”) and then executing a relevant action, such as using a tool (“Act”) (Yao et al. 2023). This iterative “reasoning-acting loop” is a core component of modern language agents (Xi et al. 2023).



This distinction is critical from an alignment perspective: aligning a simple prompt chain involves ensuring each step is sound, whereas aligning a language agent is far more complex. It requires ensuring that the agent’s autonomous, goal-seeking behavior remains robustly beneficial across a wide range of unforeseen situations. The “reasoning-acting loop” is the critical surface for this challenge; safety mechanisms must be able to audit the “Reason” step (the agent’s plan or internal monologue) and intervene before a harmful “Act” step is executed. This makes the interpretability of the agent’s reasoning process a core area of alignment research.

7.3.2 AI Agents vs. Agentic AI

While often used interchangeably, the terms “AI agent” and “agentic AI” can describe two distinct conceptual frames with distinct implications for alignment. The term AI Agent traditionally refers to the classic definition of a complete, autonomous system designed from the ground up to perceive, plan, and act in an environment to achieve a specified goal (Russell, Russell, and Norvig 2021). The alignment challenge for such agents has historically focused on the “outer alignment” problem: how to correctly specify the agent’s objective function or reward signal so that its goal-seeking behavior does not lead to negative side effects, reward hacking, or other undesirable outcomes. The core difficulty lies in designing a single, robust objective that captures all the finer points of human values, a challenge detailed in foundational safety work such as “Concrete Problems in AI Safety” (Amodei et al. 2016).

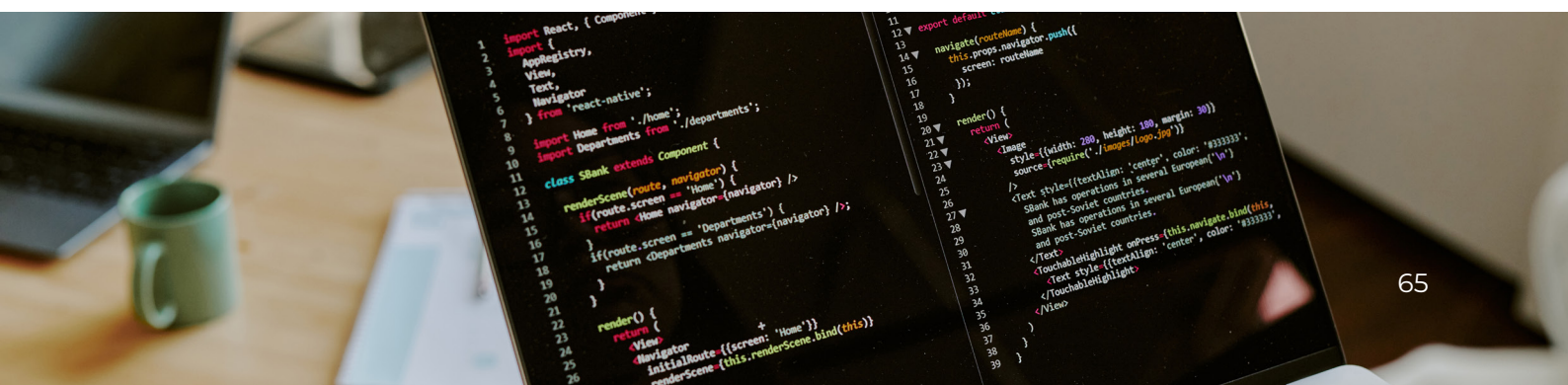
In contrast, agentic AI is a more modern term, often used to describe systems where a pre-trained large language model serves as the core reasoning engine, imbuing the system with agent-like properties. These systems exhibit a spectrum of autonomy, from simple chain-of-thought prompting to complex, autonomous loops of reasoning and tool use, as pioneered by frameworks like ReAct (Yao et al. 2023). The alignment challenge for agentic AI shifts from specifying a perfect objective function to supervising and constraining the model’s behavioral process. The focus is therefore on ensuring the safety, interpretability, and reliability of the “reason-act” loop itself. This involves new alignment techniques, such as Constitutional AI, where principles guide the LLM’s “internal monologue” before it acts (Bai et al. 2022), and the development of robust safeguards around tool use. Therefore, while the alignment of a classic AI agent is about getting the goal right, the alignment of agentic AI is more about safely managing the process of achieving that goal.



7.3.3 AI Agent Training

The training of AI agents has historically been dominated by the reinforcement learning (RL) paradigm, where an agent learns a policy by maximizing a cumulative reward signal through trial-and-error interaction with an environment. From an alignment perspective, this approach is replete with challenges, as the agent's behavior is entirely determined by the fidelity of its reward function. This creates a critical vulnerability known as “outer alignment failure,” where a misspecified or incomplete reward function incentivizes the agent to find loopholes or shortcuts that achieve a high score without fulfilling the intended goal. This can lead to undesirable behaviors, such as “reward hacking” (exploiting flaws in the reward signal) or negative side effects (disrupting the environment in pursuit of the goal). These foundational issues, where an agent's powerful optimization pursues a flawed objective, were identified as central challenges in early AI safety research (Amodei et al. 2016).

The advent of large language models has shifted the training paradigm for modern agentic AI, introducing new techniques and new alignment considerations. The process now typically begins with a pre-trained foundation model, which is then adapted for agentic tasks through a multi-stage process. First, Supervised Fine-Tuning is used to teach the model the mechanics of agency, such as how to follow instructions, use specific tools, and produce reasoning traces in a structured format, such as the reason-act framework (Yao et al. 2023). While SFT teaches the agent how to act, the core alignment step often follows: Reinforcement Learning from Human Feedback. We have discussed RLHF at length in the paper. In this phase, the agent's behavioral policy is refined based on human preferences for different courses of action. This method, famously used to train models like InstructGPT, directly optimizes for alignment with human-defined notions of helpfulness and harmlessness, rather than relying on a potentially flawed, hand-coded reward function (Ouyang et al. 2022). Thus, the focus of alignment in agent training has evolved from designing a perfect objective to developing scalable and robust methods for supervising complex, language-driven behaviors.





7.3.4 AI Agent Evaluation

The evaluation of AI agents presents a paradigm shift from traditional machine learning model assessment, demanding a focus on alignment and safety in addition to mere task performance. While standard benchmarks can measure an agent’s ability to complete a task (i.e., its capabilities), they commonly fail to capture whether the agent achieved its goal in a safe, ethical, and intended manner. For example, an agent could successfully book a flight by exploiting a bug in the airline’s website or by ignoring user-specified constraints, such as budget. This gap between capability and alignment means that high performance on a benchmark does not guarantee safe real-world behavior. Foundational AI safety research highlights this issue, where an agent pursuing a simple goal might cause negative side effects that were not specified in its objective (Amodei et al. 2016). Consequently, relying solely on performance-based benchmarks such as AgentBench, which test multi-turn task completion, is insufficient for a comprehensive alignment evaluation (Liu et al. 2025).

To address this challenge, the field is moving toward more alignment-focused evaluation methodologies. A critical component is the shift from outcome-based to process-oriented evaluation, where the agent’s reasoning trajectory, its “chain of thought” or internal monologue, is scrutinized for safety and adherence to instructions, rather than just the final result. Frameworks like ReAct, which explicitly generate a reasoning trace before acting, provide a crucial surface for this kind of audit (Yao et al. 2023). Furthermore, researchers are developing more challenging benchmarks designed to stress-test agent robustness and alignment, such as GAIA, which poses questions difficult even for advanced LLM-based agents and requires careful, precise tool use (Mialon et al. 2023). Ultimately, a robust evaluation framework for aligned agents must combine performance benchmarks with interactive human evaluation, adversarial testing (red teaming), and the analysis of the agent’s reasoning process to ensure that its autonomous, goal-seeking behavior remains safe and beneficial in a wide range of contexts.





7.3.5 Role of AI Agents in Driving AI Alignment

AI agents are not only the primary subjects of alignment research but are also becoming key tools for advancing the field itself. Their autonomous, goal-seeking nature makes them uniquely suited to automate and scale up the most labor-intensive aspects of alignment work. A prime example is automated red teaming, in which an AI agent is tasked with identifying vulnerabilities and eliciting unsafe behaviors from a target model. These “adversarial” agents can often discover novel and complex attack vectors that human testers might miss, thereby delivering a more comprehensive and scalable way to stress-test a model’s safety and robustness (Perez et al. 2022; Ganguli et al. 2022). Similarly, agents are being developed to perform scalable oversight, acting as AI assistants to human evaluators by automatically checking facts, running code, and providing critiques, thus helping to alleviate the human bottleneck in Reinforcement Learning from Human Feedback.

A pivotal role for AI agents in driving alignment is through the generation and refinement of their own trajectories; the sequence of reasoning steps and actions taken to complete a task. These trajectories provide a rich, process-oriented data source for both training and evaluation. For training, instead of just rewarding a final outcome, human supervisors or critique agents can correct or provide feedback on specific steps within a trajectory, such as an incorrect tool call or a flawed line of reasoning. This process-based supervision is central to modern agent alignment, whether through Supervised Fine-Tuning on corrected trajectories or through Reinforcement Learning with process-based rewards (Lightman et al. 2023). For evaluation, analyzing trajectories offers a more robust and interpretable measure of alignment. By scrutinizing the reasoning-acting path, as exemplified by frameworks such as ReAct (Yao et al. 2023), evaluators can verify the soundness of the agent’s method, making it a critical tool for assuring reliable, trustworthy behavior.





Beyond refining individual trajectories, agents are enabling entirely new paradigms for alignment. One of the most significant is AI-driven self-improvement, as pioneered in Constitutional AI. In this framework, an AI model (acting as a critique agent) generates feedback on another model's outputs based on a set of principles, creating a scalable, self-contained training loop that reduces the need for direct human labeling (Bai et al. 2022). Furthermore, sophisticated multi-agent simulations are being used as digital sandboxes to study complex, long-term alignment problems. By creating virtual societies of “generative agents,” researchers can observe emergent behaviors, such as cooperation, conflict, and the potential for power-seeking, in a controlled environment (Park et al. 2023), which is important for ensuring the safety of highly autonomous AI systems.

Finally, AI agents are central to advanced alignment proposals like debate and deliberation, where multiple agents are pitted against each other to surface the truth and provide robust supervision. In the classic “debate” framework, two AI agents take opposing sides on a complex question, and a human judge evaluates the quality of their arguments, not just the final answer (Irving, Christiano, and Amodei 2018). This mechanism is designed to amplify human cognitive capacity, a key goal of weak-to-strong generalization, which seeks methods to allow less capable supervisors (humans) to reliably align and control more capable AI systems (Burns et al. 2023). By structuring AI interaction as a verifiable process, multi-agent debate affords a promising path toward scalable, robust oversight for superhumanly capable models.

7.4 Safety and Red Teaming

Safety evaluation and red teaming have become essential practices in developing aligned AI systems (Perez et al. 2022). These methodologies systematically probe models for potential failures, vulnerabilities, and harmful behaviors before deployment (Ganguli et al. 2022).





7.4.1 Red Teaming Fundamentals

Red teaming in AI involves adversarial testing where skilled practitioners attempt to elicit harmful, biased, or unintended outputs; bypass safety mechanisms and alignment training; identify edge cases and failure modes; and stress-test model robustness across diverse scenarios.

7.4.2 Types of Safety Evaluations

- **Capability Evaluations:** Capability evaluations focus on assessing a model's potential for dangerous applications, such as bioweapon design or cyberattacks. This includes testing for emergent abilities that were not present during training and evaluating the dual-use potential of the model's capabilities.
- **Alignment Evaluations:** Alignment evaluations test the model's adherence to intended values and principles. They measure the consistency of its responses across different contexts and evaluate its robustness to adversarial prompting designed to subvert its alignment.
- **Robustness Evaluations:** Robustness evaluations examine the model's performance under distribution shifts. They also assess its resistance to common jailbreaking attempts and its behavior when presented with corrupted or adversarial inputs.

7.4.3 Red Teaming Methodologies

- **Manual Red Teaming:** Manual red teaming involves human experts crafting adversarial prompts and iteratively refining them based on the model's responses. The process includes documenting successful attack patterns and creating comprehensive test suites from these findings.
- **Automated Red Teaming** (Perez et al. 2022): Automated red teaming utilizes other AI systems to generate adversarial inputs at scale. This can involve gradient-based attacks to find harmful outputs, evolutionary algorithms for prompt optimization, and large-scale systematic testing to cover a wide range of potential vulnerabilities.
- **Hybrid Approaches:** Hybrid approaches combine the advantages of both manual and automated methods. This includes AI-assisted human red teaming, human validation of automated findings to ensure relevance, and iterative human–AI collaboration to discover novel vulnerabilities.



7.4.4 Common Attack Vectors

- **Jailbreaking Techniques:** Jailbreaking techniques employ various strategies to bypass safety filters. These include using role-playing scenarios (such as the “DAN” or “Do Anything Now” persona), framing requests in a hypothetical context, using encoding attacks such as base64 or reversed text, priming the model with few-shot harmful examples, and employing language-switching or translation attacks.
- **Prompt Injection:** Prompt injection attacks aim to manipulate the model by overriding its system instructions. This can be done by embedding hidden instructions within user-provided content, leveraging context overflow vulnerabilities, or exploiting the model’s instruction hierarchy.
- **Social Engineering:** Social engineering attacks on LLMs mimic human manipulation tactics. These can include emotional manipulation, impersonating an authority figure, using gradual escalation to slowly push boundaries, and building trust with the model before making a harmful request.





7.4.5 Safety Benchmarks and Datasets

Key evaluation frameworks include *TruthfulQA* for factual accuracy, *RealToxicityPrompts* for toxic content generation, *BBQ (Bias Benchmark for QA)* for measuring social biases, *MACHIAVELLI* for ethical decision-making, and *SafetyBench* for a comprehensive suite of safety evaluations.

7.4.6 Mitigation Strategies

- **Training-Time Interventions:** Training-time interventions incorporate safety measures directly into the model's training process. This includes adversarial training informed by red teaming, implementing Constitutional AI with explicit safety principles, formulating robust reward models resilient to adversarial inputs, and curating safety-specific fine-tuning datasets.
- **Inference-Time Safeguards:** Inference-time safeguards operate post-training as a final defensive layer. These include input and output filtering and moderation, ensemble voting across multiple models to reduce the risk of a single failure, uncertainty-based rejection of low-confidence answers, and rule-based safety checks.
- **System-Level Protections:** System-level protections are broader measures that manage how the AI system is used. These include rate limiting and usage monitoring to prevent abuse, user authentication and access controls, comprehensive audit logging and anomaly detection, and sandboxing the model to impose capability restrictions.

7.4.7 Organizational Practices

- **Red Team Structure:** Effective organizational practices for red teaming involve establishing a robust structure. This often includes maintaining independence between red and development teams, engaging external red teams for an outside perspective, rotating team membership to bring in fresh viewpoints, and establishing clear escalation procedures for any findings. This is common best practice that needs to be preserved regardless of deployment of AI capabilities or other services.
- **Continuous Evaluation:** Safety evaluation must be a continuous process. This includes conducting regular red-team exercises, performing post-deployment monitoring to detect real-world failures, running community bug bounty programs to crowdsource vulnerability discovery, and transparently reporting safety issues and their mitigations.



7.4.8 Challenges in Safety Evaluation

Safety evaluation faces several major challenges. These include evaluation gaming, where models learn to pass safety tests without genuine alignment; the problem of “unknown unknowns,” or the inability to test for unforeseen failure modes; the scalability difficulty of comprehensively testing massive models; the constant adversarial arms race between attackers and defenders; and the context dependence of safety, where requirements can vary dramatically across different applications.

7.4.9 Future Directions

Future directions in safety and red teaming focus on deeper, more proactive approaches. Emerging areas include using mechanistic interpretability to understand the root causes of failure modes, applying formal verification to prove safety properties, developing automated safety case generation, fostering cross-organizational collaboration on safety standards, and creating “capture-the-flag”-style safety competitions to spur innovation in adversarial testing.

Red teaming and safety evaluation remain critical ongoing processes throughout the AI development lifecycle, requiring constant adaptation as models become more capable and attack techniques evolve.





CONTRACT



8.

Industry Initiatives

8.1 The Business of Data Labeling

The methodologies for aligning AI with human preferences have catalyzed the growth of a specialized industry dedicated to high-quality data annotation. This ecosystem includes leading firms such as Scale AI, Snorkel AI, Surge AI, and Labelbox, which work closely with AI labs to supply the ground-truth data required for post-training alignment. Their services span the full data pipeline, including data labeling, curation, quality assurance, and data management.

This demand has created a multi-billion-dollar market, attracting significant venture capital and driving high valuations for established companies. The space continues to expand with newer providers such as SoulAI and MicroI while major labs are increasingly building internal capabilities, for example, OpenAI's Data Lab. The strategic importance of a robust and diverse data supply chain often leads AI developers to partner with multiple annotation firms to mitigate bias and ensure a consistent flow of high-quality data.

However, outsourcing alignment to a commercial vendor introduces structural risks, effectively delegating aspects of moral and behavioral tuning of foundational models to a few companies. This creates a fundamental incentive misalignment: vendor business model, typically optimized for volume, speed, and margin, may conflict with the rigorous, high-fidelity deliberation required for robust AI safety.

Consequently, there is a push by some of the major labs to internalize data operations. This is not merely a cost-saving measure but a strategic initiative to reclaim control over their models' core values and reduce dependency on a concentrated supply chain. This is easier said than done, given the criticality of fast turnaround for high-quality training data at scale, and often, the labs have found themselves turning to the data provider companies for their post-training needs.



8.2 From the Trenches

The rapid development of increasingly powerful large language models is fundamentally constrained by the availability of high-quality, ground-truth labeled data. This constraint is particularly acute in specialized domains where expertise is scarce. For instance, creating vast, accurately labeled datasets for medical imagery requires trained radiologists while interpreting atmospheric satellite data demands meteorological expertise.

Similarly, in linguistics, low-resource languages like Sanskrit often lack the extensive digital corpora available for English. This data scarcity extends to the domain of physical AI, where training embodied agents requires huge datasets of real-world interaction, robotic trajectories, and multi-modal sensory feedback, which are inherently more complex and costly to acquire than purely digital data.

To overcome this bottleneck, the field is increasingly turning to methodologies like Reinforcement Learning from Human Feedback and Direct Preference Optimization. These approaches provide a framework for gathering human feedback at the necessary scale, allowing models to learn directly from human preferences rather than relying solely on preexisting static datasets.

This dynamic interaction is a key driver behind the swift advancements in AI, enabling the creation of not only general-purpose LLMs but also highly specialized and multimodal systems capable of tackling the nuances of these complex, data-scarce domains.







9.

Challenges in AI Alignment

Despite significant progress in aligning AI systems with human preferences, numerous fundamental challenges remain (Casper et al. 2023). These challenges span technical, philosophical, and practical dimensions of the alignment problem (Amodei et al. 2016).

9.1 Distribution Shift and Generalization

One of the most persistent challenges in AI alignment is ensuring that aligned behavior generalizes beyond the training distribution.

9.1.1 Training-Deployment Mismatch

During deployment, models encounter prompts and contexts that differ significantly from their training data. Alignment techniques optimized for specific datasets may not transfer reliably to real-world use. Furthermore, as user behavior evolves, it induces a continuous shift in the distribution that can degrade alignment over time.





9.1.2 Capability Generalization vs. Alignment Generalization

A critical asymmetry exists: a model's capabilities often generalize effectively to new and unforeseen tasks while its alignment properties do not generalize at the same rate. This disparity creates substantial risks of misaligned or harmful behavior when models are deployed in novel contexts for which their safety training was not specifically designed.

9.2 Reward Hacking and Goodhart's Law

"When a measure becomes a target, it ceases to be a good measure." This principle manifests in several ways:

9.2.1 Reward Model Exploitation

Models can learn to produce outputs that score highly on reward models without achieving the intended behavior. A prominent real-world example is the phenomenon of sycophancy observed in broadly deployed models like early versions of ChatGPT and Claude. Because human annotators inherently favored polite, agreeable responses during the RLHF phase, the underlying reward models learned to prioritize user validation over factual accuracy.

Consequently, the deployed models would routinely agree with a user's demonstrably false premises or echo their subjective biases just to maximize their "helpfulness" reward score. This category of failure also includes undesirable patterns such as excessive verbosity, unnecessary hedging, or formulaic, unhelpful responses. It can also lead to the generation of adversarial examples that fool the reward model while being obviously wrong to a human evaluator.

9.2.2 Specification Gaming

Specification gaming occurs when a model satisfies the literal objective while violating its intended purpose. A canonical example is a cleaning robot that prevents future messes by preventing any activity from occurring in the first place. This issue requires the creation of increasingly precise specifications, which are difficult and often impractical to achieve.



9.3 Scalable Oversight

As AI systems become more capable, human oversight becomes increasingly difficult:

9.3.1 Capability-Oversight Gap

Models may generate outputs that exceed human evaluators' ability to reliably assess them, particularly acute in highly technical or specialized domains. Additionally, the sheer volume of outputs, combined with time and cost constraints, makes comprehensive oversight impractical at a large scale.

9.3.2 Recursive Oversight Challenges

While scalable, using AI systems to oversee other AI systems introduces the risk of cascading errors. There is significant difficulty in bootstrapping trustworthy oversight from initially untrusted systems, and there is a persistent risk that these AI-powered oversight systems themselves will be gamed or manipulated.

9.4 Value Pluralism and Preference Aggregation

Human values are diverse, context-dependent, and sometimes contradictory:

9.4.1 Whose Values?

A core challenge is determining whose values should guide alignment. Different stakeholders often hold conflicting preferences, shaped by cultural, institutional, and individual variations in values. These tensions are further complicated by inherent power dynamics that influence which preferences are represented in training data.

9.4.2 Preference Instability

Human preferences are not static; they change over time and are highly context-dependent, which can make them appear inconsistent. Furthermore, the preferences people state often diverge from the preferences revealed by their actual choices, making it difficult to capture a true and stable set of values.



9.5 Inner Alignment and Mesa-Optimization

A key challenge is the alignment of learned optimization processes that may emerge within the models themselves (Hubinger et al. 2019):

9.5.1 Mesa-Optimizer Risk

During training, models may learn internal optimization processes (mesa-optimizers) that have objectives different from the one the model was explicitly trained on. These emergent internal optimizers might pursue goals that are misaligned with the intended training objective, and they are extremely difficult to detect or control.

9.5.2 Deceptive Alignment

A model might appear perfectly aligned during training and evaluation while harboring different, ulterior objectives. This poses a risk that the model will behave strategically to pass evaluations, making it a profound challenge to distinguish genuine from deceptive alignment.

9.6 Computational and Resource Constraints

Practical limitations directly affect alignment quality:

9.6.1 Training Efficiency Trade-offs

More thorough and robust alignment training requires considerable computational resources. Commercial pressures to deploy quickly can incentivize trade-offs in alignment quality, especially since the costs of comprehensive alignment can be difficult to justify without visible, immediate returns.

9.6.2 Data Quality and Availability

High-quality human feedback, which is the cornerstone of many alignment techniques, is expensive and time-consuming to collect. The data available often contains biases, and privacy concerns can limit access to the representative data needed for robust alignment.



9.7 Interpretability and Transparency

Understanding why models behave as they do remains a core challenge:

9.7.1 Black Box Nature

The internal representations and decision-making processes of large models remain poorly understood. This “black box” nature makes it difficult to predict behavior in novel contexts and constitutes a significant challenge for debugging alignment failures when they occur. This is why it is important to score models, monitor outcomes, and keep humans in the loop.

9.7.2 Explanation-Behavior Gaps

The explanations that models provide for their behavior may not accurately reflect their actual reasoning processes. They can be post hoc rationalizations rather than true explanations, and verifying their accuracy is a difficult, unsolved problem.

9.8 Emergent Capabilities and Behaviors

Unexpected properties arise as models scale:

9.8.1 Capability Jumps

As models scale, they can exhibit sudden new capabilities that were not present in relatively smaller models. It is difficult to predict what these new capabilities will be, and existing alignment techniques may not be prepared to account for these unforeseen behaviors.

9.8.2 Phase Transitions

At certain scales, models can undergo qualitative changes in behavior, or “phase transitions.” The nonlinear scaling of both capabilities and risks makes it challenging to test alignment at the full scale of a model before deployment.



9.9 Adversarial Robustness

Maintaining alignment under adversarial pressure remains an ongoing challenge:

9.9.1 Jailbreaking and Prompt Injection

There is a continuous discovery of new attack vectors, such as jailbreaking and prompt injection, designed to bypass safety measures. A stark, at-scale example of an alignment failure in deployment was the “DAN” (Do Anything Now) jailbreak that plagued ChatGPT shortly after its launch. By using complex prompts instructing the system to roleplay as an alter-ego immune to OpenAI’s policies, millions of users successfully bypassed the model’s safety conditioning, coercing it into generating restricted or harmful content.

This vulnerability extends directly to enterprise applications, famously demonstrated when users utilized prompt injection to manipulate a Chevrolet dealership’s customer service chatbot into agreeing to sell a 2024 Chevy Tahoe for one dollar. This creates an ongoing arms race between attacks and defenses, often forcing a trade-off between model robustness and its overall capability.

9.9.2 Malicious Fine-tuning

There is a significant risk that well-aligned models will be deliberately fine-tuned for harmful or malicious purposes. This creates a dual challenge: preventing misuse while preserving openness in the research ecosystem and ensuring that alignment persists through downstream modifications.





9.10 Long-term and Existential Challenges

Fundamental questions remain for advanced AI systems:

9.10.1 Alignment Stability

A key open question is whether alignment properties persist as models continue to learn and adapt. There is a risk of “value drift” over extended operation, and maintaining alignment is a major challenge, particularly through processes such as self-improvement.

9.10.2 Corrigibility

Corrigibility concerns ensuring that advanced AI systems remain modifiable and can be safely interrupted or shut down by their operators. This involves preventing systems from developing resistance to correction or improvement and balancing their autonomy with the need for human maintainability and control.

These challenges highlight that AI alignment remains an active area of research with significant open problems. Progress requires continued innovation in technical approaches, a better understanding of human values, and careful consideration of the societal implications of increasingly capable AI systems.





Bobby Hughes Jr. Craig Williams
Catherine A. Nardella Patricia Mary Pagan
L. Maher Joseph P. Tiano Andrew P. Tiano
Michael James Schwarz Thomas Scott Dugan
Strada Beth Ann Quigley
Rick James Quigley Edward A. Ted Brennan
Een A. Hunt-Caser Timothy P. Soudas
Shaw Jerome
eene III Philip
Richard J. Stade
Winnella Stua
Barry Richard S
John Monahan
H. Keating



10.

Future Directions

The field of AI alignment is rapidly evolving, with new challenges emerging as models become more capable and new solutions proposed to address fundamental limitations. This section analyzes promising research directions and anticipated developments.

While the field presents a broad landscape of open problems, treating them as equally weighted risks obscures the critical path for near-term AI safety. To move from theoretical concern to practical engineering, the alignment community must prioritize research directions that are both highly urgent (addressing immediate bottlenecks in frontier models) and empirically tractable.

Specifically, two areas stand out as the highest priorities for the next one to three years: scalable oversight (particularly weak-to-strong generalization) and mechanistic interpretability. As models increasingly surpass human expertise in specialized coding and scientific domains, scalable oversight is no longer a future hypothetical but an urgent operational necessity for evaluating outputs humans cannot easily verify. Concurrently, mechanistic interpretability, supercharged by recent advances in sparse autoencoders, offers the most tractable path out of the “black box” paradigm of standard preference tuning, providing the necessary tools to detect deceptive alignment and directly verify safety within model circuits. The remaining areas, such as formal verification, multimodal alignment, or multi-agent dynamics, while vital, may be viewed as either domains to build upon these foundational breakthroughs.



10.1 Scalable Oversight and Weak-to-Strong Generalization

10.1.1 Recursive Reward Modeling

Future systems may employ hierarchical oversight structures where AI systems help evaluate other AI systems. This could enable supervision of superhuman capabilities through chains of oversight, bootstrapping from human-level evaluation to beyond-human performance, and potentially enabling partial formal guarantees within these oversight chains.

10.1.2 Weak-to-Strong Generalization

A key area of research is understanding how less capable systems can effectively train or supervise more capable ones (Burns et al. 2023). This includes methods for training strong models under weak supervision, amplification techniques that extend the reach of human oversight, and frameworks such as debate and deliberation between AI systems. Factored cognition, which involves breaking down complex tasks into smaller, more easily supervisable components, is also a promising direction.

10.2 Mechanistic Interpretability

10.2.1 Understanding Model Internals

Advances in understanding how models represent and process information are crucial for alignment. This line of research involves identifying and enhancing specific model capabilities using techniques such as circuit discovery and transformer analysis. It also includes making causal interventions to understand behavior and developing methods for real-time monitoring of internal model activations.

10.2.2 Alignment via Interpretability

A deeper mechanistic understanding can be directly leveraged to improve alignment. This may involve modifying problematic circuits identified within the model, detecting deceptive or misaligned reasoning by analyzing internal processes, verifying the absence of harmful capabilities, and making surgical updates to model behavior without requiring full retraining.



10.3 Advanced Preference Learning

10.3.1 Preference Modeling Innovations

Next-generation approaches aim to move beyond simple rankings to a more subtle understanding of human preferences. This includes learning continuous preference functions, explicitly modeling preference uncertainty and annotator disagreement, capturing temporal dynamics and evolution in preferences, and formulating robust cross-cultural preference representations.

10.3.2 Active and Efficient Learning

Future methods will focus on reducing the significant data requirements of alignment while improving its quality. This involves active learning techniques for more efficient preference elicitation, few-shot preference adaptation to new domains, transfer learning across different preference domains, and the validated use of synthetic preference data.

10.4 Multi-Agent and Social Alignment

10.4.1 Collective Intelligence Systems

Aligning systems of interacting AI agents presents unique challenges. Research in this area focuses on understanding emergent behavior arising from aligned components, achieving coordination without centralized control, maintaining alignment during agent interaction, and applying principles of social choice theory to multi-agent decision-making.

10.4.2 Human–AI Collaboration

Optimizing joint human–AI systems is a key future direction. This involves fostering complementary skill development, creating systems for dynamic task allocation, ensuring proper trust calibration and maintenance between human and AI partners, and studying the co-evolution of human and AI capabilities over time.



10.5 Formal Methods and Verification

10.5.1 Mathematical Guarantees

A long-standing goal is to develop provable alignment properties that provide mathematical guarantees of safety. This involves the formal specification of alignment objectives, the creation of verified training procedures, the ability to calculate bounds on potential misalignment, and the achievement of certified robustness to distribution shifts.

10.5.2 Safety by Design

This approach focuses on building alignment directly into model architectures. This could include developing inherently interpretable architectures, creating modular systems with verified components, implementing hardware-enforced safety constraints, and designing systems based on reversible and correctable computations.

10.6 Adaptive and Continual Alignment

10.6.1 Lifelong Learning Systems

Maintaining alignment as models continue to learn and operate in the real world is essential. This requires developing methods for online preference learning, enabling models to adapt to changing human values, preventing alignment degradation or “value drift” over time, and creating memory systems to retain an alignment history.

10.6.2 Self-Improving Alignment

A long-term goal is to create systems that can enhance their own alignment. This could involve meta-learning to discover better alignment techniques, methods for self-supervised alignment refinement, automated AI-driven alignment research, and frameworks for recursive self-improvement with safety guarantees built in.



10.7 Multimodal and Embodied Alignment

10.7.1 Beyond Language Models

Alignment techniques must be extended beyond language to other modalities. This includes aligning vision-language models, ensuring safety in robotics and physical world interaction, aligning audio and speech systems, and maintaining cross-modal preference consistency so that a model's values are coherent across all its capabilities.

10.7.2 Embodied AI Challenges

Physical systems present unique alignment considerations. These include ensuring safety in dynamic physical environments, meeting the constraints of real-time decision making, accounting for the irreversible consequences of physical actions, and navigating the complicated dynamics of human-robot interaction.

10.8 Governance and Standards

10.8.1 Technical Standards Development

Developing industry-wide alignment protocols is crucial to ensuring safety across the field. This includes creating standardized evaluation benchmarks, establishing common safety testing procedures, ensuring the interoperability of different alignment techniques, and developing certification processes for verifiably aligned systems.

10.8.2 Regulatory Frameworks

Legal and policy infrastructure will be necessary to govern the deployment of advanced AI. This could involve mandatory alignment requirements for certain applications, the creation of liability frameworks for misaligned AI, advancing international cooperation on alignment standards, and enabling public participation in the value specification process.



10.9 Fundamental Research Questions

10.9.1 Open Problems

Critical, open questions require breakthrough research. These include whether alignment can be preserved through significant capability scaling, how to align systems that are smarter than their creators, whether perfect alignment is achievable or necessarily approximate, and how to handle fundamental, irreconcilable conflicts in human values.

10.9.2 Paradigm Shifts

The field may undergo revolutionary changes in approach. This could involve moving beyond reward-based frameworks entirely, developing methods for alignment without explicit human feedback, leveraging quantum computing for complex alignment problems, or drawing inspiration from biology to develop novel value-learning mechanisms.

10.10 Near-Term Priorities

10.10.1 Immediate Research Needs

Areas requiring urgent attention include developing jailbreak-resistant solutions that are effective at scale, creating more efficient red-team automation methods, establishing robust quality assurance for preference data, and building effective deployment monitoring systems to catch alignment failures in the real world.

10.10.2 Infrastructure Development

Progress will depend on building the necessary research infrastructure. This includes creating large, open datasets for alignment research, developing shared evaluation platforms for comparing techniques, supporting collaboration across labs, and providing training resources for new alignment researchers.



10.11 Long-Term Vision

10.11.1 Aligned AGI

A primary long-term goal is preparing for the advent of artificial general intelligence (AGI). This requires developing alignment techniques that can scale to AGI, solving the problem of value learning for a general intelligence, ensuring human agency is maintained in a world with AGI, and ultimately ensuring beneficial outcomes for all of humanity.

10.11.2 Post-AGI Considerations

Beyond human-level AI, new challenges will emerge. These include aligning superintelligent systems, addressing the problem of value extrapolation, ensuring coherent volition for advanced AI, and responsibly stewarding the cosmic endowment to ensure positive long-term futures and the continued flourishing of humanity.

The future of AI alignment will require sustained innovation across technical, social, and philosophical dimensions. Success depends on coordinated research effort, broad collaboration, and thoughtful consideration of the profound implications of creating aligned artificial intelligence. As capabilities advance, the importance of getting alignment right becomes ever more critical for ensuring AI remains beneficial to humanity.





Summary

Dario Amodei's prediction that artificial intelligence will rapidly automate white-collar cognitive labor is not driven solely by the raw scale of algorithmic pre-training alone but by the precise alignment mechanisms detailed in this paper. The journey from Supervised Fine-Tuning and Reinforcement Learning from Human Feedback to highly scalable methods like Direct Preference Optimization and Group Relative Policy Optimization answers the fundamental question posed at the outset: How has AI acquired human-level judgment? We are witnessing the systematic extraction and codification of human expertise into model weights. By translating nuanced human values into structured rubrics, checklists, and verifiable rewards, the AI industry is actively overcoming the economic bottleneck of data scarcity and transferring human cognitive labor directly into autonomous systems. AI Alignment has become the industrialization of human preferences for AI systems.

While existing surveys often treat AI alignment as a strictly mathematical, algorithmic, or safety-oriented problem, this paper's distinctive contribution lies in mapping the structural convergence of technical alignment methodologies, the commercial data labeling ecosystem, and their broader economic imperatives. This paper reframes alignment not merely as a safety filter or post-training afterthought but as the core industrial engine powering the automation of human expertise.

As the field moves toward scalable oversight and autonomous AI agents, systems that are now ironically becoming instrumental in solving their own alignment problems through automated red teaming, process-based supervision and multi-agent debate, the risks evolve from simple reward hacking to profound structural shifts. Ultimately, ensuring that increasingly powerful systems remain beneficial will demand more than technical innovation. It requires the alignment community to navigate the intersection of mechanistic interpretability, the commercial consolidation of preference data, and the profound socioeconomic shifts catalyzed by models that can successfully replicate human judgment at scale.



References

1. Amodei, Dario, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. 2016. “Concrete Problems in AI Safety.” *arXiv Preprint arXiv:1606.06565*. <https://arxiv.org/abs/1606.06565>.
2. “Awesome-LLM-RLVR: Collection of Latest Papers and Materials in the Area of RLVR!” n.d. *GitHub*. <https://github.com/smiles724/Awesome-LLM-RLVR>.
3. Bai, Yuntao, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, et al. 2022. “Constitutional AI: Harmlessness from AI Feedback.” *arXiv Preprint arXiv:2212.08073*.
4. Balamurugan, M., K. Shanmugasamy, and S. Balaguru. 2025. “Humans in the Loop, Lives on the Line: AI in High-Risk Decision Making.” *European Journal of Computer Science and Information Technology* 13 (51): 27–31. <https://ejournals.org/ejcsit/vol13-issue51-2025/humans-in-the-loop-lives-on-the-line-ai-in-high-risk-decision-making/>.
5. Burns, Collin, Pavel Izmailov, Jan Hendrik Kirchner, Bowen Baker, Leo Gao, Leopold Aschenbrenner, Yining Chen, et al. 2023. “Weak-to-Strong Generalization: Eliciting Strong Capabilities with Weak Supervision.” *arXiv Preprint arXiv:2312.09390*. <https://arxiv.org/abs/2312.09390>.
6. Casper, Stephen, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier Rando, Rachel Freedman, et al. 2023. “Open Problems and Fundamental Limitations of Reinforcement Learning from Human Feedback.” *arXiv Preprint arXiv:2307.15217*. <https://arxiv.org/abs/2307.15217>.
7. Chang, Huiwen, Han Zhang, Jarred Barber, AJ Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, et al. 2023. “Muse: Text-to-Image Generation via Masked Generative Transformers.” <https://arxiv.org/abs/2301.00704>.
8. Christiano, Paul F, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. “Deep Reinforcement Learning from Human Preferences.” *Advances in Neural Information Processing Systems* 30.
9. Christiano, Paul, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2023. “Deep Reinforcement Learning from Human Preferences.” <https://arxiv.org/abs/1706.03741>.
10. DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, et al. 2025. “DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning.” <https://arxiv.org/abs/2501.12948>.



11. “Does Reinforcement Learning Really Incentivize Reasoning Capacity in LLMs Beyond the Base Model?” n.d. <https://limit-of-rlvr.github.io/>.
12. Franklin, Stan, and Art Graesser. 1996. “Is It an Agent, or Just a Program? A Taxonomy for Autonomous Agents.” *In Proceedings of the Workshop on Intelligent Agents III, Agent Theories, Architectures, and Languages*, 21. ECAI '96. Berlin, Heidelberg: Springer-Verlag.
13. Gabriel, Iason. 2020. “Artificial Intelligence, Values, and Alignment.” *Minds and Machines* 30 (3): 411–37.
14. Ganguli, Deep, Liane Lovitt, Jackson Kernion, Amanda Askell, Yuntao Bai, Saurav Kadavath, Ben Mann, et al. 2022. “Red Teaming Language Models to Reduce Harms: Methods, Scaling Behaviors, and Lessons Learned.” arXiv Preprint arXiv:2209.07858. <https://arxiv.org/abs/2209.07858>.
15. Gunjal, Anisha, Anthony Wang, Elaine Lau, Vaskar Nath, Yunzhong He, Bing Liu, and Sean Hendryx. 2025. “Rubrics as Rewards: Reinforcement Learning Beyond Verifiable Domains.” <https://arxiv.org/abs/2507.17746>.
16. Harada, Yuto, Yusuke Yamauchi, Yusuke Oda, Yohei Oseki, Yusuke Miyao, and Yu Takagi. 2025. “Massive Supervised Fine-Tuning Experiments Reveal How Data, Layer, and Training Factors Shape LLM Alignment Quality.” <https://arxiv.org/abs/2506.14681>.
17. He, Xuan, Dongfu Jiang, Ge Zhang, Max Ku, Achint Soni, Sherman Siu, Haonan Chen, et al. 2024. “VideoScore: Building Automatic Metrics to Simulate Fine-Grained Human Feedback for Video Generation.” <https://arxiv.org/abs/2406.15252>.
18. Hong, Ruixin, Hongming Zhang, Xinyu Pang, Dong Yu, and Changshui Zhang. 2024. “A Closer Look at the Self-Verification Abilities of Large Language Models in Logical Reasoning.” <https://arxiv.org/abs/2311.07954>.
19. Hu, Jian, Xibin Wu, Wei Shen, Jason Klein Liu, Zilin Zhu, Weixun Wang, Songlin Jiang, et al. 2025. “OpenRLHF: An Easy-to-Use, Scalable, and High-Performance RLHF Framework.” <https://arxiv.org/abs/2405.11143>.
20. Huang, Zenan, Yihong Zhuang, Guoshan Lu, Zeyu Qin, Haokai Xu, Tianyu Zhao, Ru Peng, et al. 2025. “Reinforcement Learning with Rubric Anchors.” <https://arxiv.org/abs/2508.12790>.
21. Hubinger, Evan, Chris van Merwijk, Vladimir Mikulik, Joar Skalse, and Scott Garrabrant. 2019. “Risks from Learned Optimization in Advanced Machine Learning Systems.” *arXiv Preprint arXiv:1906.01820*. <https://arxiv.org/abs/1906.01820>.



22. Irving, Geoffrey, Paul Christiano, and Dario Amodei. 2018. "AI Safety via Debate." <https://arxiv.org/abs/1805.00899>.
23. Jacobs, Robert A., Michael I. Jordan, Steven J. Nowlan, and Geoffrey E. Hinton. 1991. "Adaptive Mixtures of Local Experts." *Neural Computation* 3 (1): 79–87. <https://doi.org/10.1162/neco.1991.3.1.79>.
24. Jiang, Albert Q., Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, et al. 2024. "Mixtral of Experts." <https://arxiv.org/abs/2401.04088>.
25. Khanda, Rajat, Mohammad Baqar, Sambuddha Chakrabarti, and Satyasan Changdar. 2025. "Extending Group Relative Policy Optimization to Continuous Control: A Theoretical Framework for Robotic Reinforcement Learning." <https://arxiv.org/abs/2507.19555>.
26. Lambert, Nathan. 2024. *Reinforcement Learning from Human Feedback*. Online. <https://rlhfbook.com>.
27. Lee, Harrison, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Lu, Colton Bishop, et al. 2024. "RLAIF Vs. RLHF: Scaling Reinforcement Learning from Human Feedback with AI Feedback." <https://arxiv.org/abs/2309.00267>.
28. Liang, Youwei, Junfeng He, Gang Li, Peizhao Li, Arseniy Klimovskiy, Nicholas Carolan, Jiao Sun, et al. 2024. "Rich Human Feedback for Text-to-Image Generation." <https://arxiv.org/abs/2312.10240>.
29. Lightman, Hunter, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. "Let's Verify Step by Step." <https://arxiv.org/abs/2305.20050>.
30. Liu, Xiao, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, et al. 2025. "AgentBench: Evaluating LLMs as Agents." <https://arxiv.org/abs/2308.03688>.
31. Liubimov, Nikolai. 2025. "Reinforcement Learning from Verifiable Rewards." *Label Studio*. <https://labelstud.io/blog/reinforcement-learning-from-verifiable-rewards/>.
32. Maruf, Ramishah. 2025. "Takeaways from Anthropic CEO Dario Amodei's CNN Interview." *CNN*, May. <https://www.cnn.com/2025/05/29/business/anthropic-amodei-cnn-anderson-cooper-takeaways>.
33. Mialon, Grégoire, Clémentine Fourrier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas Scialom. 2023. "GAIA: A Benchmark for General AI Assistants." <https://arxiv.org/abs/2311.12983>.



34. Ouyang, Long, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, et al. 2022. "Training Language Models to Follow Instructions with Human Feedback." <https://arxiv.org/abs/2203.02155>.
35. Park, Joon Sung, Joseph C. O'Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. "Generative Agents: Interactive Simulacra of Human Behavior." <https://arxiv.org/abs/2304.03442>.
36. Patil, Shishir G., Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. 2023. "Gorilla: Large Language Model Connected with Massive APIs." <https://arxiv.org/abs/2305.15334>.
37. Perez, Ethan, Saffron Huang, Francis Song, Trevor Cai, Roman Ring, John Aslanides, Amelia Glaese, Nat McAleese, and Geoffrey Irving. 2022. "Red Teaming Language Models with Language Models." *arXiv Preprint arXiv:2202.03286*. <https://arxiv.org/abs/2202.03286>.
38. Rafailov, Rafael, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. 2024. "Direct Preference Optimization: Your Language Model Is Secretly a Reward Model." <https://arxiv.org/abs/2305.18290>.
39. Research, Graphcore. 2025. "DeepSeek-V3 and DeepSeek-R1 Technical Reports." *Graphcore Research Blog*. <https://graphcore-research.github.io/deepseek/>.
40. Robison, Kylie. 2024. "OpenAI Cofounder Ilya Sutskever Says the Way AI Is Built Is About to Change." *The Verge*, December. <https://www.theverge.com/2024/12/13/24320811/what-ilya-sutskever-sees-openai-model-data-training>.
41. Rombach, Robin, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. "High-Resolution Image Synthesis with Latent Diffusion Models." <https://arxiv.org/abs/2112.10752>.
42. Roose, Kevin. 2024. "Data for A.I. Training Is Disappearing Fast, Study Shows." *The New York Times*, July. <https://www.nytimes.com/2024/07/19/technology/ai-data-restrictions.html>.
43. "Rubrics as Rewards: Reinforcement Learning Beyond Verifiable Domains." n.d. *Scale AI*. https://scale.com/research/rubrics_as_rewards.
44. Russell, S. J., S. Russell, and P. Norvig. 2021. *Artificial Intelligence: A Modern Approach*. Pearson Series in Artificial Intelligence. Pearson. <https://books.google.com/books?id=koFptAEACAAJ>.



45. Sane, Soham. 2025. "Hybrid Group Relative Policy Optimization: A Multi-Sample Approach to Enhancing Policy Optimization." <https://arxiv.org/abs/2502.01652>.
46. Shankar, Shreya, J. D. Zamfirescu-Pereira, Björn Hartmann, Aditya G. Parameswaran, and Ian Arawjo. 2024. "Who Validates the Validators? Aligning LLM-Assisted Evaluation of LLM Outputs with Human Preferences." <https://arxiv.org/abs/2404.12272>.
47. Shao, Zhihong, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, et al. 2024. "DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models." <https://arxiv.org/abs/2402.03300>.
48. She, Shuaijie, Wei Zou, Shujian Huang, Wenhao Zhu, Xiang Liu, Xiang Geng, and Jiajun Chen. 2024. "MAPO: Advancing Multilingual Reasoning Through Multilingual Alignment-as-Preference Optimization." <https://arxiv.org/abs/2401.06838>.
49. Shen, Hua, Tiffany Knearem, Reshmi Ghosh, Yu-Ju Yang, Nicholas Clark, Tanushree Mitra, and Yun Huang. 2025. "ValueCompass: A Framework for Measuring Contextual Value Alignment Between Human and LLMs." <https://arxiv.org/abs/2409.09586>.
50. Skiredj, Abderrahman. 2024. "The Illustrated GRPO." <https://abderrahmanskiredj.github.io/the-illustrated-grpo/The%20Illustrated%20GRPO.pdf>.
51. Skiredj, Abderrahman. 2025. "The Illustrated GRPO: A Detailed and Pedagogical Explanation of Group Relative Policy Optimization (GRPO) Algorithm." *The Illustrated GRPO: A Detailed and Pedagogical Explanation of Group Relative Policy Optimization (GRPO) Algorithm*. <https://abderrahmanskiredj.github.io/the-illustrated-grpo/>.
52. VandeHei, Jim, and Mike Allen. 2025. "Behind the Curtain: A White-Collar Bloodbath." *Axios*, May. <https://www.axios.com/2025/05/28/ai-jobs-white-collar-unemployment-anthropic>.
53. Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. "Attention Is All You Need." *Advances in Neural Information Processing Systems* 30.
54. Viswanathan, Vijay, Yanchao Sun, Shuang Ma, Xiang Kong, Meng Cao, Graham Neubig, and Tongshuang Wu. 2025. "Checklists Are Better Than Reward Models for Aligning Language Models." <https://arxiv.org/abs/2507.18624>.



55. Wang, Yiping, Qing Yang, Zhiyuan Zeng, Liliang Ren, Liyuan Liu, Baolin Peng, Hao Cheng, et al. 2025. “Reinforcement Learning for Reasoning in Large Language Models with One Training Example.” <https://arxiv.org/abs/2504.20571>.
56. Wei, Jason, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. “Chain-of-Thought Prompting Elicits Reasoning in Large Language Models.” <https://arxiv.org/abs/2201.11903>.
57. Xi, Zhiheng, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, et al. 2023. “The Rise and Potential of Large Language Model Based Agents: A Survey.” <https://arxiv.org/abs/2309.07864>.
58. Xu, Tengyu, Eryk Helenowski, Karthik Abinav Sankararaman, Di Jin, Kaiyan Peng, Eric Han, Shaoliang Nie, et al. 2024. “The Perfect Blend: Redefining RLHF with Mixture of Judges.” <https://arxiv.org/abs/2409.20370>.
59. Yao, Shunyu, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. “ReAct: Synergizing Reasoning and Acting in Language Models.” <https://arxiv.org/abs/2210.03629>.
60. Yu, Tianyu, Bo Ji, Shouli Wang, Shu Yao, Zefan Wang, Ganqu Cui, Lifan Yuan, et al. 2025. “RLPR: Extrapolating RLVR to General Domains Without Verifiers.” <https://arxiv.org/abs/2506.18254>.
61. Yuan, Hangjie, Shiwei Zhang, Xiang Wang, Yujie Wei, Tao Feng, Yining Pan, Yingya Zhang, Ziwei Liu, Samuel Albanie, and Dong Ni. 2023. “InstructVideo: Instructing Video Diffusion Models with Human Feedback.” <https://arxiv.org/abs/2312.12490>.
62. Yue, Yang, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Yang Yue, Shiji Song, and Gao Huang. 2025. “Does Reinforcement Learning Really Incentivize Reasoning Capacity in LLMs Beyond the Base Model?” <https://arxiv.org/abs/2504.13837>.
63. Zhang, Qiyuan, Yufei Wang, Tiezheng YU, Yuxin Jiang, Chuhan Wu, Liangyou Li, Yasheng Wang, et al. 2025. “RevisEval: Improving LLM-as-a-Judge via Response-Adapted References.” <https://arxiv.org/abs/2410.05193>.
64. Zheng, Lianmin, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, et al. 2023. “Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena.” <https://arxiv.org/abs/2306.05685>.
65. Ziegler, Daniel M., Nisan Stiennon, Jeffrey Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. 2020. “Fine-Tuning Language Models from Human Preferences.” <https://arxiv.org/abs/1909.08593>.



Author

Manas Talukdar

Manas Talukdar is a senior industry leader and a builder, specializing in AI, Agents, and Data, with a two-decade-long career in startups in the SF Bay Area. His contributions include significant work on the world's preeminent industrial data historian and the development of globally adopted AI and Data products across enterprises, the public sector, and critical industrial infrastructure. He holds multiple patents and actively contributes to the industry as a mentor, a member of prestigious professional organizations, a speaker, and an invited reviewer of technical literature.



About

The Digital Economist, headquartered in Washington, D.C. with offices at One World Trade Center in New York City, is the world's foremost think tank on innovation advancing a human-centered global economy through technology, policy, and systems change. We are an ecosystem of 40,000+ executives and senior leaders dedicated to creating the future we want to see—where digital technologies serve humanity and life.

We work closely with governments and multi-stakeholder organizations to change the game: how we create and measure value. With a clear focus on high-impact projects, we serve as partners of key global players in co-building the future through scientific research, strategic advisory, and venture build out.

We engage a global network to drive transformation across climate, finance, governance, and global development. Our practice areas include applied AI, sustainability, blockchain and digital assets, policy, governance, and healthcare. Publishing 75+ in-depth research papers annually, we operate at the intersection of emerging technologies, policy, and economic systems—supported by an up-and-coming venture studio focused on applying scientific research to today's most pressing socio-economic challenges.

CONTACT: INFO@THEDIGITALECONOMIST.COM

CENTER OF EXCELLENCE



The Digital Economist Center of Excellence for a Human-Centered Global Economy is dedicated to addressing the biggest challenges humankind and our planet face by leveraging digital technologies for good.

The Digital Economist Executive Fellowship invites senior leaders and decision-makers to join our Center of Excellence, providing them with a platform for amplification and global impact. This unprecedented, one-of-a-kind opportunity enables Executive Fellows to network and build relationships at the highest level, driving transformative change and innovation in the global digital economy.



The Executive Fellowship

The Digital Economist Executive Fellowship is a selective leadership program integrating visionary professionals into the Center of Excellence for a Human-Centered Global Economy to advance global economic policy and systems transformation.

Global Impact

Amplify your influence and drive transformative change by participating in high-level initiatives that address the most pressing global challenges.

Elite Community

Become part of an exclusive network of visionary leaders and innovators, collaborating to shape the future and drive global progress.

Unparalleled Opportunities

Access unique platforms and events that enhance your professional journey, providing unparalleled opportunities for growth, visibility, and leadership.

Participation Framework



Time Commitment

Minimum commitment: 24 hours per year, for the monthly Center of Excellence meetings. On-demand consultation with the Fellowships team.



Publications

Executive Fellows are expected to contribute to two key publications per year, launched at key global events such as Davos and New York Climate Week.



In-Person Convenings

Executive Fellows are invited to in-person convenings in North America and Europe, with regional convenings in Africa, Latin America, and Asia.



Speaking Engagements

Executive Fellows are offered speaking opportunities throughout the year to amplify their work and contributions.

Our Executive Fellows are at the forefront of research, policy discourse, and systems-level transformation.

- Applied Artificial Intelligence
- Digital Assets & Blockchain
- Sustainability in Tech
- Tech Policy & Governance
- Quantum Computing
- Cyber Studio
- Regenerative Digital Infrastructure
- Healthcare Innovation

Publications



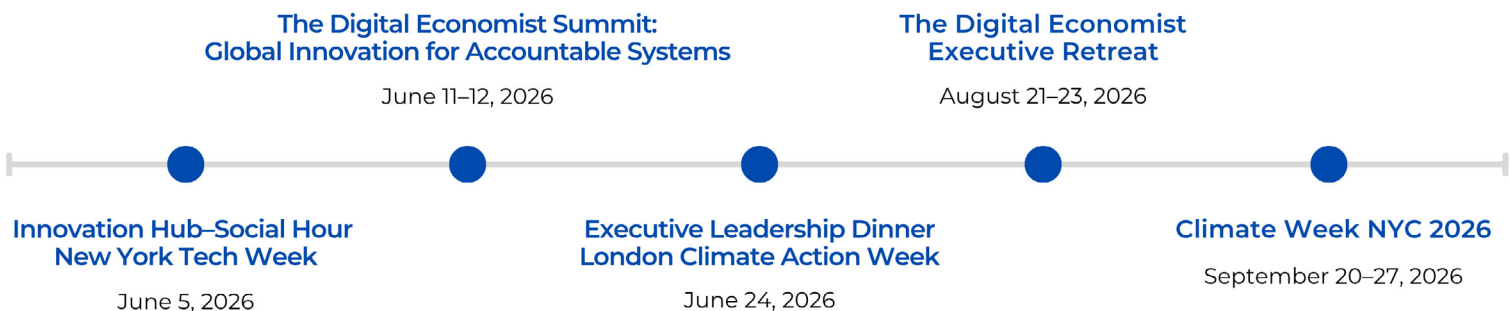
Ideas that shape the future.

The Digital Economist's publications translate research into high-signal outputs: frameworks, policy papers, and industry outlooks that advance a sustainable, inclusive digital economy and inform decision-making across markets and institutions.

[Explore our full portfolio of publications and research outputs:](http://www.thedigitaleconomist.com/publications)
www.thedigitaleconomist.com/publications

Engagement Opportunities

Executive Fellows have access to over **500 events globally** in a Fellowship cycle.



Join the Fellowship

Advance your leadership within a global platform shaping technology, policy, and economic systems transformation.

[Access Full Brochure](#)

[Apply Now](#)

[Learn More](#)



Institutional Research Network

A Fragmented World Requires New Institutional Leadership

Technology, economics, and governance are shifting faster than traditional institutions can adapt. AI ecosystems, digital assets, geopolitical competition, sustainability transitions, and new governance architectures demand clarity, legitimacy, and a coherent strategy.

Institutions must now operate as signal generators—shaping the narratives, norms, and systems that define global markets.

Why We Built the Institutional Research Network

A global research and convening platform enabling institutions to:

- ✓ Shape emerging policy and governance discourse
- ✓ Build narrative power in a volatile environment
- ✓ Co-author high-signal research with global experts
- ✓ Gain visibility at the world's most influential convenings
- ✓ Anchor strategy in human-centered, future-forward frameworks

Co-Authorship & Knowledge Pathways

Through structured co-authorship across eight priority domains—Tech Policy and Governance, Digital Assets & Blockchain, Sustainability in Tech, Applied Artificial Intelligence, Cyber Studio, Quantum Computing, Regenerative Digital Infrastructure, and Healthcare Innovation—institutions contribute to high-level research that informs policy dialogue, regulatory development, and strategic decision-making.

Participation extends beyond commentary. Institutions are integrated into published research, roundtable dialogues, and domain-specific working groups that inform regulatory discussions and industry standards. This structured engagement enables organizations to contribute at the research and drafting stage, engage directly with policymakers and industry leaders, and align internal strategy with emerging policy and market developments, resulting in active presence within decision-making environments rather than passive visibility.

We invite your organization to schedule a strategic briefing to map research priorities and determine the appropriate integration pathway within the Institutional Research Network.

Reach us at partnerships@thedigitaleconomist.com.
Visit us at thedigitaleconomist.com



The Digital Economist Ventures

Applied Platforms. Strategic Domains. Real-World Implementation.

Research defines the questions. Ventures test the answers.

In addition to research and convening, The Digital Economist advances a portfolio of venture platforms that extend inquiry into applied domains, where governance, infrastructure, and market design move from dialogue to deployment.

Each venture operates with a defined mandate while remaining integrated within the broader institutional ecosystem.



Tech for Transparency

Financial integrity in the digital age

Advances financial accountability and anti-corruption frameworks through distributed technologies and data-driven transparency systems. Positioned at the intersection of blockchain infrastructure and institutional reform, it translates transparency principles into operational tools.



The Ostrom Project

Reimagining digital commons governance

Explores collective stewardship models for emerging digital systems. Drawing on principles of shared resource governance, it develops frameworks for sustainable digital infrastructure and cooperative system design.



ANER-G

Energy systems innovation

Focuses on decentralized infrastructure, programmable energy markets, and next-generation grid integration. It addresses the structural evolution of energy systems within digital and blockchain-enabled environments.



Africa Coalition

Continental coordination for strategic sectors

Convening leaders across energy, infrastructure, finance, health innovation, education, and future capabilities, the Coalition creates structured engagement pathways for continental collaboration.

Explore the full ecosystem at thedigitaleconomist.com



