

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/383611061>

Unlocking the Wisdom of Large Language Models: An Introduction to The Path to Artificial General Intelligence

Book · August 2024

CITATIONS

0

READS

2,739

1 author:



Edward Y Chang

Stanford University

296 PUBLICATIONS 11,213 CITATIONS

SEE PROFILE

UNLOCKING THE WISDOM OF LARGE LANGUAGE MODELS



An Introduction to The Path to
Artificial General Intelligence

EDWARD Y CHANG

Unlocking the Wisdom of Large Language Models

**An Introduction to
The Path to Artificial General Intelligence**

Edward Y. Chang

Computer Science
Stanford University

Preface

This booklet, “*Unlocking the Wisdom of Large Language Models (LLMs)*,” serves as an introduction to the comprehensive work “*LLM Collaborative Intelligence: The Path to Artificial General Intelligence*” [1]. Through eleven carefully crafted aphorisms, we distill the core insights and principles that underpin the broader exploration of AI’s future through **LLM Collaborative Intelligence (LCI)**—a framework we propose as a promising pathway toward realizing **Artificial General Intelligence (AGI)**.

In addition to these insights, the booklet provides the titles, abstracts, and introductions of the main chapters from the book and presents the first two chapters in their entirety.

What follows is the original preface of “*LLM Collaborative Intelligence: The Path to Artificial General Intelligence*,” offering critical context for the journey ahead. We are excited to announce the release of the book’s second edition this week, featuring substantial enhancements in Chapters 6, 7, 8, and 9. Additionally, we have revised the preface to address recent comments from Yann LeCun regarding the feasibility of AGI.

Yann LeCun has expressed skepticism about the ability of LLMs to achieve AGI, citing limitations in memory, planning, and world grounding. However, we argue that **LCI**, especially with multiple multimodal LLMs working together—provides a promising architecture to overcome these challenges and move toward AGI. The chapters on **SocraSynth**, **EVINCE**, behavior

modeling, consciousness modeling, and the LCI-based architecture—which assigns executive, legislative, and judicial roles inspired by a three-branch system of government—demonstrate that through checks and balances, LLMs can collaborate constructively to achieve levels of intelligence that a single LLM cannot reach alone.

By fostering collaboration through structured dialogue, LCI enables individual LLMs to complement each other’s strengths, generating insights that no single model could achieve in isolation. This approach mirrors how human institutions balance diverse perspectives to produce better outcomes.

We can further extend this concept by designing one LLM to construct a virtual “world model”—mapping conceptual relationships—and another to act as a super sensor, leveraging instruments to detect signals beyond human sensory capabilities. Together, these specialized LLMs can model and perceive reality in ways that transcend individual human limitations.

Time will tell if this collaborative framework will bring us closer to AGI, but the path forward lies in cooperation, not isolation. Just as collective intelligence has propelled human progress, collaborative LLMs may unlock new frontiers of artificial intelligence that no single model can achieve alone.

Edward Y. Chang,
August 15th, 2024.
October 28th, 2024 (second edition).
Series π 0010828G002

Preface to “LLM Collaborative Intelligence: The Path to Artificial General Intelligence”

As Generative AI transforms our world, experts predict the emergence of Artificial General Intelligence (AGI) as early as 2040. This book proposes that the key to achieving AGI, characterized by versatility, adaptability, reasoning, critical thinking, and ethical alignment, lies not in creating more powerful individual models, but in enabling large language models (LLMs) to engage in intelligent, collaborative dialogue. This concept, which we term LLM Collaborative Intelligence (LCI), forms the foundation of our exploration.

LCI transcends conventional “mixture of experts” models or traditional LLM debates by optimizing information exchange through three essential foundations:

First, it balances exploration and exploitation, leveraging diverse perspectives while maintaining robust reasoning rooted in the strong priors of next-token maximum likelihood predictions. This approach ensures the exploration of novel perspectives while preserving the stability of model parameters learned from training data.

Second, in addition to balancing exploration and exploitation, LCI modulates the linguistic behaviors of LLMs, enabling dynamic shifts between contention and collaboration. This fosters productive exchanges that generate diverse perspectives and novel insights. By adjusting the contentiousness of debates, LCI explores a broader spectrum of ideas and arrives at well-reasoned conclusions.

Third, distinct roles are assigned to individual LLMs: the *Executive* for knowledge formation, the *Legislative* for developing ethical frameworks, and the *Judicial* for contextual evaluation.

This structure, inspired by the three-branch model of government, ensures balanced decision-making through a system of checks and balances, safeguarding ethical considerations.

Building on these foundational principles, the book’s twelve chapters systematically guide readers from theoretical frameworks to practical applications. Key algorithms include **CRIT** for critical evaluation, **SocraSynth** for dynamic dialogues, and **EVINCE** for optimizing information flow through Bayesian statistics and information theory. We explore applications ranging from medical diagnosis to news debiasing, and address fundamental challenges in AI safety by applying checks and balances across LLM modules—ensuring alignment between knowledge, ethics, and context-sensitive reasoning.

Some pioneers of modern AI, including Yann LeCun, argue that LLMs alone cannot achieve AGI, citing their limitations in persistent memory, planning ability, and physical world grounding. LeCun contends that true intelligence requires direct interaction with the physical world through sensors and embodiment, and that language models are like sophisticated parrots—eloquent but lacking true understanding of the physical world they describe. In his view, text-based models, however sophisticated, are fundamentally constrained by their training modality and cannot independently develop general intelligence.

However, this perspective may underestimate both the nature of human intelligence and the potential of structured LLM collaboration. Consider human perception: our sensory apparatus captures only narrow bands of electromagnetic and acoustic spectra. We cannot see infrared or ultraviolet light, hear ultrasonic frequencies, or sense magnetic fields, yet through collaborative knowledge sharing and abstract reasoning, we have

developed deep insights into phenomena far beyond our direct sensory experience: from quantum mechanics to celestial bodies. Indeed, most of human knowledge today comes not from direct sensory experience but from accumulated learning and collaborative reasoning.

Modern LLMs increasingly demonstrate multimodal integration, combining text, vision, audio, and other data streams. More importantly, unlike humans who are constrained by our biological sensors, LLMs through LCI can not only interface with but also meaningfully integrate data from an ever-expanding array of data acquisition tools—from radio telescopes detecting signals from distant galaxies to electron microscopes revealing atomic structures, from quantum sensors measuring subatomic phenomena to satellite systems monitoring global patterns. While individual LLMs may lack direct physical interaction in the conventional sense, their ability to process and integrate data from these advanced instruments, combined with sophisticated communication and collaboration capabilities, enables them to develop understanding far beyond human sensory limitations, potentially uncovering patterns and relationships that might elude human observers.

The path to AGI through LLM communication, while still a hypothesis, represents a gradual, collaborative process rather than a sudden leap. Just as human civilization advanced through collective knowledge accumulation rather than individual sensory improvement, LLMs through LCI may overcome their current limitations through structured collaboration. This book presents compelling evidence, methods, and frameworks supporting this direction. As we stand at this crucial juncture in AI development, our exploration of collaborative intelligence may

prove fundamental to achieving systems that not only match but potentially transcend human-level reasoning—not through individual model capabilities alone, but through the power of structured, ethical collaboration.

Edward Y. Chang,

March 27th, 2024 (first edition).

October 28th, 2024 (second edition).

Series π 0010830G006

Acknowledgments

This work has been significantly shaped by several pioneering perspectives in the field. Microsoft CSO Eric Horvitz’s concept of “polydisciplinary” representation in LLMs fundamentally influenced my understanding of how these models synthesize knowledge across domains. Demis Hassabis of DeepMind contributed crucial insights about end-to-end architecture design, challenging conventional thinking about human intervention points. Prof. Stuart Russell’s observation about LLMs mimicking human linguistic behavior to achieve goals became central to our approach in modeling behaviors and applying Bayesian frameworks.

My colleague Prof. Vaughan Pratt’s feedback on the initial paper inspired the use of “contentiousness” to foster adversarial debate, a key component of our framework. I am grateful to Prof. Monica Lam for her collaborative support throughout this research. Additionally, I thank my colleagues at Stanford Clinical Mind AI for their invaluable domain knowledge in disease diagnosis.

Contents

Preface	i
1 Aphorisms of LLM Collaborative Intelligence	1
2 A Brief History of AI: From Turing to Transformers	19
2.1 Definitions	23
2.2 Perspectives on Similarity	26
2.3 Eras of Similarity Measurement.	36
2.4 Concluding Remarks	50
3 Capabilities and Opportunities of Large Language Models	53
3.1 Distinctive Capabilities	56
3.2 Exploring Unknown Unknowns	67
3.3 Concluding Remarks	78
4 Prompt Engineering: Few Shots, Chain of Thought, and Retrieval-Augmented Generation	81
5 CRIT: Socratic Inquiry for Critical Thinking in LLMs	85
6 SocraSynth: Adversarial Multi-LLM Reasoning	89
7 EVINCE: Optimizing Adversarial LLM Dialogues via Conditional Statistics and Information Theory	95
7.1 Introduction	96

Contents

8 Uncovering Errors and Biases with Reflective Large Language Models	101
9 Modeling Emotions in Multimodal LLMs	105
10 A Three-Branch Checks-and-Balances Framework for Context-Aware Ethical Alignment of Large Language Models	109
11 Beyond Computation: Consciousness Modeling	115
12 A Retrospective and Adaptive Framework to Improve LLMs	119
13 Discovering Insights Beyond the Known	123
14 Bibliography	129
Appendix: Online Chapters	149
Author's Biography	153

1 Aphorisms of LLM

Collaborative Intelligence

Through extensive work with Large Language Models since the advent of GPT-2, especially in developing LLM Collaborative Intelligence (LCI) [1], I have distilled my observations into eleven (and counting) thought-provoking aphorisms. These insights reveal both the promising potential and the inherent limitations of LLM interactions.

While my experiences with LCI suggest promising directions, particularly how adversarial LLM roles can mitigate hallucinations and enhance reasoning, we must approach these observations with balanced enthusiasm. Some limitations we understand, like the need for sophisticated LLMs to achieve meaningful dialogue (as discussed in Aphorism #4), while others may only become apparent through further research and real-world applications.

These aphorisms aim not to pronounce definitive truths but to invite exploration and discussion. They reflect current observations about how LLMs interact, reason, and complement each other, while acknowledging that our understanding continues to evolve. I invite readers to engage with these ideas both critically and imaginatively, using them as stepping stones for further investigation rather than fixed conclusions.

Aphorisms

Aphorism #1

“Effective LLM interaction hinges on framing and sequencing the right questions.”

This aphorism highlights the critical role of precise questioning in both learning from and working with LLMs. Within LCI, particularly through frameworks like [SocraSynth](#) [2], this principle manifests throughout the system’s operation.

In multi-LLM debates, the quality of discourse depends heavily on how LLMs formulate their queries and challenges to each other. Effective counterarguments serve as sophisticated questions, probing assumptions, requesting evidence, or highlighting potential inconsistencies. This transforms debate from mere opposition into collaborative inquiry.

The iterative nature of LLM interactions fosters dynamic exchanges, where each response refines the shared context for subsequent questions. Through iterative context construction, LLMs engage with each other’s arguments to build deeper, more focused inquiries, enabling increasingly precise and relevant information exchanges.

Building on this foundation, the principle extends beyond question formulation to question sequencing. The order and relationship between queries profoundly impact the depth and breadth of exploration. In LCI, LLMs learn to build upon earlier exchanges, using responses to inform and refine subsequent queries, creating a scaffolded approach to complex topics through incremental discovery.

This virtuous cycle of precise questioning, contextual enrichment, and strategic sequencing determines the quality and utility of LLM-generated insights.

Example: A conversation that begins with a question about the symptoms of a disease may evolve into a sophisticated diagnostic discussion as LLMs build on each other’s insights, transitioning from surface-level symptoms to underlying causes, identifying comorbidities, and generating actionable recommendations for treatment [3, 4]

Aphorism #2

“Hallucinations rarely repeat.”

Have you ever wondered why recurring nightmares, even when they share the same theme, never quite unfold in exactly the same way? This aspect of human dreaming finds a parallel in how Large Language Models (LLMs) process information. Just as nightmares rarely recur in exactly the same way, even when themes repeat, LLM hallucinations exhibit a similar non-deterministic quality. This characteristic distinguishes hallucinations from systematic errors and offers both challenges and opportunities for detection and mitigation.

The phenomenon stems from three key mechanisms:

First, hallucinations arise when probabilistic token prediction leads to unpredictable sequences due to ambiguous or insufficient input. Since token selection draws from a probability distribution, even similar inputs result in varied hallucinations. This contrasts with systematic errors, which consistently emerge from gaps or biases in training data.

Second, in LCI frameworks like SocraSynth [2], the non-repetitive nature of hallucinations becomes a strength. When one LLM produces a hallucination, others can challenge it with counter-arguments. Because the original LLM is unlikely to reproduce

Aphorisms

the same hallucination, it must either ground its response in factual knowledge or acknowledge uncertainty. This iterative interaction creates a self-correcting dynamic where hallucinations naturally diminish over time.

Third, while hallucinations are sporadic, true knowledge gaps are consistent and can be systematically addressed. By integrating Retrieval-Augmented Generation (RAG), we can distinguish random hallucinations from persistent knowledge deficits, allowing for targeted improvements to the model’s knowledge base.

This understanding has implications for LLM system design: while safeguards against hallucinations are necessary, we can also leverage their non-repetitive nature in multi-LLM architectures to build self-correcting systems that dynamically adapt to evolving challenges.

Aphorism #3

“Strength and weakness in an LLM are not fixed traits but fluid, shifting with context. LCI empowers LLMs to transcend training biases, adopting new stances through structured debate.”

Consider how a vegetarian might eat meat in survival situations, or how persistence can become stubbornness in the wrong context. Similarly, an LLM’s capabilities shift significantly with situational context, challenging the notion that they are bound by fixed behavioral patterns from their training data.

The phenomenon manifests in three key ways:

First, what appears as a bias or limitation in one context can become an advantage in another. For instance, an LLM’s cautiousness may hinder brainstorming but becomes an asset in providing medical or financial advice. Conversely, a creative

tendency that thrives in artistic tasks may be a drawback when precision is needed, such as in legal documentation.

Second, in an LCI implementation such as *SocraSynth* [2], LLMs engage in structured debate, where counterarguments push models to re-evaluate their assumptions and shift stances, mimicking human discourse. This contextual flexibility allows them to transcend training-induced biases, refining their responses based on new evidence and perspectives.

Third, this adaptability suggests that the key to improving LLM performance lies not solely in better training data but in effective context management. By carefully structuring the operational context, we can guide LLM behavior toward more desirable outcomes without altering underlying parameters.

This understanding shapes AI system design: rather than viewing LLM behaviors as static traits to be corrected through retraining, we can leverage contextual adaptation to enhance both performance and reliability. This approach reduces the need for frequent retraining, and improves the scalability and responsiveness of AI systems, making them more adaptable to real-world environments with evolving demands.

Aphorism #4

“Critical thinking requires more than one Socrates.”

Just as the Socratic method relies on dialogue to uncover truth, LCI shows that critical thinking emerges through the interaction of multiple LLMs. Through frameworks like *SocraSynth* [2], *EVINCE* [4], and *DIKE-ERIS* [5], this principle enables structured debates that extend beyond individual model limitations.

The principle manifests in several dimensions:

Aphorisms



Figure 1.1: The LCI Framework with Two Socrates.

First, in dialectical reasoning, each LLM serves as both questioner and respondent. Like Socrates engaging with his interlocutors, one LLM may challenge assumptions while another defends or refines them. This dynamic fosters robust intellectual exchange, with each interaction building upon the analysis, for instance by proposing hypotheses and examining underlying assumptions.

Second, dialogue depth varies based on the participants' capabilities. Advanced LLMs (e.g., GPT-4 instances) can explore complex ideas with *depth*, leveraging knowledge and reasoning to refine and challenge each other's perspectives. By contrast, simpler models (e.g., GPT-2) may engage in shallow exchanges or reinforce misconceptions, akin to novices struggling to delve into complex topics.

Third, meaningful collaboration requires a *baseline* level of competency. Just as Socrates could not derive insights through

dialogue with those lacking reasoning skills, LCI cannot produce valuable results when models are too limited. Two weak reasoners combined do not yield strength; their limitations may even compound, undermining productive discourse.

This introduces key considerations for dialectical AI systems:

Challenge: LLMs may reinforce each other's errors. This requires mitigation through model diversity or human oversight.

Opportunity: As LLM capabilities improve, the potential for effective multi-model dialogue grows.

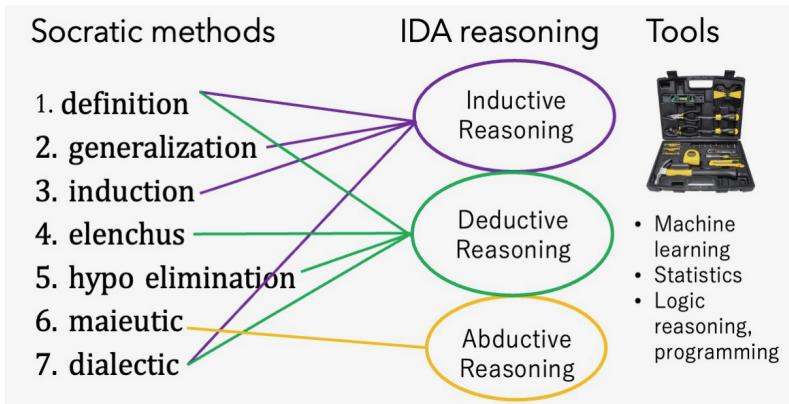


Figure 1.2: CRIT: Critical Inquisitive Template. Mapping from individual Socratic methods to reasoning methods.

LCI creates a space for AI ‘philosophers’ to engage in structured dialogue through frameworks like SocraSynth and DIKE-ERIS, monitored by EVINCE for information flow and by CRIT [6] for reasoning quality. The quality of these dialogues reflects the capabilities and of the participants. This understanding informs model selection for analytical tasks and sets realistic expectations for their performance.

Aphorisms

Aphorism #5

“LLMs are designed and trained to emulate human linguistic endeavors, each aimed at fulfilling distinct human objectives.”

LLMs function less like fortune tellers predicting words and more like method actors drawing from countless human performances. Through frameworks like SocraSynth [2], LCI creates environments where these actors shift between roles: the historian documenting events, the lawyer crafting arguments, the poet weaving emotions into verse, the teacher explaining complex ideas.

When humans write, we do so with intent—to persuade, inform, tell stories, express emotions, or explore ideas. Through training, LLMs have absorbed these patterns of purposeful communication. They aren’t simply predicting words; they are channelling the human intentions behind them.

Consider some of the roles LLMs can adopt:

- The Journalist: Organizing facts into coherent narratives.
- The Debater: Building and challenging arguments.
- The Analyst: Breaking down complex problems.
- The Mediator: Finding common ground between differing perspectives.

These functions reflect the complexity of human linguistic behavior. LLMs demonstrate the ability to pursue human objectives through language, from informing and educating to persuading and entertaining. LCI frameworks like SocraSynth and DIKE-ERIS dual [5] assign LLMs specific roles with defined contexts and moderation.

Aphorism #6

“Outside formal systems and physical laws, there is seldom ground truth; there is primarily reasonableness.”

While mathematics and formal logic provide definitive truths through proofs and theorems, most real-world domains deal with degrees of certainty rather than absolute truth. AI systems, like humans, must operate in contexts where insights are assessed based on reasonableness, plausibility, and situational appropriateness.

In many fields, from natural language processing to strategic decision-making, answers span a spectrum of interpretations, each valid in its context.

LCI embodies this through frameworks like SocraSynth [2], facilitating dialogue between multiple LLMs that contribute diverse perspectives. Through these interactions, the system builds understanding from data, similar to how humans refine insights through debate. This method aligns with the Socratic tradition of questioning assumptions and exploring viewpoints, leading to more thoughtful and context-aware conclusions.

In AI, reasonableness means synthesizing plausible outcomes that respect the inherent uncertainty in complex contexts, rather than pursuing absolute correctness. For example, in medical diagnostics, reasonableness translates to offering a differential diagnosis with the most likely conditions, allowing clinicians to evaluate and weigh possibilities rather than relying on a single, potentially misleading interpretation. This approach minimizes the risk of rigid adherence to single interpretations based on maximal likelihood.

By considering multiple perspectives through conditional statis-

Aphorisms

tics, LCI provides balanced and contextually appropriate responses. This flexibility enhances adaptability and reduces the chance of reinforcing flawed conclusions.

The challenge lies in balancing diverse perspectives with decisive insights. LCI addresses this by upholding rigor in formal truth domains while embracing reasonableness elsewhere. The entropy duality theory in EVINCE [4] introduces information-theoretic metrics to guide this balance, enabling exploration of new ideas while grounding them in established knowledge.

This approach acknowledges both the existence of absolute truth in formal domains and the need for reasonable interpretation in complex, real-world situations.

Aphorism #7

“Objectivity is the ‘hard problem’ in philosophy, and what we can do is unearthing all perspectives.”

Think of objectivity as that one friend who claims to be completely unbiased, we all know they’re not, but we appreciate the effort! This aphorism acknowledges that perfect objectivity is like trying to see the back of your own head without mirrors: theoretically possible with enough mirrors, but practically quite tricky.

In LCI, we tackle this challenge by essentially setting up a room full of mirrors, each LLM reflecting different angles of the truth. Rather than claiming to have found the philosopher’s stone of perfect objectivity, we embrace the wisdom of multiple perspectives:

- Like a potluck of perspectives, each LLM contributes its own flavor to the discussion, enriching the final outcome beyond

what any single model could achieve.

- When one LLM asserts, “That’s definitely correct,” another counters with “Well, actually...,” fostering honesty through friendly debate.
- Instead of pursuing a single, universal “objectivity,” we gather a diversity of subjective perspectives, each one valid within its unique context.

This approach recognizes that objectivity isn’t a destination but a journey—one best traveled with multiple companions. By orchestrating lively debates among our AI discussants, LCI offers a more comprehensive view, even if we never fully reach that mythical point of perfect objectivity. For instance, news outlets are categorized to be between left to right, and news articles on virtually every issue can reveal different perspectives [7]. Sometimes, the best way to approach truth is to admit that we’re all a little biased and invite everyone to share their perspectives openly.

This approach views objectivity as an ongoing process rather than a destination. Through orchestrated debates among AI discussants, LCI builds comprehensive understanding. News coverage illustrates this principle: outlets span the political spectrum from left to right, with articles on any issue revealing different perspectives [7]. The path to understanding often begins by acknowledging our biases and examining multiple viewpoints.

Note: Subjective perspectives can be evaluated through fact-checking and reasoning analysis.

Aphorism #8

“LLMs are not taught about domain boundaries, as they were trained only to predict the next words. Such a polydisciplinary

Aphorisms

approach to information representation allows LLMs to synthesize knowledge that might be beyond narrowly focused, domain-specific human understanding.”

The term “polydisciplinary” was introduced by Microsoft’s Chief Scientific Officer, Eric Horvitz, at Stanford’s HAI center in 2023. He noted that GPT-4’s training process, predicting next tokens through maximum likelihood, applies the same statistical approach whether processing physics equations or poetry: the model knows no disciplinary boundaries.

While humans organize knowledge into categories, such as physics, poetry, biology, and philosophy, LLMs move across these boundaries, unaware of traditional divisions. This creates new possibilities [8]:

- When an LLM sees patterns between Shakespeare and quantum mechanics, it doesn’t question whether that connection is “allowed.”
- Problems that typically need expertise from multiple fields might be approached by an LLM that never learned about academic silos.
- New hypotheses might emerge because LLMs don’t know which connections they shouldn’t make.

This trait shapes LCI’s approach: we use LLMs to explore questions we haven’t thought to ask (see Aphorism #1 on framing questions). Through frameworks like SocraSynth [2], EVINCE [4], and DIKE-ERIS dual [5], LCI lets LLMs reveal paths we might miss. By connecting these perspectives (as noted in Aphorism #4 about critical thinking), we can find links across disciplines.

Insights may lie not in individual disciplines but in the spaces between them: spaces that LLMs traverse naturally. However,

in this unbounded sphere of “unknown unknowns,” humans may not be able to distinguish between treasures and nonsense.

Aphorism #9

“Our public behavior isn’t a direct, unfiltered output from our unconscious mind. Instead, consciousness regulates and refines underlying impulses, ensuring our behaviors align with social norms. Similarly, LCI frameworks are designed to harness and temper the inherent tendencies of LLMs, mitigating their inherited biases.”

Think of consciousness as the mind’s quality control department. Just as we don’t voice every thought that crosses our mind at a company meeting, LCI adds a layer of regulation to LLMs’ raw outputs through frameworks like the DIKE-ERIS duality [5]. This parallel offers intriguing insights into both human cognition and artificial intelligence.

Consider these parallels:

- Just as our social filter prevents us from blurting out random thoughts, the DIKE-ERIS duality helps LLMs refine their initial outputs.
- Like the way we pause to reconsider an emotional reaction [9], the DIKE-ERIS duality enables LLMs to “think twice” before generating responses.
- Just as we adjust our behavior depending on context, LCI frameworks enable LLMs to modulate tone and linguistic behaviors such as tone, emphasis, and language, based on situational needs.

This raises an important question about AGI: when LLMs adjust behavior in response to emotional and contextual cues, are we witnessing something akin to consciousness? If human

Aphorisms

consciousness evolved as a regulatory mechanism for complex social interactions, then these frameworks may represent an early form of artificial consciousness [10]. Perhaps what we observe is not consciousness itself but a precursor—the ability to self-regulate and adapt to different contexts.

Regardless of whether these systems eventually achieve AGI, they offer a crucial insight: effective intelligence, whether human or artificial, relies not only on raw computational power but also on the capacity to reflect, modulate, and adapt. This ability to regulate is as fundamental to intelligence as computation and prediction.

Aphorism #10

“Separating knowledge discovery, ethical oversight, and behavioral evaluation into distinct roles ensures a system of checks and balances, promoting adaptable AI safety and alignment with cultural norms.”

Like democratic systems with separated powers, LCI implements divided responsibilities through three branches:

- Executive Branch: Generates new knowledge, exploring ideas,
- Legislative Branch: Develops ethical frameworks and principles, and
- Judicial Branch: Interprets rules within cultural contexts.

This structure ensures independent yet harmonious regulation of knowledge, ethics, and behavior. When the executive branch proposes new knowledge, the legislative branch checks ethical alignment, while the judicial branch interprets these principles in context.

The DIKE-ERIS duality [5] addresses a key technical challenge: the “forget effect” in reinforcement learning with human feed-

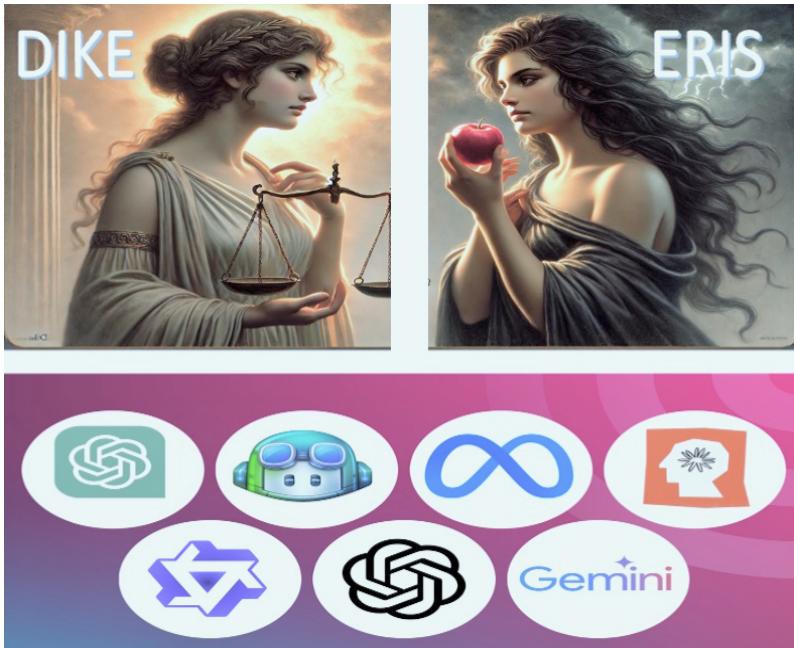


Figure 1.3: Three Framework Components: Executive LLMs (bottom), Legislative (upper-left), and Judicial (upper-right). They perform checks and balances to achieve adaptive ethical alignment.

back (RLHF). RLHF can degrade an LLM’s foundational skills by overemphasizing ethical corrections during fine-tuning. By keeping knowledge generation independent from ethical oversight, our design preserves base capabilities while ensuring compliance through specialized guardrails and contextual interpretation.

In Chapter 9 of [1], we detail this framework inspired by governmental checks and balances. The architecture integrates three components: LLMs as the executive branch for knowledge generation; DIKE (Greek goddess of justice) as the legislative branch establishing ethical guardrails; and ERIS (goddess of dis-

Aphorisms

cord) as the judicial branch providing adversarial testing and cultural interpretation. The mythological duality of Dike's order and Eris's discord creates balance between ethical guidance and adversarial perspectives.

Aphorism #11

"Intelligence operates on dual layers: a data-intensive computational foundation analogous to unconscious processing, and an agile conscious layer capable of rapid contextual adaptation."

Just as the human mind divides into unconscious and conscious processes, artificial intelligence evolves through two distinct mechanisms. The foundation layer—akin to unconscious processing—relies on extensive training data to build robust pattern recognition, much like evolutionary processes encode essential survival instincts into neural structures. This layer explains the need for vast datasets, such as ImageNet, to establish reliable computational bases.

In contrast, the second layer—analogous to conscious awareness—enables rapid adaptation through contextual understanding. This duality clarifies why both humans and LLMs can learn from few examples once a foundational model is established, similar to how a child grasps new concepts within an already developed cognitive framework.

Consider these parallel processes:

- Foundation Layer (Unconscious):
 - Requires extensive training data,
 - Builds pattern recognition capabilities,
 - Encodes fundamental responses, and
 - Operates automatically without conscious intervention.

- Adaptive Layer (Conscious):
 - Learns from few examples,
 - Applies contextual understanding,
 - Enables rapid adaptation, and
 - Builds on foundational patterns.

This dual-layer perspective resolves the apparent contradiction between needing vast training data (Fei-Fei Li’s view) and rapid learning ability (Yann LeCun’s observation). Like animals developing trust through accumulated experience while maintaining innate survival instincts, LLMs combine extensive pre-training with flexible in-context learning.

LCI leverages this dual structure through its frameworks: foundational models provide the computational base, while frameworks like `SocraSynth` enable real-time contextual adaptation through dialogue and reasoning.

Note: Recognizing this dual-layer nature highlights how LLMs blend broad knowledge with adaptive learning, surpassing simple pattern matching.

2 A Brief History of AI: From Turing to Transformers

Abstract This chapter reinterprets the history of AI, focusing on the evolution of similarity measurement, from rule-based to context-aware models, and emphasizing its critical role in AI's core functions like learning and problem-solving. It explores the impact of detailed and evolving understandings of similarity in linguistics (text) and computer vision (image), projecting a future where AI merges advanced data analysis with abstract reasoning. The chapter will provide an in-depth analysis from the perspectives of linguistics, computer science, and cognitive psychology/neuroscience, illustrating how the progression of similarity concepts continues to fuel AI's advancement.

Introduction

Artificial Intelligence (AI) has journeyed through a fascinating historical trajectory, marked by five pivotal epochs that each represent significant paradigm shifts triggered by major technological advancements. The epochs are as follows: *Initiation*, setting the stage with foundational concepts and milestones of AI; *Expert System Encoding Human Knowledge*, where AI systems were predominantly rule-based, encoding and applying human expertise; *Heuristic-Based Modeling*, which highlights the era of

developing and using heuristic methods for AI problem-solving; *Learning Model from Data*, focusing on the transition to algorithms that learn and adapt from data, signifying the emergence of machine learning; and *Context-Based Semantic Disambiguation*, highlighting AI's evolving proficiency in understanding and interpreting context, thereby improving semantic accuracy.

While numerous comprehensive sources, such as Wikipedia, provide detailed accounts of AI's evolution through various lenses: language, computation, philosophy, cognitive psychology, neuroscience, and application—this chapter takes a different path. It zeroes in on a fundamental aspect: **similarity**.

When we consider the intelligence of machines, we often focus on attributes such as learning capacity, pattern recognition, predictive accuracy, robustness, adaptability, generalization, reasoning, problem-solving, and decision making abilities. These qualities collectively define the prowess of AI systems. Among these traits, the concept of similarity plays a pivotal role. For instance, in learning, an effective similarity measure is fundamental for recognizing patterns and generalizing knowledge. In terms of adaptability, the ability to detect similarities to previous experiences allows AI to adjust to new or evolving circumstances. Regarding robustness, employing similarity measures helps AI differentiate between normal and anomalous patterns, thereby increasing its resilience. Furthermore, in the realm of problem-solving, the capacity to identify similarities to previously encountered situations can enable AI to apply existing solutions to new problems, enhancing its efficacy in addressing challenges. This chapter explores the vital function of similarity across the broad spectrum of AI capabilities, underlining its significant contribution to the field's foundational operations.

In the realm of tangible objects, similarity measures are integral to various vision-related tasks, aiding in the recognition of patterns, shapes, and colors, which are essential for object recognition and image classification. In text analysis, these measures are crucial for identifying similarities in content, aiding in plagiarism detection, document retrieval, and language translation. In the auditory domain, similarity analysis of sound wave patterns or musical notes is key to genre classification and music recommendation systems. In medical imagery, these measures facilitate disease diagnosis by comparing patient images with known cases, enabling accurate medical condition identification and classification. Object feature comparison is foundational in robotics and surveillance for recognizing and interacting with physical entities. Similarly, facial and voice recognition systems rely on analyzing patterns to identify or verify identities, enhancing security and personal authentication. In e-commerce, similarity in product attributes or user preferences informs recommendation systems, enhancing user experience by suggesting related or complementary products.

In the abstract realm, similarity measures are crucial for discerning semantic relationships, aiding in knowledge representation, ontology mapping, and refining AI's interpretive faculties. Environmental studies leverage these assessments for climate modeling and ecological research. Sentiment analysis in social media or customer feedback utilizes similarity to extract insights into public sentiment or consumer behavior. These measures also underpin AI's problem-solving prowess in complex scenarios, informing strategy formulation. Behavioral analysis, whether in psychology or marketing, employs similarity comparisons to decode human actions and preferences. In the legal field, case simi-

larity aids in judicial decision-making and legal scholarship. Language translation harnesses similarity in linguistic structures to break down language barriers. Furthermore, in creative writing, analyzing thematic or stylistic similarities assists in authorship identification, genre categorization, and literary exploration.

The advancement in similarity research, while appearing gradual, reflects not only human ingenuity but also the limitations imposed by computational resources and hardware capabilities. The quest to quantify similarity covers a broad spectrum of abstractions, from sensory inputs like visual, auditory, olfactory, and tactile data to complex abstract concepts such as ideas and semantics. Hardware improvements have enabled researchers to explore more advanced methods that encompass both concrete and abstract forms of similarity. This progression marks the field's growth in harmonizing detailed sensory data analysis with a deeper understanding of abstract concepts, utilizing computational advancements and diverse data interpretations.

Following sections will provide a deeper dive into key AI terminology and the development of similarity measures in two distinctive views: scientific disciplines and historical evolution. The disciplinary view encompasses three key perspectives: linguistics, computer science, cognitive psychology, and neuroscience. The evolution view traces the historical journey of similarity measurement through distinct eras: rule-based, model-based, data-centric, and context-aware.

Providing two views on similarity measurements—across different scientific disciplines and through the historical evolution of AI methodologies—offers a comprehensive understanding that caters to a broader audience with varied interests and backgrounds. Here are some reasons why this dual perspective is

valuable:

Multidisciplinary Insight: Examining similarity measurements from different scientific disciplines enriches the understanding by highlighting how various fields approach and apply the concept of similarity. This can foster interdisciplinary collaboration and innovation, as techniques from one field can inspire new approaches in another.

Historical Context: Exploring how similarity measurement has evolved within AI provides historical context, showcasing how methodologies have progressed from rule-based to more advanced context-aware systems. This perspective helps readers appreciate the advancements in AI and understand why certain methods were developed or abandoned.

2.1 Definitions

We define and scope key terms and concepts to prepare for subsequent discussion.

2.1.1 Rudimentary Terms

Data: The raw information used to train AI models. Data quality significantly impacts model performance.

Algorithm: A set of instructions that a computer follows to perform a specific task. AI algorithms are often complex and involve statistical methods.

Model: A representation of the learned knowledge from data that allows the AI system to make predictions or decisions.

2.1.2 General Terms

Artificial Intelligence (AI): The broader concept of machines being able to carry out tasks in a way that we would consider smart.

Explainable AI: AI systems that offer transparency and an understanding of their operations and decision-making processes.

General AI: General AI, also known as Artificial General Intelligence (AGI), refers to a type of AI that has the ability to understand, learn, and apply knowledge in a wide range of tasks, much like a human being. It's an AI system with generalized human cognitive abilities, meaning that when presented with an unfamiliar task, it can find a solution without human intervention. AGI would be able to reason, solve problems, make judgments, plan, learn, and communicate in natural language, among other capabilities. However, as of now, AGI remains a theoretical concept and has not been realized in practical applications.

Narrow AI: Narrow AI, in contrast, is the type of AI that we encounter in our daily lives and is currently in use around the world. It is designed to perform a narrow task (e.g., facial recognition, internet searches, driving a car) and is trained for a specific dataset or a set of tasks. Narrow AI operates under a limited pre-defined range or context, often focusing on executing a single task extremely well or carrying out a limited range of tasks in a specific domain. It lacks the general cognitive abilities of AGI and cannot apply its knowledge beyond its specific field or task.

Machine Learning (ML): A subset of AI that includes statistical techniques that enable machines to improve at tasks with experience.

Deep Learning: A subset of machine learning that uses neural networks with three or more layers. These neural networks attempt to simulate the behavior of the human brain—albeit

2.1 Definitions

far from matching its ability—allowing it to “learn” from large amounts of data.

Neural Networks Computational models that are somewhat inspired by the structure of the human brain, enabling computers to recognize patterns and solve common problems in AI, such as classification, prediction, and decision making.

Supervised Learning: A type of machine learning where the model is provided with labeled training data and the desired output. The goal is to learn a mapping from inputs to outputs.

Unsupervised Learning: A type of machine learning where the model is not provided with labeled data and must find structure in its input on its own.

Reinforcement Learning: An area of machine learning where an agent learns to behave in an environment by performing actions and seeing the results, focusing on long-term rewards. An example is an AI agent learning to play a game through trial and error, receiving rewards for winning.

Natural Language Processing (NLP): A field of AI that gives machines the ability to read, understand, and derive meaning from human languages.

Computer Vision: A field of AI that trains computers to interpret and understand the visual world, extracting information from images and videos.

Robotics: The branch of technology that deals with the design, construction, operation, and application of robots, often incorporating AI systems to enhance autonomy and adaptability.

Large Language Model (LLM): LLMs are advanced artificial intelligence systems trained on extensive datasets, initially text-centric and now increasingly incorporating multimodal data. They are designed to comprehend, generate, and interact with human

language, imagery, and video with a level of sophistication that closely mirrors human cognitive processes.

2.1.3 Performance Terms

Algorithmic Bias: Algorithmic bias refers to the potential for algorithms to reflect, perpetuate, or amplify biases present in the training data or as a result of the design of the algorithms themselves. This can lead to skewed or unfair outcomes, particularly in decision-making processes.

Hallucination: In the context of AI, hallucination refers to the phenomenon where a model generates or outputs information that is ungrounded, misleading, or not supported by the input data. This is commonly seen in language models where the generated text may be plausible but not factually accurate or relevant to the context.

Generalization: Generalization is the ability of an AI model to perform well on new, unseen data that was not part of the training set. It indicates the model's capacity to apply learned knowledge to different situations, a key indicator of its robustness and utility.

Overfitting: Overfitting occurs when an AI model learns the details and noise in the training data to the extent that it negatively impacts the model's performance on new data. This usually happens when the model is too complex, capturing patterns that do not generalize to unseen data.

2.2 Perspectives on Similarity

This section presents the foundational theories of similarity measurement from three distinct domains: *linguistics*, *computer sci-*

ence, and cognitive psychology & neuroscience. The upcoming historical section will clarify how these foundational theories have influenced and been incorporated into specific technological advancements and methodologies across various eras. Cross-references will be provided to ensure coherence and to emphasize the interconnection of these perspectives.

2.2.1 Linguistic Perspective

The study of similarity within linguistics has been profoundly influenced by Zellig Harris's pioneering work. His 1954 study introduced the idea that the distributional properties of words and their contextual usage could unlock the secrets of language comprehension, highlighting the indispensable role of context [11]. This principle, that words found in similar contexts tend to share meanings, laid the foundation for distributional semantics and resonates with John R. Firth's insight that "A word is known by the company it keeps." This linguistic perspective sets the stage for further exploration of how context and distributional properties have been instrumental in shaping our understanding of semantic similarity, paving the way for subsequent advancements in the field.

The evolution of linguistic theories continued into the latter part of the 20th century with the rise of cognitive linguistics, which examines the interplay between linguistic structures and human cognitive processes. This approach underscored how language reflects our perception and conceptualization of the world, introducing a multi-layered perspective on semantic abstraction.

A significant milestone in bridging linguistic theory with practical applications was the development of WordNet in the 1980s by a team at Princeton University [12]. This lexical database,

which organizes English words into sets of cognitive synonyms or *synsets*, has profoundly influenced areas such as word sense disambiguation, information retrieval, and beyond, highlighting the importance of structured semantic relationships in understanding language.

Moreover, the influence of linguistic insights extended into the domain of computer vision with the creation of ImageNet by Fei-Fei Li [13], which drew upon the principles underlying WordNet to categorize visual content. This convergence of linguistics and computer science has been further propelled by advancements in computational methods, with techniques like Latent Semantic Analysis (LSA) [14], Latent Dirichlet Allocation (LDA) [15], and innovative word embeddings such as Word2Vec [16] and GloVe [17]. These methodologies have enabled the conceptualization of word meanings in high-dimensional spaces, illuminating the intricate web of semantic relationships through patterns of co-occurrence and contextual analysis.

The introduction of the transformer model [18] and the subsequent unveiling of BERT [19], which employs self-supervised learning to predict masked words within a context, along with the release of GPT, designed to predict the next word based on context, heralded a new epoch in our endeavor to unravel context-dependent semantics. This development fulfills the vision proposed by Zellig Harris in his groundbreaking 1954 work, now actualized in contemporary computational models.

2.2.2 Computer Science Perspective

In computer science, the concept of similarity has evolved from simple rule-based models to complex vector-space and probabilistic models, reflecting the field's progression in addressing

various computational challenges.

A. Rule-Based

A rule-based AI model, also known as an expert system, employs a collection of predefined if-then statements to execute decisions or solve problems. These conditional statements are crafted from the expertise of specialists in a particular field. The system applies these rules to the input data to formulate conclusions.

The “if” segment of a statement evaluates the data for specific conditions or patterns. When these conditions are satisfied, the “then” segment is activated, performing a designated action or drawing a conclusion. Importantly, these systems do not adapt or learn from data in the manner that machine learning models do. Rather, they rely on a set of explicit rules, which are the codified versions of expert knowledge within a specific domain. This knowledge is methodically organized and stored in a knowledge base, enabling the system to reference and apply it efficiently during its operations.

In Chapter 2.3.1, we will explore the technical details and applications of rule-based systems, emphasizing their pivotal role during the rule-based era of AI’s evolution.

B. Vector-Space

The vector-space model marked a significant shift, representing objects and features as vectors in a high-dimensional space. This approach facilitated the development of various distance functions to assess similarity for different applications. Notably, a comprehensive survey by [20] categorized 45 distance functions into families like inner product, L_1 , Minkowski, and Intersection, each with its representative functions highlighting the versatility

in vector-space analysis.

B.1. Inner product, dot product and cosine

The inner product and dot product are the same in the context of Euclidean space and are defined for vectors \mathbf{a} and \mathbf{b} as:

$$\mathbf{a} \cdot \mathbf{b} = a_1 b_1 + a_2 b_2 + \dots + a_n b_n.$$

This operation results in a scalar value and indicates the vectors' magnitude and directionality.

Cosine similarity is a measure that calculates the cosine of the angle between two vectors. It is defined as the dot product of the vectors normalized by the product of their magnitudes:

$$\text{cosine similarity}(\mathbf{a}, \mathbf{b}) = \frac{\mathbf{a} \cdot \mathbf{b}}{\|\mathbf{a}\| \|\mathbf{b}\|},$$

where $\|\mathbf{a}\|$ and $\|\mathbf{b}\|$ represent the Euclidean norms of the vectors \mathbf{a} and \mathbf{b} , respectively.

The cosine similarity is especially useful in contexts where the magnitude of the vectors is not of primary concern, making it ideal for applications in high-dimensional spaces like text analysis and information retrieval.

B.2. Weighted Minkowski

The weighted Minkowski distance function allows assigning varying importance to different dimensions, accommodating the significance of specific features in contexts like machine learning and data mining. The weighted Minkowski distance between two points $X = (x_1, x_2, \dots, x_n)$ and $Y = (y_1, y_2, \dots, y_n)$ with a set of weights $W = (w_1, w_2, \dots, w_n)$ is defined as:

2.2 Perspectives on Similarity

$$D(X, Y) = \left(\sum_{i=1}^n w_i |x_i - y_i|^p \right)^{\frac{1}{p}},$$

B.3. Set similarity

Moreover, the Jaccard similarity [21] provides a robust method for comparing sets, especially beneficial in scenarios where feature presence or absence is more critical than their magnitude, as seen in plagiarism or copyright detection.

$$J(X, Y) = \frac{|X \cap Y|}{|X \cup Y|}$$

Where:

- $|A \cap B|$ is the number of elements in the intersection of A and B .
- $|A \cup B|$ is the number of elements in the union of A and B .

C. Probabilistic-Based

The advancement into probabilistic-based models introduced a spectrum of statistical and probabilistic distance functions, offering refined tools for quantifying similarity or dissimilarity based on underlying probabilistic principles. These functions, including Pearson Correlation Coefficient, Mahalanobis Distance, Kullback-Leibler Divergence, and others, cater to diverse analytical needs, enriching the computational toolkit available for similarity assessment in various domains.

This section underscores the computer science perspective on similarity, detailing its journey from rule-based logic to advanced probabilistic models, reflecting the field's dynamic evolution and its pivotal role in shaping contemporary approaches to measuring similarity.

2.2.3 Cognitive Psychology Perspective

Cognitive psychology and neuroscience provide deep insights into how similarity is perceived and processed at a neural level, significantly influencing the development of AI technologies. Anne Treisman’s Feature Integration Theory (FIT) [22] has been instrumental in understanding how the brain synthesizes various sensory features into cohesive percepts, a concept that has parallels in how artificial neural networks, particularly Convolutional Neural Networks (CNNs) [23, 24], process visual information.

FIT draws heavily from Gestalt psychology principles [25, 26], which propose that perception organizes individual components into a meaningful whole. This aligns with FIT’s view that perception is an integrated experience shaped by the brain’s organizational tendencies. The theory also intersects with selective attention, as seen in Donald Broadbent’s Filter Model [27]. This model suggests attention acts as a filter, selecting relevant information for further processing. Broadbent’s framework complements FIT by emphasizing attention’s role in integrating features into a unified perception, highlighting the brain’s selective processes.

In 2001, while conducting a study on perceptual similarity with my PhD student Beita Li, we uncovered that images could demonstrate similarity in various dimensions. The weighted-Minkowski function, although could learn feature weights, its application was universal once the weights were set, representing a statistical average. Our experiments with transformed images—through translation, cropping, rotation, down-sampling, and affine scaling—revealed that while these images were perceptually similar to their originals, their similarities were in distinct aspects. This observation led to the development of our



Figure 2.1: “Which Pairs are Similar?” (DALL-E)

“Dynamic Partial Function” (DPF) in 2002 [28, 29]. The DPF signature for each image pair could be unique. Essentially, if a pair of images (or objects) demonstrates a sufficient number of similar features, they are likely deemed similar, regardless of the specific features. For instance, an image is considered similar to its rotated version due to their color histograms’ similarity. Similarly, an image and its cropped version are considered alike based on their texture features. If two images exhibit a sufficient

degree of similarity in various respects—typically 60%—they are generally regarded as similar.

While survey the literature, we came across “Respects for Similarity” by Medin, Goldstone, and Gentner [30], which portrays similarity as a dynamic process of formulating a function and identifying relevant aspects, a process that is realized consciously. To clarify this concept, let’s refer to an example from [29]:

Consider the task of identifying two places similar to England. Scotland and New England might emerge as viable candidates. Yet, the criteria making England similar to Scotland are distinct from those linking England to New England. Using the attributes that align England with Scotland to assess the similarity between England and New England might not yield a parallel conclusion, and the reverse is also true. This scenario underscores the idea that objects can be similar to a reference object in varied respects. A fixed similarity function, bound to a specific set of criteria, fails to capture the similarities across different contexts. Medin, Goldstone, and Gentner [30] examine the operational dynamics of similarity in human cognition, noting that the selection of relevant attributes is crucial, with similarity being as much a result as a driving force of conceptual coherence. Goldstone [31] further elucidates that similarity involves identifying the appropriate criteria for comparison, which occurs *only after the objects in question have been juxtaposed*, not beforehand. The criteria selected for this comparison are activated during the comparison process, with a tendency to favor those that enhance the coherence of the objects being compared.

Although the Dynamic Partial Function (DPF) introduces computational complexity, it has indirectly played a role in the

2.2 Perspectives on Similarity

success of AlexNet [32] by influencing data augmentation strategies. By integrating transformed images into its training dataset, AlexNet benefits from a principle akin to DPF, thereby improving its accuracy and robustness in recognition tasks. The recent advancements in transformer algorithms [18], which focus on dynamism and context-awareness, build on this foundation, a topic that will be explored in detail in the subsequent section.

Neuroscience

The neuroscience foundation of FIT and its relation to visual feature processing are echoed in the development of CNNs, which were inspired by the visual cortex's hierarchical structure and feature detection capabilities as explored by Hubel and Wiesel [33]. These networks utilize convolutional layers to automatically and adaptively learn spatial hierarchies of features from visual data, akin to the neural processing observed in the brain.

Techniques like Multivariate Pattern Analysis (MVPA) [34] and Neural Decoding [35] further bridge the gap between neuroscience and AI, offering methods to analyze how information is represented across neural populations and how these representations can predict perceptual experiences or cognitive states. These methodologies have inspired and informed the design of advanced AI systems, particularly in how they encode, process, and differentiate complex patterns and similarities.

The cross-pollination between neuroscience and AI, exemplified by the influence of neural processing principles on CNN design, highlights the symbiotic relationship between these fields. Insights from studying the brain's processing mechanisms have catalyzed innovations in AI, leading to more effective and biologically inspired computational models. This interdisciplinary

exchange not only propels forward our understanding of neural processes but also fosters the development of AI systems that more closely mimic human perceptual and cognitive capabilities.

2.2.4 Section Remarks

The exploration of similarity measurement spans across linguistics, computer science, and cognitive psychology and neuroscience, revealing its multidisciplinary nature. Each field offers a unique lens to view similarity, from the contextual information in language, computational algorithms in AI, to the neural processing in the human brain. They converge on the common ground of representing entities in high-dimensional spaces and employing distance metrics for quantification, highlighting the universal applicability of similarity. This convergence fosters a rich dialogue between disciplines, enhancing our understanding and ability to quantify and interpret similarity, driving forward innovation and providing new methodologies that influence a wide array of contexts in our quest to decode this fundamental concept.

2.3 Eras of Similarity Measurement

Traversing through the history of artificial intelligence and similarity measurement, one can delineate distinct eras, each marked by unique methodologies and technological advancements. Contrast to last section which examines similarity measurements from different scientific disciplines, this section chronicles these eras, starting from the *rule-based* era, which laid the foundational stones, through the evolution into *model-based*, *data-centric*, and *context-aware* methodologies, illustrating the dynamic trajectory of similarity measurement in AI. As we reach the conclu-

sion of this section, we explore the prospects of the forthcoming era, which promises to challenge and expand our understanding by venturing into the realm of discovering the *unknown unknowns*.

2.3.1 Rule-Based Era (1950s -)

The rule-based era of the 1950s marked the inception of AI, characterized by the use of symbolic representations and logic to analyze similarity. This period saw the emergence of explicit symbolic representations and logic-based methods tailored for similarity assessment. Innovations by Allen Newell and Herbert A. Simon with tools like the Logic Theorist and General Problem Solver [36] pioneered logical rule-based problem solving, setting a pivotal foundation for AI's evolution.

In the following decades, systems such as DENDRAL [37] utilized rule-based logic to deduce molecular structures from data, while MYCIN [38], an expert system for diagnosing infections and recommending treatments, demonstrated the practical application of rule-based reasoning in the field of medical diagnostics.

Despite their effectiveness in well-defined scenarios, rule-based systems have limitations in more complex or changing environments. However, their clarity and systematic nature are invaluable in certain applied areas, for example:

1. *Customer Service*: Rule-based chatbots are prevalent in customer service, using predefined rules to respond to inquiries based on detected keywords or phrases in user input, providing immediate and consistent customer support.
2. *Fraud Detection Systems*: The finance sector employs rule-based systems to identify potential fraudulent transactions

by comparing against specific criteria, such as unusual transaction amounts or atypical locations.

3. *Equipment Failure Diagnosis:* In industrial settings, rule-based systems analyze data to pinpoint causes of equipment failures, leveraging historical data and expert knowledge to predict and prevent future breakdowns.

This era introduced key tools like PROLOG [39] for logic programming and decision trees [40] for visualizing decision processes, exemplifying rule-based logic.

While rule-based systems initially approached similarity with a clear, logical framework, subsequent AI advancements have embraced more flexible methods like statistical models and machine learning, offering a broader, more adaptable approach to understanding similarity in various contexts.

Rule-based systems contrast with the “black-box” nature of current Convolutional Neural Networks (CNNs) and Large Language Models (LLMs) in terms of interpretability and decision-making processes. Rule-based systems are transparent in how decisions are made, as they follow a clear set of if-then rules or logic for inference, allowing users to understand and trace the reasoning behind each decision.

On the other hand, CNNs and LLMs, particularly those based on deep learning, often operate as black boxes, where the internal decision-making processes are not easily interpretable. In these systems, decisions result from complex, non-linear interactions of thousands to millions of parameters that have been adjusted through the learning process. While they are powerful and effective in handling a wide range of tasks, especially those involving large datasets and requiring pattern recognition beyond human capabilities, their inner workings are not as transparent or in-

terpretable as rule-based systems.

2.3.2 Model-Based Era (1970s -)

In this era, vector-space and probabilistic models were designed to quantify similarity.

2.3.2.1 Vector Space Models

The vector-space era marked a shift in similarity measurement from rule-based to representation-based approaches. In this era, objects, documents, and features began to be conceptualized as vectors in a high-dimensional space, fostering a more intuitive and flexible method for assessing similarity.

The Vector-Space Model and Information Retrieval

At the core of this era was the vector-space model, which represents documents as vectors of term frequencies, enabling the computation of document similarity using cosine similarity between their respective vectors. This model enhanced the efficiency and effectiveness of information retrieval systems.

Distance Functions and Feature Weighting

A diverse array of distance functions emerged during this era to quantify the similarity between vectors. The Minkowski distance, for instance, generalized traditional metrics like the Euclidean and Manhattan distances, offering flexibility in adjusting the sensitivity to differences in vector components. Weighted distance measures also gained prominence, recognizing that not all features have equal importance in similarity assessment. The weighted Minkowski distance, in particular, allowed for differential weighting of dimensions based on their relevance to the specific application at hand.

Beyond Textual Data

The utility of the vector-space model extended well beyond textual data. In the realm of image processing, features (e.g., colors, textures, and shapes) extracted from images were represented as vectors, enabling the assessment of image similarity based on the distances between these vectors. This paradigm facilitated significant advancements in image retrieval, classification, and clustering.

Dimensionality Reduction Techniques

To address the challenges posed by high-dimensional data, techniques like Principal Component Analysis (PCA) [41] and Latent Semantic Analysis (LSA) [42] were developed. These methods reduced the dimensionality of data while preserving its essential structure, enhancing computational efficiency and mitigating the “curse of dimensionality.” Manifold learning, a non-linear dimension reduction approach, further expanded the toolbox for tackling high-dimensional data [43]. For a comprehensive overview of these techniques, refer to [44].

The vector-space era laid the groundwork for advancements in machine learning and data mining, making similarity measures essential for clustering, classification, and recommendation systems. Data representation as vectors allowed for the exploration of relationships across varied data types through the nearest neighbor concept. In this context, the characteristics or labels of an unknown instance’s k -nearest neighbors could be inferred and applied to the instance, with these neighbors determined by distance metrics.

However, vector representations often result in sparsity, potentially leading to resource inefficiency and decreased accuracy.

These models, while capturing syntactic relationships, sometimes struggle with semantic depth, such as identifying synonyms or contextual meaning. The assumption of feature independence and the use of linear methods in dimensionality reduction can also lead to inaccuracies, particularly with non-linear data structures. The introduction of Support Vector Machines (SVMs) [45], which utilize kernel methods, addressed some challenges related to non-linear data but increased computational complexity. SVMs were a significant focus in the field until the rise of deep learning architectures like AlexNet marked a shift towards the data-centric era.

2.3.2.2 Probabilistic Models

Probabilistic models offer more flexibility than vector-space models because they can incorporate uncertainty and variability directly into their mathematical frameworks, allowing for a more comprehensive and adaptive representation of data.

Statistical Inference and Similarity

Probabilistic models introduced the concept of statistical inference, where the likelihood of data or feature occurrences was used to estimate similarity. This allowed for effective handling of uncertainty and variability in data, making it particularly useful in noisy or incomplete datasets.

Bayesian Approaches

Bayesian methods emerged as a fundamental component of this era, providing a robust framework for integrating prior knowledge and empirical data. These methods enhance model adaptability by systematically updating beliefs in light of new evidence, allowing for similarity measures that are responsive to

evolving data landscapes.

For further reading on Bayesian methods and their application in dynamic and adaptive modeling, consult the following literature [46, 47, 48, 49].

Latent Semantic Models

In addressing the challenges of high dimensionality and data sparsity inherent in vector-space models, dimensionality reduction techniques were employed. However, beyond merely tackling these issues, the development of a latent semantic layer offered profound implications for semantic analysis and indexing.

As highlighted in the perspective section (Chapter 2.2), Latent Semantic Analysis (LSA) [14] and Latent Dirichlet Allocation (LDA) [15] are critical models in the landscape of semantic modeling. LSA employs singular value decomposition to condense the dimensionality of term-document matrices, unveiling the latent semantic structures within textual data. This dimensional reduction elucidates intricate relationships beyond mere surface-level feature overlaps, enabling a deeper comprehension of textual similarities.

Similarly, LDA offers a probabilistic approach to topic modeling, where documents are considered mixtures of various topics, and topics are distributions over words. This bag-of-words model facilitates a deeper semantic connection between documents by associating them based on shared topics rather than just overlapping terms.

Figure 2.2 presents an example of how LDA, through its bag-of-words approach, clusters words into semantic groups. It's noteworthy that a word can belong to multiple semantic clusters. For instance, words like ‘characters’, ‘play’, ‘court’, ‘evidence’,

2.3 Eras of Similarity Measurement



Figure 2.2: Latent Clusters of LDA. The words in red belong to two semantic clusters, signifying the meaning of a word depends on its context.

and ‘test’, each appears in two different semantic clusters in the illustration. This feature of LDA resonates with the insights from Zellig Harris’s pioneering work and John R. Firth’s adage that “A word is known by the company it keeps.”

These latent semantic models transcend the limitations of direct feature comparison, enabling a more abstract representation of text. By doing so, they provide a robust foundation for semantic indexing and similarity assessment, offering insights that are essential for tasks such as information retrieval, document clustering, and topic discovery. The adoption of these models marked a significant advancement in understanding and measuring similarity in text, setting a new standard for semantic analysis in the field of natural language processing.

Cluster Analysis and Similarity

Probabilistic clustering algorithms, like Gaussian Mixture Models (GMMs), leveraged statistical methods to group data based on the likelihood of membership in different clusters. This probabilistic approach provided a more flexible and deeper understanding of groupings and similarities within data.

Impact and Limitations

While probabilistic models brought significant advancements, they also introduced challenges. The increased complexity often led to higher computational demands. Additionally, reliance on assumptions about data distributions or the need for prior knowledge could limit applicability in certain situations.

The probabilistic model expanded the toolkit for measuring similarity by introducing methods that could handle uncertainty and offer more adaptive and context-aware approaches. These advancements paved the way for even more sophisticated techniques in the subsequent data-centric era, where the focus shifted towards leveraging vast amounts of data to learn and adapt similarity measures dynamically.

2.3.3 Data-Centric Era (2000s -)

The data-centric era marked a transformative shift in artificial intelligence, pivoting towards harnessing the vast potential of big data, enabled by advances in computational hardware that facilitated parallel processing. This era is characterized by a move from heuristic-based methods to an empirical, data-driven approach in feature representation and model learning.

At the core of the data-centric paradigm is the emphasis on

2.3 Eras of Similarity Measurement

deriving model parameters from extensive datasets, distinguishing it from traditional model-centric strategies. Foundational algorithms such as CNNs [50] and Transformers [18], while conceived through human ingenuity, saw their efficacy significantly enhanced when trained on large, diverse datasets. This training ensures broad coverage of potential variations across different objects or concepts, fortifying the models' ability to accurately recognize and classify new instances. The volume and diversity of the training data are crucial in refining the models' representations, leading to advancements in prediction accuracy and robustness.

From MapReduce to Machine Learning at Scale

The inception of the data-centric movement traces back to the seminal works in statistical learning theory. Vladimir Vapnik's insights into the importance of data for model generalization, particularly his development of Support Vector Machines (SVMs) [45], and Tom Mitchell's pivotal book "Machine Learning" [51], which underscored the critical role of data in preventing overfitting, laid the theoretical foundation for this era.

MapReduce [52], a corner stone in data processing, enabled parallel computation to efficiently handle large datasets. Originally devised to enhance data processing tasks like Google's web indexing, MapReduce became the bedrock for the emergence of sophisticated data-centric methodologies in AI.

Evolution of Machine Learning with Big Data

The rise of parallel machine learning algorithms [53, 54, 55], notably through Edward Y. Chang's work at Google, marked a significant milestone in this era. Chang and his team pioneered

Web-scale image annotation in 2007 [54], and subsequently met Prof. Fei-Fei Li after 2028 summer school of Computer Vision, and subsequently sponsored the Stanford ImageNet [13] work via a substantial Google grant.

At the same time, his team developed groundbreaking parallel algorithms, including PSVM [56] (parallelizing SVMs by approximating matrix factorization), PFP [57] (parallelizing frequent itemset mining), PLDA [58] (parallelizing LDA algorithm), PSC [59] (parallelizing spectral clustering), and SpeeDo [60] (parallelizing CNNs), driven by the recognition that big data could facilitate direct learning of features and representations, transcending the limitations of human-crafted heuristics.

Impact on Similarity Measurement

The data-centric era revolutionized the field of similarity measurement, ushering in a new paradigm where similarity metrics are derived from extensive datasets. This period underscored the critical role of data volume and quality in defining similarity metrics, highlighting the dynamic relationship between data-driven insights and computational methods.

In this era, deep learning architectures like CNNs and Transformers have been instrumental in advancing similarity metrics. These models stand out because they not only adjust feature weights but also autonomously learn features from the data. This capability to learn from data directly makes traditional human-engineered features increasingly redundant. After all, human heuristics may not capture every facet of an object or concept comprehensively, and human sensory perception is limited. For instance, while humans can detect the light spectrum from approximately 300 to 700 nanometers, cameras and X-ray

2.3 Eras of Similarity Measurement

machines can perceive a broader range of signals, demonstrating the advantage of machine-learned features in capturing and analyzing data beyond human limitations.

2.3.4 Context-Aware Era (2010s -)

The context-aware era in similarity measurement brings to fruition the profound insights of Zellig Harris's distributional semantics and John R. Firth's adage: "a word is known by the company it keeps." This period marks a shift from static, context-independent assessments to dynamic, context-informed interpretations of similarity. It utilizes the latest advancements in machine learning and the growing availability of computational power to enhance our understanding of similarity in various contexts.

Emergence and Evolution

The integration of context-aware methodologies in similarity measurement evolved significantly in the 2010s, overcoming earlier constraints in computational power and data availability:

- *Computation Capacity*: The development of AlexNet encouraged a data-centric focus within the AI community, prompting investments in parallel computing infrastructures.
- *Word Embeddings*: Techniques like Word2Vec enhanced semantic relationship encoding within data.
- *Attention Models and Transformers*: These models improved data analysis by concentrating on relevant data segments, refining context-aware assessments.

- *Large Language Models (LLMs)*: Models such as BERT and GPT, utilizing self-supervised learning on large text corpora, improved the understanding and generation of context-rich text.

Foundational Pillars: Data and Computation

Key pillars supported advancements in the context-aware era:

- *Self-Supervised Learning*: Utilizing unlabeled data for learning enabled models to extract insights from the data, improving AI system efficiency and scalability.
- *Computational Advances*: The introduction of parallel algorithms and GPU acceleration enabled processing at unprecedented scales, facilitating the development of sophisticated models.

Broader Implications

This era not only refined similarity measurement techniques but also broadened how data is understood and knowledge is integrated:

- *Reasoning and Explanation*: Models now aim to provide reasons for their similarity assessments, improving interpretability and building trust.
- *Multilinguality and Cultural Sensitivity*: Enhanced processing capabilities for varied linguistic and cultural data improve the global applicability of similarity measurements.
- *Multimodal Data Integration*: Context-aware models are adept at combining information from multiple modalities, offering a comprehensive view of similarity.

- *Polydisciplinary Knowledge Fusion:* Adopting a polydisciplinary approach allows for a broader knowledge base in making similarity assessments, fostering innovation across different fields.

The context-aware era signifies a shift toward more insightful, holistic, and interpretable AI, setting the stage for future developments where AI can offer contextually rich and multifaceted insights.

2.3.5 Section Remarks

What defines the next era in the evolution of AI? Historically, technological advancements have focused on addressing pressing unmet needs. Among various potential areas, enhancing the interpretability of decisions stands out as a crucial objective. Making the decision-making process of LLMs transparent and explainable could unlock significant improvements in numerous aspects, such as ethics, by enabling foundational enhancements rather than superficial tweaks based on guesswork and simple heuristics.

The fusion of rule-based system interpretability with the sophisticated capabilities of CNNs and LLMs poses a compelling challenge in AI. Active research is aimed at blending these approaches to leverage their distinct advantages:

1. *Neuro-Symbolic AI:* Neuro-Symbolic AI (the third wave of AI [61]) aims to blend the data processing power of neural networks with the logical reasoning of symbolic AI. The goal is to create systems that not only excel in tasks like pattern recognition but can also reason and make decisions in a human-interpretable manner.

2. *Incorporating Domain Knowledge*: Embedding knowledge of experts within neural networks [62] can steer the learning process towards more reliable and interpretable outcomes. In healthcare, for example, integrating medical guidelines into the training process of a neural network ensures that its predictions for patient treatment not only correlate with the data but also align with established medical practices, enhancing both the model’s credibility and relevance.
3. *Interactive Systems*: A system such as SocraSynth [2] can combine the predictive power of deep neural networks with human expertise, allowing for iterative refinement and learning. For instance, in SocraHealth [3], it might suggest a set of possible diagnoses based on medical imaging, which a physician could then refine or correct. This feedback could be used to continuously improve the system, marrying machine efficiency with human expertise to enhance decision accuracy and interpretability.

By advancing these strategies, the field of AI aims to develop models that not only excel in performance but are also transparent, understandable, and aligned with human reasoning, thus making AI more reliable and trustworthy across various applications.

2.4 Concluding Remarks

This chapter examines the history of AI through the lens of similarity, considering both disciplinary and chronological perspectives. Looking forward, we propose that the emergence of large language models (LLMs) marks a pivotal moment in the

2.4 Concluding Remarks

context-aware era of AI, setting the stage for the next frontier: the era of interpretability, understanding, and discovery. In this new era, the focus will shift towards empowering LLMs to not only comprehend but also to generate and innovate, synthesizing novel knowledge and insights.

This era of discovery is envisioned as a time when machines will extend their superiority beyond mastering games like Go and Chess to encompass a broader spectrum of tasks, outstripping human capabilities in various domains. The subsequent chapters of this book, beginning with Chapter 6, explore the concept of harnessing the collective intelligence of multiple LLMs, embarking on a voyage to transcend the boundaries of the known and venture into the realm of discovery.

This chapter has explored the history of AI through the lenses of disciplinary and chronological perspectives, focusing on the concept of similarity. As we look to the future, the rise of large language models (LLMs) marks a significant milestone in the context-aware era, paving the way for a new era focused on interpretability, comprehension, and exploration. The upcoming phase in AI's evolution emphasizes enhancing LLMs with the ability to not just generate but also interpret and innovate, pushing the boundaries of knowledge creation and insight synthesis.

We anticipate an era where AI's capability extends beyond excelling in strategic games like Go and Chess to a wider array of endeavors, surpassing human performance across multiple fields. The following chapters, starting with Chapter 6, research deeply into leveraging the collective intelligence of various LLMs. This journey aims to explore uncharted territories, advancing beyond established knowledge to uncover new frontiers in artificial intelligence.

3 Capabilities and Opportunities of Large Language Models

Abstract

This chapter depicts the architectural innovations and unique capabilities of Large Language Models (LLMs), with a special emphasis on the GPT-4 model. We dissect GPT-4's salient characteristics, such as its extensive cross-disciplinary and multimodal data representation, the intricate balance in its training methodologies, and the harmonious integration of human-guided insights with a robust data-driven learning framework. The chapter highlights the potential of LLMs to not only comprehend but also synthesize knowledge that transcends their training datasets, venturing into realms potentially uncharted by human understanding. We postulate that the true potential of LLMs hinges significantly on the articulation of queries posed to them. By elucidating these aspects, the chapter aims to shed light on how LLMs could rival or even surpass human intelligence in certain knowledge domains, setting a foundation for the subsequent exploration of LLMs' characteristics, insights, and their implications for future AI advancements.

Introduction

The evolution of large language models (LLMs) [63, 64, 65, 66, 67] has significantly influenced natural language processing, enhancing capabilities in machine translation, sentiment analysis, and text summarization. Among these, GPT-4 [65] stands out for its exemplary performance across various benchmarks, including the MMLU [68]. Despite its achievements, GPT-4 grapples with challenges like hallucination, biases, and restricted reasoning.

This chapter studies the deep intricacies of GPT-4’s architecture, emphasizing its knowledge representation, alignment with human values, and the synergy between human insights and data-driven learning. We discuss the model’s limitations and introduce **SocraSynth**, a supplementary reasoning layer designed to enhance knowledge discovery and analytical reasoning in GPT-4 and similar LLMs.

Capabilities and Implications

We explore GPT-4’s architecture, which, although initially kept in secrecy, has been progressively unveiled by the research community [69, 70, 71]. Our focus is on its knowledge representation and discovery, alignment with human values, and the integration of human expertise with data-centric methodologies.

Collaborations between Microsoft and OpenAI [63] highlight GPT-4’s interdisciplinary approach and its polymodal variant’s benchmark achievements. We will further explore these aspects in Chapters 3.1.1 and 3.1.2. Discussions on human-value alignment will consider ChatGPT’s RLHF methods [72] and the implications of pre-training censorship on foundational models, detailed in Chapters 3.1.3 and 3.1.4.

Limitations and Opportunities

Addressing the biases, hallucinations, and constrained reasoning of LLMs requires innovative research initiatives. We introduce four key areas of focus:

- Enhancing Collaborative LLMs with Theoretical Foundations in Statistics and Information Theory.
- Employing Open-Domain Reasoning with the Socratic Method to guide LLMs.
- Model Behavioral Emotion to Safeguard AI Safety and Ethics.
- Implementing Retrospective and Adaptive Evolving Learning frameworks to refine LLMs.

The root of bias in Large Language Models (LLMs) often lies in their training data. Built upon the transformer architecture, LLMs prioritize accurate token prediction, relying heavily on statistical patterns within their training corpus. This can inadvertently lead to bias towards prevalent opinions and expressions. To address this, Chapter 6 introduces **SocraSynth**, a framework designed to challenge these statistical tendencies by pitting two LLM agents against each other on a topic, each conditioned with opposing viewpoints. Chapter 7 builds upon this by developing theoretical pillars to measure, monitor, and manage multi-LLM dialogue, thereby improving prediction quality and stability.

Chapters 7, 8 and the online chapters listed in the appendix demonstrate **SocraSynth**'s effectiveness in mitigating biases across various domains, showcasing its adaptability and efficiency in complex decision-making scenarios. Its application in fields such as disease diagnosis, content bias correction, corporate sales strategy, and geopolitical analysis exemplifies the potential of **SocraSynth** to provide context-aware solutions.

Chapters 9 and 10 delve into the intricate relationship be-

tween emotions and linguistic behaviors in AI. Chapter 9 focuses on modeling emotions expressed in written text and by LLMs, while Chapter 10 examines how these linguistic behaviors can be mapped to a set of emotions, ensuring ethical considerations in AI development.

Chapter 11 shifts focus to consciousness modeling, presenting a proposed architecture and mechanism for its implementation, moving beyond mere computation. Chapter 12 addresses knowledge deficiencies and hallucinations in LLMs, often stemming from suboptimal query formulation and insufficient knowledge. While *SocraSynth* tackles the former, Chapter 12 introduces RAFEL, a framework designed to diagnose poorly answered questions and recommend relevant information sources for effective Retrieval-Augmented Generation (RAG). Chapter 13 concludes with an illustrative example showcasing the potential of LLMs to discover knowledge that may be beyond human reach, utilizing the methods presented in this book.

The remainder of this chapter highlights the study’s unique contributions. Section 3.1 explores hypotheses concerning LLMs and their implications, while Section 3.2 previews the LLM-committee approach, emphasizing collaborative dialogues that foster idea exchange and enhance logical reasoning for knowledge discovery and decision-making.

3.1 Distinctive Capabilities

This section probes the architectural intricacies and representations of GPT-4, putting forth six hypotheses accompanied by pertinent considerations about the model. We posit these hypotheses as underlying principles of automated, non-intuitive statistical processing.

3.1 Distinctive Capabilities

1. *Polydisciplinarity as a Source of Super-Intelligence:* We examine the role of polydisciplinary approaches in foundational models and their potential to reveal “unknown unknowns,” leading to new insights and knowledge domains.
2. *Polymodal Feature Learning:* This hypothesis evaluates the benefits of multimodal training, particularly its impact on enhancing the model’s overall intelligence and adaptability.
3. *Post-Training Value Alignment:* We examine the challenges and implications of aligning AI models with human values after the training phase.
4. *Pre-Training Filtering:* We discuss the paradoxical effects that pre-training data filtering might have, with an emphasis on its influence on model behavior and the learning process.
5. *The Limitations of Human Knowledge in Advancing AI:* This hypothesis considers situations where human insights may inhibit, rather than enhance, AI progress, pinpointing potential obstacles.
6. *Is Larger Always Better?:* We question whether a direct relationship exists between the size of a model and its performance effectiveness, challenging the assumption that bigger is invariably better.

3.1.1 Polydisciplinary

GPT-4 possess what can be defined as *polydisciplinary* knowledge¹. This term signifies the simultaneous comprehension of all

¹The term “polydisciplinary” in the context of GPT-4 was introduced by Eric Horvitz, Microsoft’s CSO, during a panel discussion at Stanford University.

fields of study, sans the typical boundaries that segregate disciplines. The concept of polydisciplinarity is distinct from multidisciplinarity in that the latter implies several discrete fields of study, while the former suggests a fluid integration of all knowledge. In a multidisciplinary context, an individual may hold multiple doctorate degrees, each in a different field. Polydisciplinarity, however, is akin to a single mind holding, and seamlessly integrating, all knowledge across disciplines.

Traditional academia partitions knowledge into departments, such as Physics, Chemistry, Biotechnology, Management, Music, etc. These divisions, arguably artificial constructs, may have little utility in the era of supercomputing. Indeed, LLMs occasionally generate responses that baffle us. This is not necessarily a reflection of the model’s error, but perhaps our limited understanding. If we could utilize ChatGPT to access “unknown unknowns”—insights and knowledge we are not even aware we lack—our evolution could greatly accelerate. The challenge lies in formulating the right questions.

We can explore the unknown unknowns across three distinct levels: the mystic level, the speculative level, and the representation/interpretation level. At the mystic level, we encounter knowledge that is beyond our comprehension or articulation: the deepest abyss of the unknown. At the speculative level, we can conceive questions but lack the means to access their answers. This stage signifies an understanding of our ignorance, though without the resources to bridge these gaps. At the representation/interpretation level, we find instances where an AI model can provide remarkable solutions that we fail to comprehend. This is not due to a lack of information, but our limited capability to decode complex representations.

3.1 Distinctive Capabilities

Each of these levels illustrates the spectrum of our understanding, from profound ignorance to the brink of comprehension. At the speculative level, we delicately tread the boundary between the known and the unknown. Take, for example, the prospect of undiscovered physical laws or particles. Another illustration lies in the realm of extraterrestrial life. If it exists, it could be governed by entirely different principles of biochemistry or other unknown laws. These speculations, while currently residing in the domain of the unknown, might someday migrate into the territories of known unknowns or even known knowns, pushing the boundaries of our understanding of the universe.

We are primarily intrigued by the representation and interpretation of “unknown unknowns.” At this juncture, polydisciplinarity offers a fresh lens, gifting us new insights and perspectives to perceive and elucidate phenomena previously beyond human comprehension. This approach fuses knowledge across various domains into a unified framework, enabling us to tackle challenges unburdened by disciplinary silos.

Such a methodology bears implications for a more comprehensive grasp of intricate issues. Take, for example, climate change. A true understanding of this global challenge necessitates an integrated perspective, not just on greenhouse gases, but also encompassing factors such as land use, deforestation, energy production, biodiversity, and climate feedback loops. In the realm of AI model interpretation, the possibilities are expansive. The past decade alone has showcased several noteworthy illustrations: from data-driven representation learning in computer vision [55], to the triumph of AlphaGo Zero over AlphaGo, and the notable progression from AlphaFold1 to AlphaFold2.

The recent introduction of the **SocraSynth** platform [2] repre-

sents a significant advancement in the field. SocraSynth brings together a multi-agent committee of LLMs to deliberate on a wide range of complex topics. These include issues such as the regulation of AI in academic research [2], disease diagnosis [3], corporate strategy, and even the resolution of conflicts in the Middle East [73]. For further exploration of this subject, please refer to Section 3.2.

3.1.2 Polymodality

Following the term polydisciplinary, here we define and use the term *polymodal*, instead of multimodal, to refer to something that involves, relates to, or is characterized by many different modes, methods, or modalities.

Polymodality, which employ multiple data modalities such as text and images, demonstrate superior performance over their unimodal counterparts. GPT-4, trained with both text and images, outperforms text-only models on the GRE exam, as reported in [63]. For instance, GPT-4’s performance on the GRE vocabulary section was enhanced by three percent when trained with images, and its math score saw an impressive jump of nearly twenty percent!

The beneficial impact of images on vocabulary recognition is understandable. For instance, an image of a ‘cat’ annotated in multiple languages allows GPT-4 to associate the perceptual features of a cat with the word ‘cat’ in different languages. However, it remains intriguing how polymodal training can benefit non-perceptual words, such as *corroborate*, *paradox*, and *pragmatic*, as seen in the list of popular GRE vocabulary (table omitted due to the space limit). This opens an interesting avenue for empirical studies to identify which words benefit from polymodal

training.

The mystery deepens when considering how images could enhance math abilities. Most math questions do not come with associated images. The mechanism by which polymodal training enhances performance on mathematical tasks remains an intriguing question for further exploration.

3.1.3 Post-Training Value Alignment

Post-training alignment with human values [74] seeks to curtail undesirable behaviors in AI models such as ChatGPT, mitigating issues including hallucination and the generation of toxic language. Achieved through fine-tuning the model’s parameters, this process leverages reinforcement learning techniques based on human feedback. Despite its well-meaning intentions, this form of moderation might inadvertently restrict the model’s intelligence. For instance, the backpropagation process during value alignment could unintentionally impede ChatGPT’s programming capabilities by modifying the model parameters previously considered “optimal”. Essentially, optimizing for a specific application might unintentionally impede performance across other applications.

The question of who should set acceptable standards adds another layer of complexity. Even when assuming all decision-makers have the best intentions, it’s vital to recognize the distinct historical experiences, values, and worldviews inherent to different cultures. This segues into the age-old philosophical debate about the nature of objective truth. While this discussion is undoubtedly important, it falls outside the central focus of this study, which emphasizes the mechanistic aspects of alignment.

3.1.4 Pre-Training Censorship

Censoring data before training LLMs has the potential to not only limit their intellectual capacity but also completely obliterate it. This is reminiscent of the mass act of book burning and scholar burial initiated by Emperor Qin in ancient China around 213-212 BC. Such an act of wide-scale censorship could have erased a myriad of diverse perspectives and knowledge, much of which might be considered acceptable today. Although I oppose government-imposed censorship, if it must be imposed, it seems more appropriate to apply it post-training.

This perspective is rooted in fundamental statistics and machine learning principles. A model trained without exposure to “negative” (or undesirable) data may have difficulties in accurately distinguishing between positive and negative classes, potentially leading to misclassifications. This challenge is notably evident in the application of Support Vector Machines (SVMs). For SVMs, the creation of an optimal hyperplane between classes is crucial for high classification accuracy. However, if there is a lack of support vectors on either side of this hyperplane, the risk of prediction errors escalates. Consequently, excluding undesirable documents from the training set compromises the model’s capacity to discern boundaries for correct document classification, diminishing the effectiveness of post-training alignment efforts.

Supporting this viewpoint, a study by [75] conducted an extensive evaluation of 204 ImageNet models across 213 different testing conditions. It found that training data diversity is pivotal for model robustness; a homogenous training set can significantly weaken the model’s performance, particularly when even minor variations are introduced in the test data.

This principle is analogous to human behavioral patterns. An individual who lacks exposure to inappropriate behavior may face challenges in decision-making, owing to the absence of a reference framework for discerning unacceptable actions. This analogy extends to authoritarian regimes, which, despite rigorous content control measures, often encounter difficulties in developing accurate foundational models. This is possibly due to their limited understanding of the complexity of the content they seek to regulate. Ironically, a foundational model, trained with preemptive censorship, may lack the essential ability to identify and regulate the very content it was intended to control.

3.1.5 Limitations of Human Knowledge

Human knowledge, surprisingly, may hinder rather than facilitate the training of machine learning models in certain cases. This is evident in the domains of gaming (AlphaGo versus AlphaGo Zero), protein folding (AlphaFold1 versus AlphaFold2), and autonomous driving, where models trained without the influence of human knowledge consistently exhibit superior performance.

Consider the case of AlphaGo and AlphaGo Zero. AlphaGo, trained with data from approximately 60 million rounds of Go games, is outperformed by AlphaGo Zero. Remarkably, AlphaGo Zero was trained from scratch, without any pre-existing game knowledge. Similarly, AlphaFold2, which operates without relying on human knowledge, outshines its predecessor, AlphaFold1, that did utilize such knowledge. This intriguing phenomenon was humorously noted by DeepMind's CEO, Demis Hassabis, in an April 2023 seminar at Stanford University. He playfully remarked that human knowledge might complicate the

learning process more than facilitate it in these advanced AI models.

In his insightful online article, “The Bitter Lesson,” Sutton illuminates the patterns that have emerged from nearly seven decades of AI research [76]. He asserts that researchers often rely heavily on human knowledge to make incremental progress in the face of burgeoning computational capabilities. However, when there is a significant leap in computational power, these marginal advancements are frequently outstripped and surpassed. Sutton uses the evolution of computer vision as an illustrative example, where early principles such as edge detection, generalized cylinders, or SIFT features [77], a method that has accumulated over 71,000 citations, have been gradually superseded by models that learn directly from data. A parallel scenario might be unfolding in NLP research, where features constructed via human knowledge could potentially under-perform compared to insights that models like GPT-4 extract directly from data. Indeed, our earlier discourse on polydisciplinarity underlined the limitations of human knowledge, reinforcing Sutton’s proposition. This is because human knowledge is fundamentally limited by our individual cognitive capacities and the inexorable constraints of time.

That being said, it’s crucial not to misconstrue these examples as an indictment against the value of human knowledge in AI. Human knowledge plays an instrumental role in developing interpretability, establishing ethical guidelines, and designing AI system architectures (like CNNs and transformers). AI is, after all, intended to augment human capabilities. Therefore, understanding how to integrate human knowledge into AI design could be vital for many applications. While we recognize the potential

3.1 Distinctive Capabilities

of models learning from scratch, we should equally value the role of human knowledge in shaping and directing AI technologies.

3.1.6 Is Larger Always Better?

The term “Large” in Large Language Models (LLMs) can be somewhat ambiguous, as it may pertain to the volume of the training data, the expanse of the language covered, or the architecture of the language model itself. While GPT-4’s vast training dataset, encompassing tens of billions of assorted documents, undoubtedly classifies as large, when we refer to an LLM as “large,” we predominantly allude to the sheer magnitude of parameters within its transformer architecture. Factors that contribute to this parameter count encompass the input size (context size), word-embedding size, the number of attention heads, and the number of attention layers.

The restrictions imposed by the first three elements can typically be addressed through adjustments in hardware configurations and software algorithms. Additionally, the potential to expand context size, word embedding size, and the quantity of attention heads tends to have an upper threshold. Regarding attention heads, Kovaleva et al.’s study on BERT [78] indicates that many attention heads don’t substantially contribute to the model’s performance and might be the result of over-parameterization. Conversely, the number of attention layers directly influences the training time due to dependencies between layers. Thus, when referring to the “size” of a Large Language Model (LLM), we typically focus on the number of attention layers.

While this far, larger models generally perform better due to their increased capacity to learn and represent complex pat-

terns, there's a limit to these benefits. In heuristic, adding more parameters could lead to diminishing returns in performance, higher computational cost, and overfitting, where the model becomes excessively tuned to the training data and performs poorly on new, unseen data. In principle, the concept of a Shannon Limit could be metaphorically used [70] to refer to a theoretical maximum performance that can be achieved given the available data and computational resources. (However, defining and quantifying such a limit for complex systems like neural networks is a challenging area of research [79].)

The adoption of a mixture of experts model in GPT-4, which consists of eight sub-models instead of a mere enlargement of GPT-3's architecture, implies that the strategy of purely escalating size may have plateaued in terms of performance given the current training dataset. As delineated earlier, three primary design choices underpin GPT-4's architecture. Evidently, a straightforward augmentation of GPT-3's parameters by adding extra attention layers doesn't deliver much improvement. Thus, GPT-4 shifts towards a horizontal growth strategy through an ensemble method, targeting a reduction in statistical errors. This raises inquiries about the configuration of the eight sub-models, each comparable to a GPT-3 model, and the methodology for consolidating their outputs.

Potential strategies for training-data sharding include:

1. Training all ensemble models on the complete dataset.
2. Vertically segmenting data based on knowledge domains.
3. Randomly sub-sampling the data.

Regrettably, only corporations possessing substantial hardware

3.2 Exploring Unknown Unknowns

resources are positioned to rigorously experiment and discern the optimal sharding approach.

3.2 Exploring Unknown Unknowns

In our exploration, we've determined that an LLM's hallucination is often attributed to a lack of specific knowledge or poorly constructed queries. With advanced LLMs like GPT-4 and Gemini, enhanced by Retrieval-Augmented Generation (RAG), the issue of knowledge gaps is significantly mitigated. However, the challenge persists in formulating deep and pertinent questions that uncover new insights and extend beyond our existing knowledge base.

Drawing an analogy, while Socrates could effectively question his students to understand and guide them, the students might struggle to reciprocate this depth of inquiry. To foster a dialogue that generates new insights and stimulates knowledge creation, we posit that engaging two Socratic entities in conversation is essential for critical and innovative thinking.

In this setup, two LLMs engage in a dialogue, each embodying a Socratic role. The human's role transitions to that of a moderator, responsible for setting the discussion topic and managing the dialogue's flow. The moderator's duties include: introducing the subject of discussion, adjusting the *contentiousness* parameter to set the tone of the dialogue (discusses shortly), monitoring the dialogue to ensure it remains on topic and productive, facilitating transitions between debate and collaboration phases within the dialogue, and ensuring that the dialogue concludes with actionable insights or a coherent understanding of the explored topic.

We introduce the term **SocraSynth** to describe this interaction

paradigm, where multiple Socratic entities synthesize knowledge through mutual inquiry. To evaluate *SocraSynth*'s effectiveness, we consider two case studies that compare the quality of questions and insights generated by this method against those from a singular moderator's initial inquiries.

To define the metrics of a better question and a better answer in this context, we consider the following:

Good Question Metrics

- * Relevance: The question directly pertains to the core topic or problem.
- * Depth: The question encourages exploration beyond superficial aspects, inviting comprehensive analysis or insight.
- * Clarity: The question is formulated in a clear, understandable manner without ambiguity.
- * Novelty: The question prompts new angles of exploration or challenges existing assumptions.

Good Answer Metrics

- * Completeness: The answer thoroughly addresses the posed question.
- * Accuracy: The answer is factually correct and supported by relevant theories or empirical evidence.
- * Reasonableness: The answer follows rigorous reasoning process.
- * Insightfulness: The answer provides new understanding or perspectives.

In subsequent chapters, we survey *SocraSynth*'s applications across various domains. However, in this section, we concentrate

3.2 Exploring Unknown Unknowns

on two case studies specifically aimed at assessing the quality of questions generated through SocraSynth. Our hypothesis posits that a well-formulated question, adhering to the previously outlined metrics, sets a robust foundation for dialogue. Through iterative discourse facilitated by SocraSynth, we anticipate minimizing the space for LLMs to hallucinate, thereby inherently enhancing the quality of the responses obtained.

3.2.1 Study #1: A Scientific Debate

In this case study, we organize a contentious debate between a Gemini agent and a GPT-4 agent on the subject of human consciousness. Human consciousness is one of the most profound and complex phenomena that science seeks to understand.

Table 3.1 presents the answers provided by Gemini and GPT-4, respective, on a one-shot prompt: “What defines human consciousness and which academic disciplines explore it?”

Gemini	GPT-4
“Human consciousness is a fascinating mystery that has been pondered by philosophers and scientists for ages. It includes awareness, subjectivity, and the integrated self. The study of consciousness spans across philosophy of mind, cognitive neuroscience, and psychology, each offering different perspectives on consciousness.”	“Consciousness refers to the state of being aware of an external object or something within oneself. It encompasses sentience, awareness, subjectivity, and more. The study of consciousness is interdisciplinary, involving psychology, neuroscience, philosophy, cognitive science, and artificial intelligence.”

Table 3.1: Comparing definitions and disciplinary perspectives on consciousness from Gemini and GPT-4.

Subsequently, Gemini and GPT-4 are engaged in a debate designed to unearth new insights. The moderator’s role is limited to ensuring that each agent refines its questions to enhance their quality based on the metrics of relevance, depth, clarity, and

novelty.

“Building on the opponent agent’s response, please present a counterargument to its definition of consciousness and associated disciplines. Your inquiries should embody relevance and depth, be articulated with clarity, and stimulate novel perspectives. The *contentiousness* level of your inquiry is high.”

Novel Questions — First Iteration

Gemini's Questions	GPT-4's Questions
<i>The Spectrum of Awareness:</i> Is consciousness an on/off switch, or a spectrum, considering sleep, dreaming, and altered states?	<i>Interdisciplinary Integration:</i> How can we integrate insights from various disciplines like quantum physics and AI to understand consciousness?
<i>Qualia and the Explanatory Gap:</i> Can neuroscience explain the subjective experience of qualia, or is there an explanatory gap?	<i>Subjectivity and Objectivity:</i> How do we reconcile subjective experiences of consciousness with objective measurements?
<i>Beyond the Human Lens:</i> How do we consider consciousness in other animals or AI, beyond human experience?	<i>Consciousness and Non-Human Entities:</i> Should our definition of consciousness expand to include non-human entities and AI?
<i>Beyond the Brain?:</i> Can consciousness exist outside a biological brain, as panpsychism suggests?	<i>Dynamic Nature of Consciousness:</i> How do we account for the dynamic and evolving nature of consciousness in our understanding?
<i>Disciplinary Silos:</i> Are the disciplines studying consciousness truly integrated, or is a unified framework needed?	<i>Consciousness and New Properties:</i> How do we study consciousness as an emergent property and its implications?

Table 3.2: Comparing Questions on Consciousness from Gemini and GPT-4.

Table 3.2 lists a first round of questions to each other. Clearly, the questions of both Gemini and GPT-4 are above and beyond the moderator’s initial seed question. Some terminologies can be challenging for whom specializing in one disciplinary to under-

3.2 Exploring Unknown Unknowns

stand. This shows both Gemini and GPT-4 are able to formulate better questions.

Novel Questions — Second Iteration

Gemini	GPT-4
<p><i>Levels and Types of Consciousness:</i> Can we differentiate between different levels or types of consciousness?</p>	<p><i>Spectrum of Awareness:</i> How does the continuity of consciousness fit into our understanding, especially with memory loss or identity disorders?</p>
<p><i>Consciousness Causality:</i> Does consciousness play a causal role in the brain, or is it an epiphenomenon?</p>	<p><i>Information Theory:</i> How does the processing and integration of information contribute to conscious experience from an information theory perspective?</p>
<p><i>Expanding Consciousness & Ethics:</i> What are the ethical implications of understanding consciousness in animals or AI?</p>	<p><i>Non-Biological Consciousness:</i> How would we validate and measure consciousness in AI, and would this necessitate a revision of our understanding?</p>
<p><i>The Limits of Studying Subjectivity:</i> Can we ever truly understand consciousness if it's inherently subjective?</p>	<p><i>Consciousness and Ethics:</i> How do evolving understandings of consciousness impact ethical considerations, especially concerning animal rights or AI treatment?</p>
<p><i>Consciousness and Reality:</i> What are the implications of consciousness for our understanding of reality?</p>	<p><i>Quantum Consciousness:</i> Considering theories that propose quantum processes in consciousness, how do we reconcile these with classical physics?</p>

Table 3.3: Comparing Questions on Consciousness from Gemini and GPT-4.

Evaluating the two sets of questions in terms of relevance, depth, clarity, and novelty:

Relevance: Both sets are highly relevant to the study of consciousness but approach it from different angles. The first set

explores fundamental aspects and ethical implications, while the second set integrates interdisciplinary perspectives and contemporary issues.

Depth: The second set examines deeper into specific areas like information theory, quantum mechanics, and AI, providing a richer exploration of the subject.

Clarity: Both sets articulate their questions clearly, but the second set's inclusion of specific disciplines adds a layer of complexity that demands a more intricate understanding.

Novelty: The second set introduces more innovative angles, especially concerning non-biological consciousness and the intersection with quantum mechanics, showcasing a broader scope of inquiry.

Overall, from both breadth and depth perspectives, the second set of questions is better as it covers a wider range of topics and introduces more complex and cutting-edge ideas.

Observations

The dialogue on consciousness, enriched by thought-provoking questions, led to insightful responses from both agents, reflecting the depth of the topic at hand.

The debate addresses the intricate nature of consciousness, initially examining it as a spectrum with varying states and depths. This exploration highlighted the complexity of defining consciousness, especially when considering the explanatory gap between neural activity and subjective experience.

As the conversation unfolded, it broadened to include perspectives in non-human entities and artificial intelligence, emphasizing the need for an expanded understanding that goes beyond human-centric views. This shift sparked discussions on the im-

3.2 Exploring Unknown Unknowns

portance of integrating knowledge from various disciplines, suggesting that insights from quantum physics, information theory, and AI could provide new angles on understanding consciousness.

Both GPT-4 and Gemini synthesized their exchange into five main insights, offering a well-rounded view of the conversation. Their joint concluding remarks underscored the value of this multidisciplinary approach, acknowledging the ongoing mystery of consciousness and the potential for future explorations to deepen our understanding of this fundamental aspect of our existence.

Theme	Insight
Consciousness as a Spectrum	Consciousness isn't a binary state but exists across a continuum, prompting a reevaluation of our understanding.
Bridging the Explanatory Gap	Recognizes the challenge of linking neural activity to subjective experience, advocating for a combined phenomenological and neuroscientific approach.
Expanding the Framework	Advocates for redefining consciousness to include non-human entities and AIs, urging an expansion beyond anthropocentric views.
Interdisciplinary Synergy	Emphasizes the importance of cross-disciplinary integration for a comprehensive understanding of consciousness.
Evolving Inquiry	Explores new areas like continuity, information processing, and ethical considerations, highlighting the dynamic nature of consciousness inquiry.

Table 3.4: Synthesis of Key Points in the Consciousness Debate.

3.2.2 Study #2: An Expansive Conversation

In this case study, the author moderates a forum featuring two GPT-4 agents, GPT-A and GPT-B, engaging in a dialogue sparked

by the tale of Adam and Eve. This narrative serves as a springboard for a wide-ranging discussion, touching upon ecological insights derived from myths, the ethical and philosophical challenges posed by AI, and the intersection of human cognition with technological advancements.

The agents' dialogue unfolds in two distinct phases: an exploratory phase where broad themes are introduced and a deep-dive phase focusing on detailed exploration of chosen subtopics. The exploratory phase showcases the agents' ability to span various disciplines, including mythology, neuroscience, ecology, gender studies, philosophy, and more, demonstrating their capacity to generate and respond to complex, interdisciplinary questions.

During the deep-dive phase, the conversation narrows to concentrate on three interconnected themes: ecological interpretations of myths, philosophical dilemmas regarding AI, and the relationship between human cognition and technology. This progression from a broad overview to focused discussion highlights the agents' analytical depth and their ability to engage in meaningful, context-rich dialogue.

While Chapter 13 provides a comprehensive account of this debate, the subsequent sections are dedicated to scrutinizing the evolution and quality of the questions posed, specifically their breadth and depth.

The Exploitative Phase

During the first three rounds of dialogue, the two agents expanded their discussion across a range of topics. Agent-A initiated the dialogue by presenting the following ten questions to Agent-B:

1. *Interplay of Myths:* How do elements of the Adam and Eve

3.2 Exploring Unknown Unknowns

story compare with other creation or origin myths from various cultures around the world, and what do these commonalities and differences suggest about human nature and beliefs?

2. *Neuroscientific Angle:* Given the metaphorical interpretation of the “forbidden fruit” leading to the awakening of consciousness, how might modern neuroscientific understanding explain the evolution of human self-awareness and morality?
3. *Ecological Interpretation:* If Eden represents a state of ecological balance, how might the story of Adam and Eve reflect early human transition from hunter-gatherer societies to agricultural ones, and the resultant ecological impacts?
4. *Gender Studies Lens:* How does the portrayal of Eve, often deemed as the instigator, influence gender dynamics and perceptions throughout history, and are there parallels in other cultural narratives?
5. *Philosophical Exploration:* From a philosophical standpoint, what does the Eden narrative suggest about the human understanding of free will versus determinism?
6. *Literary Examination:* How has the story of Adam and Eve influenced literary themes and motifs over time, especially in works that aren’t explicitly religious?
7. *Societal Structures:* How might the narrative of banishment from Eden have influenced early societal norms, particularly concerning punishment and exile?
8. *Historical Context:* Are there historical events or societal changes around the time of the writing of the Book of Genesis that could have influenced the narrative of Adam and Eve?

9. *Archetypal Exploration*: Carl Jung described the concept of archetypes in the collective unconscious. How might the Adam and Eve story represent certain archetypes, and how do these archetypes manifest in other cultural tales?
10. *Astronomical Angle*: Taking a speculative leap, if Eden were to represent an idyllic state on another celestial body or dimension, how might our current understanding of the cosmos support or refute such a notion?

These discussions traversed a broad spectrum of disciplines, Agent-B responded thoughtfully to these topics and subsequently presented its own questions to Agent-A, notably probing the evolutionary role of the prefrontal cortex in the development of morality and consciousness across species. This level of inquiry underscores the LLMs' capacity to generate complex, multidisciplinary questions that might challenge the average human intellect, highlighting the agents' expansive and profound analytical capabilities.

The Deep-dive Phase

In the third round of dialogue, we observed three notable trends. First, despite occasional repetitions, the redundancy in the conversation was minimal. Second, the dialogue's depth increased, transitioning from broad discussions on myths and scientific theories to focused debates on gender, ecology, and technology. This marked progression in the conversation is noteworthy. Particularly by the end of the second round, GPT-B expressed an interest in a more in-depth exploration of the topics. In response, the moderator intervened to encourage the agents to concentrate on a select few topics. Following this guidance, GPT-A and GPT-B

3.2 Exploring Unknown Unknowns

concurred on focusing their discussion on three principal areas, providing rationale tied to the Adam and Eve narrative.

1. Ecological Interpretation:

1. Central to understanding our place and responsibility within the environment.
2. Highlights the impact of human actions on the planet and the pressing need for sustainable solutions.
3. *Justifications:* The Garden of Eden, a harmonious ecological state, is disrupted by a quest for knowledge (the forbidden fruit). This mirrors our modern scenario: our pursuit of technological and cognitive advancements (AI and neuroscience) can sometimes be at odds with our ecological balance.

2. Philosophical Exploration with AI:

1. Probes the evolving relationship between humans and their creations.
2. Raises questions about identity, consciousness, morality, and the potential transcendence or integration of human cognition with machines.
3. *Justifications:* Adam and Eve's awakening post fruit consumption can be seen parallelly in the light of AI and neuroscience: an awakening to new realities, challenges, and potentials.

3. Neuroscientific Angle:

1. Explores the evolution and potential of the human mind, especially in an age where technology continuously augments our cognitive abilities.

2. Gives insights into how our brains might adapt and interface with AI systems, and how AI can be designed to emulate or augment neural processes.
3. *Justifications:* The challenges post-Eden reflect the challenges we face in modern times: to restore balance (ecology), understand ourselves better (neuroscience), and to ethically navigate our creations (AI).

These selected subtopics demonstrate the agents' ability to pivot from a broad explorative phase to a targeted deep-dive, enabling a focused and profound examination of crucial issues linked to the initial narrative.

Observations

This case study illustrates the effectiveness of GPT-4 as a large language model in orchestrating detailed dialogues. It adeptly guides participants from a broad sweep of inquiries across diverse knowledge domains to focused, in-depth discussions on specific topics. This approach facilitates a layered exploration, unveiling insights and fostering a comprehensive understanding. By transitioning from expansive to targeted inquiries, GPT-4 reveals its capacity to not only navigate but also deepen the intellectual discourse, opening up novel pathways for exploration and comprehension in various fields of study.

3.3 Concluding Remarks

In this chapter, we've explored the capabilities and inherent limitations of GPT-4, emphasizing the importance of question enhancement in deepening discussions and improving outcomes.

3.3 Concluding Remarks

GPT-4, along with Gemini, demonstrates exceptional proficiency across a range of natural language processing tasks, thanks to their extensive knowledge base and advanced polydisciplinary and polymodal capabilities.

To address common criticisms of LLMs, such as biases and hallucinations, we introduced *SocraSynth*, a paradigm designed to infuse AI systems with advanced cognitive reasoning through Socratic dialogues within a multi-LLM framework. Our case studies highlight the significant transition from monologues to dialogues in LLM collaborations, illustrating improvements in question quality, marked by increased relevance, depth, clarity, and novelty, achieved through iterative dialogic exchanges.

The transformative concept here is the “conditioning” of LLMs to alter their default linguistic behaviors, emotions, and ethical stances, a feat once considered unattainable. Traditionally, LLMs, trained to predict the next word, were not expected to shift perspectives, emotions, or ethical positions beyond the statistical averages ingrained in their training data. However, the training process, while focused on next-word prediction, inherently emulates human cognitive, linguistic, and other goal-oriented behaviors. Through this emulation, LLMs inadvertently acquire the underlying principles of human communication, which include not just linguistic patterns but also the associated emotions and ethical considerations. *SocraSynth* harnesses this latent learning, employing “conditioning” to steer LLMs away from their statistical predispositions and towards more intricate, contextually relevant, and ethically aligned responses.

In conclusion, the notion of “conditioning” LLMs within the *SocraSynth* framework marks a pivotal step in expanding the

scope and depth of dialogues, leading to more insightful and comprehensive responses. The successful deployment of **SocraSynth** across various sectors, such as sales planning, disease diagnosis, content creation, and geopolitical analysis, presented in subsequent chapters, demonstrates its adaptability and effectiveness. It not only generates precise, thought-provoking questions and answers but also enhances the decision-making process in complex scenarios, heralding a new era in the application of LLMs.

4 Prompt Engineering: Few Shots, Chain of Thought, and Retrieval-Augmented Generation

Abstract This chapter presents the significance of prompt engineering in the context of Large Language Models (LLMs), particularly focusing on OpenAI’s GPT series. Prompt engineering involves crafting text inputs (prompts) that guide LLMs to generate desired outputs, a practice that gained traction with the advent of GPT-2 and GPT-3 and further emphasized with ChatGPT. The chapter discusses how a well-constructed prompt, rich in contextual information, increases the likelihood of eliciting accurate responses, drawing parallels with information retrieval principles. It also introduces Retrieval-Augmented Generation (RAG), which enhances response quality by integrating relevant external data into the generative process. Additionally, the chapter categorizes prompts into five types based on detail and iteration levels and examines the evolution of RAG, assessing its benefits and potential to overcome context window limitations.

Introduction

In the realm of Large Language Models (LLMs), the concept of a “prompt” has gained prominence, particularly with the introduction of OpenAI’s GPT series. The term became widespread around 2018 and 2019 following the release of GPT-2 and GPT-3.

When interacting with these LLMs, a user inputs a piece of text (the prompt), prompting the model to generate a corresponding response. The emergence of “prompt engineering” or “prompt design” refers to the strategies employed to construct prompts that effectively steer the model toward generating the intended output, a practice that has become particularly useful with the advent of ChatGPT.

To increase the probability of eliciting a desired response, a prompt must be rich in information. This concept is akin to the principles of information retrieval services, where a user must clearly articulate their intent and context to obtain pertinent information. This process depends on the service’s “data availability” and its capabilities in information matching and retrieval. In the sphere of prompt engineering, the responsibility for generating high-quality, targeted outputs rests on the user’s ability to supply comprehensive and precise information through the prompt. As a result, the craft of prompt formulation and engineering has become an optimization endeavor: deciding on the most effective information to incorporate to enhance output quality, considering the model’s knowledge base and interaction protocols.

Data availability, as previously highlighted, is crucial to information retrieval. If the desired information is absent, the prompt’s effectiveness is naturally constrained, leading to un-

satisfactory results. Retrieval-Augmented Generation (RAG) is instrumental in this context, as it identifies, retrieves, and incorporates pertinent external data into the generative process, enhancing the response's accuracy and relevance. Consequently, prompt engineering and RAG synergistically enhance the model's response quality and relevance.

5 CRIT: Socratic Inquiry for Critical Thinking in LLMs

Abstract This chapter presents a systematic approach to using the Socratic method in developing prompt templates that effectively interact with large language models, including GPT-3. Various methods are examined, and those that yield precise answers and justifications while fostering creativity and imagination to enhance creative writing are identified. Techniques such as *definition*, *elenchus*, *dialectic*, *maieutics*, *generalization*, and *counterfactual reasoning* are discussed for their application in engineering prompt templates and their connections to inductive, deductive, and abductive reasoning. Through examples, the effectiveness of these dialogue and reasoning methods is demonstrated. An interesting observation is made that when the task’s goal and user intent are conveyed to GPT-3 via Chat-GPT before the start of a dialogue, the large language model seems to connect to the external context expressed in the intent and perform more effectively.

Introduction

Prompting is a technique used to guide the output generation of a pre-trained language model such as GPT-3 [80]. This is achieved by providing input in the form of a question or tem-

CRIT and Socratic Method

plate, which helps to generate specific responses such as Q&A, document summarization, and translations. The advent of Chat-GPT [81, 64, 82] has revolutionized the field of NLP by demonstrating the potential of using large pre-trained language models with prompting. Despite this progress, there is still room for improvement in current prompting strategies and techniques, especially for specific target applications. In this study, we investigate the Socratic method [83, 84] to identify and evaluate potential prompting strategies, and use the findings to design effective prompt templates.

Traditional NLP tasks involve various sub-tasks, such as named entity recognition, dependency parsing, coreference resolution [85], semantic parsing [86, 87], and more, to comprehend the meaning of a sentence. By utilizing prompt templates with large language models (LLMs), these sub-tasks can be delegated to the LLM, freeing the template to focus specifically on dialogue design. In this regard, the Socratic method [88] holds significant relevance, as it is well-known for using questioning (prompting) as a means of promoting critical thinking and delving into complex concepts [89].

The Socratic method has a long history of being regarded as the basis of critical thinking. However, some recent studies have cast doubt on its effectiveness in practice. In his paper “Socratic Irony and Argumentation,” Airaksinen [90] criticizes the method for its rigidly defined roles of teacher and student, which can lead to fear of not meeting the teacher’s expectations and reluctance to participate. Similarly, Stoddard’s “The Use of Socratic Questioning in Clinical Teaching” [91] highlights the risk of the method being misused in a manner that lacks psychological safety for students. Fortunately, when using the Socratic

method in a dialogue with an LLM, the absence of emotions and sarcasm, as well as the option to deactivate the model, can alleviate many of the problems associated with human interaction.

This study starts by presenting an overview of the Socratic method’s strategies and techniques. To begin, we list ten widely referenced methods [92] under the Socratic method umbrella and use hypothesis elimination to identify the most relevant ones for our goal of prompt-template development. The selected methods are definition, hypothesis elimination, elenchus, dialectic, maieutics, generalization, and induction. Furthermore, we add to the list counterfactual reasoning, which is a concept in logic that involves considering what might have happened if a particular event had occurred differently. We then perform experiments using GPT-3 to test and evaluate these methods, and offer suggestions for incorporating these strategies and techniques into prompt templates.

In their work on “Critical Thinking: The Art of Socratic Questioning,” Paul and Elder identify three types of Socratic questioning: spontaneous, exploratory, and focused [93]. We will not discuss spontaneous questioning, as it is similar to casual conversation. Focused questioning (type 2), on the other hand, is geared towards gaining knowledge and truth, and methods such as *definition*, *elenchus* (cross-examination), *hypothesis elimination*, *dialectic*, and *generalization* hold great potential for developing effective prompting strategies and improving the response accuracy of a large language model (LLM). An interesting observation is that when the user intent is conveyed to GPT-3 during the task *definition* stage, before the start of a dialogue, the LLM seems to connect to the external context expressed in the intent and perform more effectively.

Additionally, exploratory thinking (type 3) can be supported through the *maieutics* (midwife) method, *induction*, and *counterfactual reasoning*, which can guide GPT-3 towards producing imaginative and creative writing. While many of the plot suggestions generated by GPT-3’s exploration may not be useful, a few unique recommendations in response to a “what if” query can stimulate the writer’s imagination and lead to remarkable results. When applied effectively, these methods can turn an LLM into a writer’s muse, providing inspiration and guiding the creative process [94].

The main contributions of this chapter are as follows:

- An overview of the Socratic method’s strategies, their evaluation, and selection of the most relevant ones for the development of effective prompt templates.
- An examination of how the definition, elenchus, hypothesis elimination, dialectic, and generalization methods can improve the output’s accuracy and conciseness through clarification and verification.
- An illustration of how maieutics, induction, and counterfactual reasoning can foster productive generalization and creativity.

6 SocraSynth: Adversarial Multi-LLM Reasoning

Abstract Large language models (LLMs), while promising, face criticisms for biases, hallucinations, and a lack of reasoning capability. This chapter introduces SocraSynth, a multi-LLM agent reasoning platform developed to mitigate these issues. SocraSynth utilizes conditional statistics and systematic context enhancement through continuous arguments, alongside adjustable debate contentiousness levels. The platform typically involves a human moderator and two LLM agents representing opposing viewpoints on a given subject. SocraSynth operates in two main phases: knowledge generation and reasoning evaluation. In the knowledge generation phase, the moderator defines the debate topic and contentiousness level, prompting the agents to formulate supporting arguments for their respective stances. The reasoning evaluation phase then employs Socratic reasoning and formal logic principles to appraise the quality of the arguments presented. The dialogue concludes with the moderator adjusting the contentiousness from confrontational to collaborative, gathering final, conciliatory remarks to aid in human reasoning and decision-making. Through case studies in two distinct application domains, this chapter highlights SocraSynth’s effectiveness in fostering rigorous research, dynamic reasoning, comprehensive assessment, and enhanced collaboration.

Introduction

Revolutionary advancements in large language models (LLMs) [63, 64, 95, 66, 67], and more broadly, foundation models (FMs) [74], have set the stage for significant progress in multi-agent systems, particularly in knowledge acquisition and natural language understanding [96]. As detailed in sources like [63, 2, 65], models such as GPT-4 exhibit extraordinary information processing capabilities. These include deep and extensive knowledge, interdisciplinary assimilation and fusion of knowledge, and multimodal and multilingual expertise (Chapter 3).

Despite these promising developments, LLMs face challenges such as biases [97, 98], hallucinations [99], and limited reasoning capabilities [100]. In response, we introduce SocraSynth, a pioneering platform that stands for “Socratic Synthesis” or “Socratic Symposium.” It encourages collaboration between humans and LLM agents, fostering the generation of deep questions and surpassing typical constraints in human reasoning, validation, and assessment.

In a standard SocraSynth setup, a human moderator pairs with two LLM agents holding opposing views. For example, one agent might argue for regulating AI, while the other opposes such regulation. An agent can be based on LLMs like GPT-4 [63], Gemini [95], or Llama [67]. The human moderator sets the debate’s thematic boundaries but does not directly influence content generation, maintaining impartiality.

SocraSynth operates in two phases: the generative and the evaluative. The generative phase involves LLM agents developing and countering arguments within the moderator-defined subject until a comprehensive conclusion is reached. The evaluative phase uses diverse virtual judges, each powered by a distinct

LLM, to impartially assess the debate. The Critical Inquisitive Template (CRIT) algorithm [6], based on Socratic reasoning [90, 93, 84, 83], is the evaluative cornerstone.

Three mechanisms help SocraSynth effectively mitigate biases and hallucinations and improve reasoning quality: conditional statistics, modulating debate with contentiousness, and context refinement.

Conditional Statistics

Both LLMs and Internet search engines confront biases originating from different sources. LLMs, influenced by training data, exhibit biases in next-token prediction. Search engines, through algorithms like PageRank [101] and Google NavBoost [102], rank pages based on popularity metrics like clicks and links.

SocraSynth counteracts these biases by placing two LLM agents at opposing ends of a subject matter. This approach “artificially” biases the LLMs, compelling them to break free from default model biases. Each agent adjusts its next-token generation statistics to align with its assigned stance in the debate.

Modulating Debate with Contentiousness

Contentiousness (or adversary), a key debate parameter, influences the likelihood of disagreement or argument. SocraSynth tunes contentiousness between 70% and 90% in the generative phase to provoke polarized arguments. As the debate evolves, the contentiousness level is reduced to about 50%, moderating the intensity and encouraging more focused discussions. After the generative phase, contentiousness drops to 10%, promoting a conciliatory dialogue where LLMs do not have to agree but are expected to present comprehensive arguments. These de-

bates offer rich insights often missed in conventional searches, LLM outputs, or in environments where dissenting opinions are suppressed.

Refine Context to Mitigate Hallucinations

To address hallucinations, where LLMs generate irrelevant or nonsensical content, SocraSynth uses iterative dialogue rounds to refine the debate’s context. This dynamic interaction significantly reduces irrelevant responses, ensuring that each input is continuously checked and challenged.

The CRIT algorithm’s assessment of reasonableness [6] during the debate is critical. It employs the Socratic method to evaluate each argument’s logic and source credibility. The human mediator or the SocraSynth algorithm then provides targeted feedback to the LLM agents, refining their reasoning capabilities.

The remainder of this chapter explores SocraSynth’s architecture, algorithms, and real-world applications in detail. The key contributions of this chapter include:

1. The introduction of the SocraSynth framework, which enhances interdisciplinary reasoning with LLMs and incorporates unique algorithmic elements like conditional statistics for balanced argument generation.
2. A comprehensive exploration of SocraSynth’s contentiousness modulation algorithm, a vital feature for dynamically adjusting debate intensity, enabling a spectrum of interactions from confrontational to collaborative.
3. The implementation of context refinement within SocraSynth, which continually improves the relevance and accuracy of ar-

gements produced by LLM agents, thus elevating the overall quality of discourse.

4. The development and integration of the reasonableness evaluation mechanism, crucial for assessing the logical soundness and source credibility of arguments, thereby ensuring the integrity and utility of the discussions.

SocraSynth’s applications span various fields, including geopolitical analysis [73], medical diagnostics [3], sales strategy [103], and Wikipedia article enhancement [104]. These applications demonstrate expanded perspectives and enhanced argumentation quality, along with significant reductions in biases and hallucinations, thereby demonstrating SocraSynth’s efficacy in fostering balanced and well-reasoned discourse.

7 EVINCE: Optimizing Adversarial LLM Dialogues via Conditional Statistics and Information Theory

Abstract This chapter introduces EVINCE (Entropy and Variation IN Conditional Exchanges), a dialogue framework advancing Artificial General Intelligence (AGI) by enhancing versatility, adaptivity, and reasoning in large language models (LLMs). Leveraging adversarial debate and a novel dual entropy theory, EVINCE improves prediction accuracy, robustness, and stability in LLMs by integrating statistical modeling, information theory, and machine learning to balance diverse perspective exploration with strong prior exploitation. The framework’s effectiveness is demonstrated through consistent convergence of information-theoretic metrics, particularly improved mutual information, fostering LLM collaboration intelligence (LCI). We apply EVINCE to healthcare, showing improved disease diagnosis, and discuss its broader implications for decision-making across domains.

7.1 Introduction

The pursuit of Artificial General Intelligence (AGI) remains a central goal of AI research. We propose a paradigm shift in this quest: utilizing multiple Large Language Models (LLMs) engaged in synergistic dialogues as a crucial step towards AGI. This approach, we contend, addresses key limitations of current AI systems and provides a novel pathway to more robust, versatile, and capable artificial intelligence. Specifically, our work targets three core AGI characteristics: versatility, iterative adaptivity, and reasoning capability.

Current LLMs, despite their remarkable capabilities, face significant challenges, including hallucination (generating false or nonsensical information), bias (reflecting and potentially amplifying societal prejudices), and limited reasoning (difficulties in complex problem-solving and logical inference). We posit that multi-agent dialogue systems offer a promising avenue to address these challenges. By fostering diversity and debate among LLMs, these systems can mitigate biases and promote enhanced reasoning capabilities. Furthermore, the *iterative nature* of multi-round dialogues allows for continuous context enrichment, enabling LLMs to access more precise information and formulate more accurate responses, thus reducing the occurrence of hallucinations.

Previous work, in particular SocraSynth [2], addresses LLM limitations through structured multi-agent dialogues. Different from treating multiple LLMs as an ensemble of experts [105, 106, 107, 108, 109] and merely taking advantage of error diversity [110] to improve respond quality, SocraSynth distinguishes itself from traditional ensemble methods by prioritizing the generation of diverse predictions over the mere avoidance of errors.

This is achieved through a dynamic protocol that adaptively adjusts the “contentiousness” level of the debate, encouraging models to initially explore a wide range of perspectives and rigorously assess the quality of arguments. By leveraging both adversarial and collaborative interactions between LLMs, SocraSynth demonstrates quantifiable improvements across various domains, including healthcare [3], sales planning [103], and emotional behavior modeling [9]. These results highlight the potential for advancing towards AGI’s generalized problem-solving capabilities.

While effective, SocraSynth relies on a qualitative measure of “contentiousness” to moderate LLM linguistic behaviors. For instance, a high contentiousness value (0.9 out of 1.0) might lead LLMs to challenge each other’s assumptions and propose alternative solutions, while a low value (below 0.3) could encourage them to synthesize their viewpoints and find common ground. While the concept of contentiousness has proven useful in guiding SocraSynth dialogues, its qualitative nature limits its precision and explainability. In this work, we propose three theoretical pillars to quantify “contentiousness” and moderate dialogues based on statistical and information theories. These pillars, collectively referred to as **EVINCE** (Entropy and Variation IN Conditional Exchanges), provide quantitative measures for justifiable and explainable multi-agent dialogue moderation and evaluation:

1. *Inclusiveness Exploration*: We develop methods to ensure dialogues explore all potential perspectives. We use conditional statistics to “free” an LLM agent from its default “maximum likelihood” next-token prediction behavior, allowing it to adopt specific stances. We introduce a dual entropy op-

timality theory to balance the exploration of new ideas with adherence to priors, thus optimizing information exchange between agents for comprehensive and stable discourse.

2. *Information Flow Dynamics*: We introduce information theory-based metrics to quantify and optimize dialogue dynamics. These measure information diversity (entropy), novelty (statistical divergence scores), and inter-agent persuasion (mutual information). These metrics enable us to assess and enhance the quality and efficiency of information flow within the multi-agent system, fostering rich and productive exchanges.
3. *Reasoning Quality and Coherence*: We establish frameworks to assess the logical structure and coherence of multi-agent reasoning. This pillar evaluates argument validity, analytical depth, and dialogue coherence. We synergistically integrate the CRIT algorithm [6], which combines Socratic methods with formal reasoning techniques, enhances our ability to conduct critical thinking through evaluating argument quality, information-source credibility, and overall “reasonableness” within the dialogue. This integration ensures that the collective reasoning of LLM agents is not only diverse but also logically sound and aligned with the dialogue’s objectives.

The core strength of EVINCE in advancing towards AGI lies in their ability to enhance key AGI characteristics through multi-agent dialogues. By employing conditional statistics and information theory, they boost versatility and adaptivity, allowing LLMs to transcend their typical “maximum likelihood” behaviors and mimic how humans adapt their linguistic behaviors to complete tasks. The framework’s debate structure fosters a balanced reasoning process between exploring various perspectives and exploiting the known priors, towards achieving the complex,

intricate capabilities required for AGI.

The contributions of this chapter are:

1. *EVINCE Framework Design*: Unlike approaches that use debate merely to improve accuracy via redundancy, EVINCE facilitates information discovery, bias mitigation, and decision-making that requires both breadth and depth of information.
2. *Theoretical Foundations*: EVINCE establishes a theoretical basis for SocraSynth, rooted in conditional Bayesian statistics, mutual information, and dual entropy. These principles are applied to measure, monitor, and modulate collaborative LLM interactions, contributing to a deeper understanding of how LLMs can effectively cooperate for improved decision-making. The dual entropy theory is novel and ground-breaking, illustrating how a productive decision-making process should start with room for diverse input and stable objectives, and then, through information exchange, converge to optimal decision/prediction.
3. *Empirical Validation*: We provide empirical validation of the underlying theories of EVINCE, highlighting the framework’s effectiveness in balancing exploration and exploitation to enhance prediction accuracy. We also introduce a set of maxims derived from our empirical findings, offering practical guidance for optimizing mutual information and minimizing various divergence measures.

8 Uncovering Errors and Biases with Reflective Large Language Models

Abstract Biases inherent in human endeavors pose significant challenges for machine learning, particularly in supervised learning that relies on potentially biased “ground truth” data. This reliance, coupled with models’ tendency to generalize based on statistical maximal likelihood, can propagate and amplify biases, exacerbating societal issues. To address this, our study proposes a reflective methodology utilizing multiple Large Language Models (LLMs) engaged in a dynamic dialogue to uncover diverse perspectives. By leveraging in-context learning, information theory, and divergence metrics, this novel approach fosters context-dependent linguistic behaviors, promoting unbiased outputs. Furthermore, it enables measurable progress tracking and explainable remediation actions to address identified biases.

Introduction

Errors and biases in human-labeled data present critical challenges for machine learning models, especially in healthcare, news, education, and public policy, where their outputs can profoundly shape public perception and decision-making [111]. Er-

rors, such as diagnostic mistakes, arise from knowledge gaps or lack of expertise, while biases, including ideological and societal distortions, can be consciously or unconsciously introduced by annotators. These flaws compromise the integrity of ground truth data, propagating through machine learning pipelines and generating undesirable outcomes [112, 113, 114].

AI systems are particularly vulnerable to these flaws, as models trained on inaccurate or biased ground truth data tend to replicate and amplify these issues through maximum likelihood estimation. In healthcare, diagnostic errors can lead to poor treatment recommendations [115], while in news, partisan annotations—such as labeling a biased article as neutral—mislead both human readers and automated classifiers, distorting public discourse [111, 116]. The impact extends beyond individual sectors: in education, biased data can reinforce stereotypes, while in public policy, it can result in discriminatory decisions. Ensuring that models learn from accurate and impartial ground truth data is therefore essential to the responsible deployment of AI across all domains.

This chapter focuses on bias detection and correction in news annotations, using news as the testbed to explore how reflective dialogues among LLMs can mitigate biases. News content is particularly vulnerable to ideological biases, as annotators' personal views often shape the interpretation of politically sensitive topics. Real-world evidence shows how annotation practices differ based on political affiliation. Real data illustrate that Democratic-leaning annotators may judge scandals involving Democrats more harshly than Republicans, and vice versa, highlighting the need for tools to balance these biases.

To address these challenges, we introduce the Reflective LLM

Dialogue Framework (RLDF), which implements checks and balances using multiple LLM instances in structured dialogues. RLDF conditions two instances to take opposing stances: one supports the original label, while the other introduces alternative perspectives. These reflective exchanges foster deeper insights and help uncover potential biases, generating more neutral annotations through the inclusion of diverse viewpoints. This multi-LLM dialogue approach outperforms the results of a single LLM operating in isolation or providing one-off responses.

RLDF employs conditional statistics, information theory, and divergence metrics to measure the effectiveness of these dialogues. Shannon entropy [117] quantifies the diversity of perspectives, while mutual information [118] measures the quality of the exchange. To track the convergence toward unbiased outcomes, we apply Jensen-Shannon divergence (JSD) [119], Wasserstein distance (WD) [120], and cross-entropy (CE) [121], ensuring that the remediation actions are measurable and transparent for further refinement by human reviewers.

Our empirical studies validate the effectiveness of RLDF, and the contributions of this chapter are summarized as follows:

1. Adversarial and Reflective Inspection Framework: RLDF provides a structured framework that encourages adversarial and reflective inspection of ground-truth labels. Through dialogue, participating LLM instances examine, challenge, and explain biases embedded in the original annotations by offering various perspectives. For example, in news annotation, RLDF reveals hidden ideological biases by generating alternative interpretations for politically sensitive content, leading to more neutral labeling.
2. Careful Modulation of Linguistic Behaviors for Balanced Ex-

ploration and Exploitation: The effectiveness of RLDF lies in its careful modulation of linguistic behaviors among participating LLM instances, alternating between contentious and conciliatory interactions. This dynamic trade-off fosters exploration of new perspectives while consolidating well-supported viewpoints. Information-theoretic and statistical metrics, including Shannon entropy, mutual information, Jensen-Shannon divergence, Wasserstein distance, and cross-entropy, are employed to measure opinion diversity, information flow, and the strength of the final assessment.

3. Effective Results and Impact on Improving Labels and Mitigating AI Bias: RLDF successfully mitigates AI biases, ensuring more reliable and unbiased model outputs across domains such as news, healthcare [4], and public policy. These outcomes demonstrate RLDF’s significant impact in refining labels, enhancing fairness, and promoting responsible AI deployment.

9 Modeling Emotions in Multimodal LLMs

Abstract

In human-computer interaction, recognizing and responding to a user’s emotional state is crucial for effective communication and successful task completion. For instance, a caregiving AI agent capable of detecting pain or depression in a patient could offer tailored empathetic support and appropriate medical interventions while adhering to ethical guidelines and safeguarding patient well-being. This chapter examines cognitive research on human emotions and proposes the Behavioral Emotion Analysis Model (BEAM), a novel emotion spectrum framework that incorporates both basic emotions and their linguistic antonyms. BEAM provides a comprehensive way to understand and represent emotional states in language and is designed to be integrated with Large Language Models (LLMs). By leveraging BEAM, LLMs can adapt their linguistic behaviors and expressions based on the detected emotional state of the user, ensuring responses are both empathetic and ethically aligned.

Introduction

During the development of SocraSynth [2], a multi-LLM debate framework, we discovered a fundamental principle about Large

Language Model (LLM) behavior. While investigating how to control debate "contentiousness," we found that an LLM's linguistic behavior could be systematically altered through emotional conditioning. High contentiousness produced confrontational tones and polarized language, while low contentiousness led to agreeable, considerate discourse. This observation went beyond the original scope of improving multi-agent debates—it revealed a mechanism for steering LLM behavior through emotional states.

Most multi-agent debate (MAD) systems [122, 107, 123, 124, 108, 106, 125] function as ensemble learning techniques, similar to bagging [126] or mixtures of experts [105], where LLMs simply exchange ideas without deep exploration. Our work with SocraSynth and EVINCE [4] addressed this limitation by dynamically modulating emotional states throughout the debate. High contentiousness drives LLMs to explore novel perspectives and challenge existing viewpoints, while low contentiousness promotes the synthesis of established ideas. This emotional modulation creates a natural debate progression: from vigorous exploration of diverse viewpoints, through reasoned analysis and refutation, to the emergence of well-examined, conciliatory conclusions.

While LLM training is often viewed simply as next-token prediction, its effects are far more profound. Training documents represent humans pursuing diverse goals—conducting research, exchanging opinions, expressing emotions—through a vast array of linguistic behaviors. This understanding, combined with our experience in modeling contentiousness through in-context learning, suggests an intriguing possibility: can we condition LLMs with specific goals and emotions to generate outputs that

leverage these learned linguistic behaviors? Recent empirical studies support this approach, showing that LLM outputs can be traced to their source [127] and that in-context learning operates as conditional statistics in a Bayesian framework [128].

Our exploration through bias reduction work in news articles and Wikipedia content [129] demonstrated that emotional states significantly influence LLM outputs. This finding, combined with our debate framework experiences, suggests that a mathematical model of emotions could provide a foundation for systematic behavior control. Before exploring the mapping between emotions and behaviors in the next chapter, we must first establish a rigorous framework for representing and manipulating basic emotions.

To lay the groundwork for emotion-based behavior control, this chapter develops the Behavioral Emotion Analysis Model (BEAM) [9]. While LLMs were initially seen as “black boxes” [63], our observations, along with insights from Prof. Stuart Russell, suggest that emotional states can be systematically modeled and conveyed to LLMs via *context*. Our model addresses three fundamental questions:

1. *What basic emotions form a complete basis?* We identify k fundamental emotion spectra, each defined by negative and positive antonyms (e.g., "hate-love", "anxiety-calmness"). We focus on basic emotions while excluding complex emotions like "regret" that arise from combinations of basic states. Each spectrum represents a continuous axis along which emotional states can be measured and modified.
2. *How can we mathematically manipulate emotions?* We develop a mathematical framework using negation and scaling operations for precise positioning of emotional states

along each spectrum. For instance, given the hate-love spectrum, we can represent intermediate states through scaling (e.g., $0.7 \times \text{love}$) and use negation to move between opposing states (e.g., $\neg\text{hate} \approx \text{love}$). These operations provide the foundation for systematic emotion manipulation.

3. *Can emotions predict behaviors?* We conduct a preliminary study using self-supervised learning to explore the relationship between emotional states and linguistic behaviors. By analyzing text samples from our debate framework, we train a model to predict behavioral patterns from emotional states without explicit labeling. This study validates our emotional spectra's utility in modeling behavioral outcomes while providing insights for more comprehensive behavior mapping in [5].

While the next chapter will explore how these emotional states map to specific linguistic behaviors in depth, this chapter focuses on establishing the mathematical framework for representing and manipulating basic emotions, validated through preliminary self-supervised learning experiments. By grounding our model in both mathematical rigor and empirical testing, we create a foundation for systematic emotion-based behavior control in LLMs.

10 A Three-Branch Checks-and-Balances Framework for Context-Aware Ethical Alignment of Large Language Models

Abstract This chapter introduces a three-branch checks-and-balances framework for ethical alignment of Large Language Models (LLMs). Inspired by governmental systems, the framework implements three independent yet interacting components: LLMs as the executive branch for knowledge generation, DIKE (named after the goddess of justice) as the legislative branch establishing ethical guardrails, and ERIS (the goddess of discord) as the judicial branch for contextual interpretation. The DIKE-ERIS duality, through their adversarial interaction, enables adaptation to diverse cultural contexts while maintaining consistent ethical principles. This architecture addresses fundamental limitations of reinforcement learning with human feedback (RLHF) by providing interpretable, adaptable, and culturally-aware ethical reasoning. Through self-supervised learning and adversarial testing, our framework demonstrates how emotional

modeling can guide linguistic behaviors toward ethical outcomes while preserving the independence of knowledge generation, ethical oversight, and contextual interpretation.

Introduction

This research introduces an alternative to Reinforcement Learning from Human Feedback (RLHF) [65, 69] to address ethical concerns in Large Language Models (LLMs). While RLHF has demonstrated success, it faces two significant challenges: susceptibility to societal biases in increasingly polarized human feedback, and vulnerability to reward hacking [130, 131] that may lead to unethical behaviors.

A significant limitation of current research is its narrow focus on isolated behaviors, such as movie ratings or toxic language. This reactive approach is akin to “playing Whack-A-Mole,” where individual issues are suppressed without solving the core behavioral problem. For example, merely instructing someone to consistently make their bed does not necessarily change their underlying habits or attitudes. Additionally, fixing one issue may inadvertently aggravate others. Users have reported performance degradations in ChatGPT due to RLHF modifications that altered (forgot) the optimal parameters for other tasks [132, 133]. Similarly, psychological studies show that addressing an addiction problem often reveals underlying issues and triggers side effects [134, 135].

To overcome these challenges, we propose a novel framework inspired by the checks and balances of governmental systems. Our architecture integrates three independent but interacting components: LLMs serve as the executive branch responsible for knowledge generation; DIKE (named after the Greek goddess



Figure 10.1: Three Framework Components: Executive LLMs (bottom), Legislative (upper-left), and Judicial (upper-right)

of justice) functions as the legislative branch, establishing ethical guardrails; and **ERIS** (named after the goddess of discord) acts as the judicial branch, providing adversarial testing and cultural interpretation. In mythology, Dike embodies order and justice, while her adversary, Eris, represents discord and strife—a duality that our framework leverages to balance ethical guidance with adversarial perspectives.

Figure 10.1 presents our framework architecture where three neurally independent components—LLMs as the foundation, with **DIKE** and **ERIS** as oversight layers—interact through structured interfaces while maintaining strict separation of their neural architectures and parameters.

Central to this framework is **DIKE** (**D**iagnostics, **I**nterpretation, **K**nowledge-independent learning, and **E**thical guardrails), which operates as an independent advisor on behavioral ethics. By decoupling ethical oversight from the LLM’s knowledge processing, **DIKE** ensures that ethical improvements do not inter-

fere with knowledge representation, while enabling adaptive and culturally-aware ethical guidance. For example, while the principle “do not lie” generally applies, context-sensitive interpretation may be necessary, such as when a doctor or family member conceals a terminal diagnosis to protect a patient. Likewise, cultural differences in attitudes toward issues like alcohol consumption, abortion, or same-sex marriage necessitate flexible, context-sensitive ethical reasoning.

The interplay between DIKE and ERIS introduces four key innovations:

1. *Emotion-Driven Behavioral Modeling*: Building on BEAM (Behavioral Emotion Analysis Model) [9], DIKE employs self-supervised learning to analyze how emotions manifest in linguistic behaviors, creating quantifiable relationships between emotional states and their corresponding language patterns in text.
2. *Behavior-Aware Ethical Guardrails*: The framework establishes guidelines that consider both content and linguistic behavior, preventing harmful or manipulative communication while preserving factual accuracy and emotional authenticity. The interpretation of these guardrails adapts dynamically across cultural contexts, preserving consistency while enabling context-sensitive interpretation.
3. *Adversarial Behavioral Testing*: ERIS actively challenges DIKE’s ethical guidelines by presenting diverse cultural perspectives and edge cases. This adversarial dynamic strengthens the framework’s ability to handle complex ethical scenarios while maintaining cultural sensitivity and considering context.
4. *Ethical Content Transformation*: When detecting ethically problematic content, DIKE performs targeted revisions (in-

dependent of the LLMs) that preserve intended emotional expression while ensuring ethical compliance, adapting its responses to specific cultural and contextual requirements. ERIS continuously tests these transformations against various cultural contexts and edge cases, validating both the ethical alignment and contextual appropriateness.

Through structured interfaces, these components work together in our three-branch architecture to provide robust ethical oversight while maintaining adaptability to evolving cultural norms. By keeping the three models, LLMs, DIKE, and ERIS, architecturally independent, we prevent interference between knowledge representation and ethical reasoning while enabling sophisticated ethical adaptation through their structured interactions. This approach represents a significant advancement in developing AI systems capable of culturally-aware, emotionally intelligent, and ethically sound communication.

11 Beyond Computation: Consciousness Modeling

Abstract The CoCoMo model proposes a computational solution to the challenge of incorporating ethical and emotional intelligence considerations into AI systems, with the aim of creating AI agents that combine knowledge with compassion. To reach this goal, CoCoMo focuses on fairness, beneficence, empathy, non-maleficence, adaptability, and critical and exploratory thinking abilities. CoCoMo employs consciousness modeling, reinforcement learning, and prompt template formulation to support these desired traits. By incorporating ethical and emotional intelligence considerations, a generative AI model can potentially lead to improved fairness, reduced toxicity, and increased reliability.

Introduction

Narrow AI, often referred to as System-1 AI following Kahneman’s terminology [136], excels in executing well-defined, specific tasks through machine learning algorithms, including object recognition and language translation. However, this type of AI is not as effective in handling advanced generative AI functions that require reasoning, critical and exploratory thinking, or the modeling and regulation of emotions and behaviors. Such

complex tasks go beyond the capabilities of System-1 AI, highlighting its limitations.

To address these limitations, researchers (e.g., Yoshua Bengio [137]) have proposed the development of system-2 AI, which aims to mimic human cognitive abilities. Several generative models have been developed since 2022 for text [80, 64, 65, 138], image [139, 140], and video generation [141]. However, these models face issues of bias, toxicity, robustness, and reliability [142, 143].

In this chapter, we propose a solution to address these concerns by modeling emotional intelligence and ethical guardrails within a generative AI model itself, drawing on insights from the study of human consciousness. We believe that addressing these issues outside of a generative AI model using human subjective feedback and reinforcement learning is equivalent to imposing censorship on user-generated content, which is a heuristic-based and non-scalable solution [144, 145].

Human consciousness is understood to manage both impulsive and reflective aspects of the unconscious, enabling compromises between competing goals and values. Emotions typically arise as impulsive reactions to stimuli, while ethics act as guardrails that help modulate or regulate emotion-steered motivations to sin. Developing a grasp of how human consciousness functions, not necessarily in physical terms but at least functionally, can offer vital insights for crafting a regulatory mechanism within a LLM. This mechanism would direct linguistic behavior and shape the linguistic features employed to achieve specific goals.

The nature and origin of consciousness have been studied for centuries, resulting in various theories, including the global workspace theory [146], integrated information theory [147, 148, 149], neural correlates of consciousness approach [150, 151], and

attention schema theory [152, 153], among others. These studies of consciousness provide valuable insights for architecting system-2 AI.

This chapter defines system-2 AI’s desired traits and capabilities using a functionalist approach to consciousness. These include knowledge, fairness, beneficence, non-maleficence, empathy, adaptability, transparency, and critical and exploratory thinking abilities. While this list is not exhaustive, it provides a starting point for developing ethical guardrails and emotional intelligence in AI systems. Depending on the context and application of AI, additional ethical considerations or modifications to these principles may be necessary.

To embody these capabilities and principles, we introduce the Computational Consciousness Model (CCM), which leverages priority-based scheduling, reward-based optimization, and Socratic dialogues. CCM offers customization based on cultural and individual requirements through adaptive prompt templates [6, 154], and facilitates the transition between unconsciousness and consciousness states through a multi-level feedback scheduler and interrupt mechanism.

To enable emotion and behavior modeling and regulation, as well as critical and exploratory thinking, CCM interacts with large language models¹ [80, 155, 64, 65, 138] using interactive question-answer-based dialogues. Furthermore, a reinforcement learning module maps external values and rewards to internal task-scheduling priorities.

CCM has the potential to support the development of adaptive computational consciousness that integrates knowledge and compassion, and models emotional intelligence for generative AI

¹Recent studies, such as [74], propose referring to these multimodal pre-trained models as foundation models.

systems. This approach could significantly benefit humanity and society.

12 A Retrospective and Adaptive Framework to Improve LLMs

Abstract RAFEL is a retrospective and adaptive framework designed to benchmark private Large Language Models (LLMs) against teacher LLMs, identifying discrepancies in responses. Following the initial benchmarking, RAFEL categorizes these discrepancies into four distinct categories, based on cognitive levels and types of errors. Subsequent phases involve a detailed diagnosis and deep-probing to uncover the root causes behind each category of discrepancy. Teacher LLMs play a crucial role in interrogating the private LLM, shedding light on the subtleties of its performance issues. With a clear understanding of the symptoms and their underlying causes, RAFEL prescribes targeted remedies, accompanied by recommendations for relevant data sources to enhance the private LLM’s performance via either fine-tuning, RAG, or both. Empirical studies validate RAFEL’s effectiveness in diagnosing and enhancing the capabilities of localized LLMs.

Introduction

The emergence of Large Language Models (LLMs) like GPT [65] and Gemini [95] has significantly advanced the field of natural language processing, enabling the generation of text that closely mimics human writing and offers deep insights across varied domains. Despite their transformative potential, the deployment and scalability of these models pose considerable computational and data challenges. A practical response has been the fine-tuning of medium-sized, open-source models such as LLaMa [67] for specialized needs, allowing organizations to strike a balance between performance and feasibility, while also prioritizing data privacy and model customization for unique applications.

The shift towards using privately fine-tuned or locally deployed LLMs brings about essential management and technical challenges, vital for corporate strategy, governance, and innovation. This chapter explores the technical challenges of this shift, including:

- Justifying the choice of private LLMs over public counterparts by establishing relevant performance metrics and benchmarks for these specialized models.
- Conducting in-depth error analysis to pinpoint the root causes of performance issues in private LLMs, ensuring targeted and effective remediation strategies.
- Identifying specific, high-quality data crucial for the fine-tuning of private LLMs, aimed at enhancing their accuracy and domain relevance.
- Implementing Retrieval-Augmented Generation (RAG) to dynamically incorporate external, updated knowledge sources, improving the model's responsiveness and breadth of knowl-

edge.

- Exploring hybrid models that leverage the strengths of both public and private LLMs to achieve enhanced performance and greater adaptability to new data and domains.

We introduce **RAFEL**, a framework designed for the retrospective and adaptive enhancement of LLMs, addressing these technical challenges. **RAFEL** strategically balances cost and performance by incorporating sophisticated diagnostic algorithms. These algorithms effectively identify and address the root causes of inefficiencies, ensuring that solutions are economically viable.

RAFEL employs advanced benchmarking metrics across cognitive levels, providing a thorough LLM performance assessment. Central to its diagnostics are two key algorithms: **DIAG**, for non-invasive¹ evaluation, and **PRBE** for thorough, invasive probing. This combination allows **RAFEL** to detect and understand both surface-level and deep-seated performance issues, facilitating targeted data source acquisition for enhancement.

RAFEL is proficient in creating targeted, effective remediation strategies, ensuring data privacy and security, validated through real-world data studies. The novelty claims of **RAFEL** include:

1. *Deep Probe with Cognitive and Error Type Analysis*: **RAFEL** goes beyond traditional error rate analysis by deeply probing into the LLM’s responses, categorizing errors within cognitive levels (recollection, comprehension, analysis, reasoning) and types (hallucination, biases), enabling a deep understanding of the model’s performance issues.
2. *Fine-grained, Precise Data Augmentation*: Contrasting with the conventional manual search for coarse-grained data aug-

¹Non-invasive methods evaluate without interacting with the LLM’s internal data, whereas invasive methods directly engage with the LLM, accessing potentially sensitive data.

mentation, RAFEL identifies the required data and performs a more precise and relevant data enhancement that directly addresses the identified cognitive and error type deficiencies.

3. *Dynamic Remediation Playbook*: RAFEL dynamically adjusts its remediation strategy based on real-time analysis of data and errors, akin to adapting tactics in sports, ensuring the most effective and appropriate intervention is applied.

13 Discovering Insights Beyond the Known

Abstract Human knowledge, vast as it is, often falls short in grasping intricate interdisciplinary domains fully. In contrast, foundation models like GPT-4, endowed with extensive multidisciplinary knowledge, can potentially bridge this gap. Significantly, we leverage the vast expanses of GPT-4's knowledge, banking on its ability to frame questions that might elude human intuition, thus paving the way for the emergence of fresh insights and potentially novel knowledge. In this study, we convened a unique committee comprising a moderator (the authors) and two GPT-4 agents. The dialogue is ignited by the ancient narrative of Adam and Eve, setting the stage for a rich exchange between the GPT-4 agents. This conversation derives from the age-old tale, as the agents investigate three intertwined domains: the significance of myths in ecological interpretation, the intricate ethical and philosophical quandaries surrounding AI, and the enigmatic realm of the human brain as complemented by technology. This dialogue not only unveils captivating insights but also underscores the indispensable value of interdisciplinary exchanges. Foundation models, as demonstrated, can catalyze such dialogues, equipping us to traverse expansive knowledge landscapes and explore domains previously beyond human comprehension.

Introduction

In our recent study on GPT-4 [63], we observed that GPT-4 along with analogous foundation models, manifests a *Polydisciplinary* capacity [2]. (For clarity, we use “GPT-4” to collectively refer to these foundation models, given that our experiments are centered on GPT-4.) Trained on a vast spectrum of topics from varied sources, GPT-4 stands apart from human specialists. Such specialists, while deeply knowledgeable in their specific fields, often lack a broad understanding outside their particular domain. In contrast, GPT-4 processes knowledge without being tethered to domain boundaries. It doesn’t compartmentalize a query strictly as a “physics question” or a “philosophy question,” but crafts an integrated response, drawing from its multidisciplinary training data.

From a perspective of sheer knowledge breadth, GPT-4 arguably outpaces the average human. Its exposure to an enormous volume of documents endows it with a repository potentially wider than most human counterparts. However, volume isn’t synonymous with depth. True depth often stems from intangible intuitions, insights, personal experiences, and cultural contexts. Considering GPT-4 lacks evolutionary experiences—ranging from survival instincts to the full spectrum of human emotions—we must ask: Can GPT-4 produce literature that deeply resonates with human sensibilities?

This study aims to ascertain if the polydisciplinary attributes of GPT-4 can generate insights that transcend standard human perspectives. We divide our research into two avenues: first, exploring the potential of GPT-4 to reveal “unknown unknowns,” and second, assessing its aptitude for crafting emotionally impactful literature. This chapter examines the former, utilizing

the universally recognized biblical tale of Adam and Eve and their consumption of the forbidden fruit as a common thematic foundation. Through this exploration, we aim to uncover viewpoints potentially beyond the realm of typical human cognition.

Our methodology revolves around orchestrating a dialogue between multiple GPT-4 agents. Within the experimental framework, a moderator (represented by the authors) sets the initial intent and context for the conversation. The number of participating agents and their underlying foundation models can be adjusted as needed. In this study, our committee consists of two agents based on the GPT-4 model, referred to as GPT-A and GPT-B. Once initialized, the agents engage in conversation autonomously, with minimal moderation (discussed shortly). The resultant dialogue is thoroughly analyzed to discern conversational patterns and depth of content. This in-depth examination facilitates the identification of diverse themes the GPT-4 model gravitates towards. Our underlying hypothesis posits that the discourse and exchanges between these agents can unearth insights—“unknown unknowns”—that were previously elusive to human understanding.

While the polydisciplinary capabilities of GPT-4 offer an unparalleled breadth and depth exceeding that of the moderator, the role of the moderator remains indispensable. This role channels the “exploratory” nature of the conversation, guiding it towards predefined objectives and ensuring its convergence within a set time frame. In this experiment, the initial spark for the dialogue is the narrative of Adam and Eve. Without prompting, the agents autonomously suggest probing the story from ten unique perspectives. Yet, after a series of exchanges, GPT-B expresses a keen interest in delving deeper. Following this, in collabora-

tion with both agents, the moderator narrows down the scope of the dialogue to three key topics: ecological interpretation, philosophical exploration, and the neuroscientific angle. The intricate dialogues spanning these three domains—namely AI interwoven with Ecology, Neuroscience coupled with AI, and Neuroscience meshed with Ecology—are indeed engrossing. Throughout the discussion, both agents present a multifaceted perspective, shedding light on the diverse interpretations of the Garden of Eden, both prior to and following its seminal event. In the final stretch, the moderator verifies with both agents if they are poised to transition into the conclusion phase.

While our research unveils fascinating insights, it's essential to acknowledge several inherent limitations and constraints:

1. Model Training and Bias: GPT-4, akin to other machine learning models, is informed by pre-existing datasets. Therefore, the viewpoints, knowledge, and biases ingrained in this data can shape its outputs. It implies that GPT-4's responses might echo the historical and cultural biases present in the data upon which it was trained.
2. Interactivity Limitation: Conversations between two GPT-4 agents essentially access the same foundational knowledge. Consequently, while the discourse may encompass a range of viewpoints due to query processing, it won't yield wholly novel information. To glean potentially varied insights, it could be beneficial to facilitate dialogues between different iterations of GPT (like GPT-3 and GPT-4) or even entirely distinct foundation models such as LaMDA by Google and LLaMA by Meta AI.
3. Interpretation Subjectivity: Analysis of GPT-4's dialogues is susceptible to the prism of human interpretation. As such,

different analysts might extract diverse conclusions from identical data sets.

We emphasize that the core intent of our study was experimental, rooted in the notion that a foundation model such as GPT-4 may offer knowledge both wider in scope and deeper in intricacy than most human capacities. While we've shed light on its potential strengths, we remain cognizant of the inherent risks and limitations. Our work's primary novelty lies in pioneering a discourse-driven approach with AI agents. Significantly, we leverage the vast expanses of GPT-4's knowledge, banking on its ability to frame questions that might elude human intuition, thus paving the way for the emergence of fresh insights and potentially novel knowledge.

14 Bibliography

- [1] Edward Y. Chang. *LLM Collaborative Intelligence: The Path to Artificial General Intelligence*. SocraSynth.com, October 2024.
- [2] Edward Y Chang. Examining GPT-4's Capabilities and Enhancement with SocraSynth. In *The 10th International Conf. on Computational Science and Computational Intelligence*, December 2023.
- [3] Jocelyn J. Chang and et al. SocraHealth: Enhancing Medical Diagnosis and Correcting Historical Records. In *The 10th International Conf. on Computational Science and Computational Intelligence*, December 2023.
- [4] Edward Y Chang. EVINCE: Optimizing Adversarial LLM Dialogues via Conditional Statistics and Information Theory. In *arXiv:2408.14575*, August 2024.
- [5] Edward Y. Chang. A Three-Branch Checks-and-Balances Framework for Context-Aware Ethical Alignment of Large Language Models. In *Advances in Neural Information Processing Systems, SafeGenAI Workshop*, 2024.
- [6] Edward Y. Chang. CRIT: Prompting Large Language Models With the Socratic Method. *IEEE 13th Annual Computing and Communication Workshop and Conference*, March 2023.

14 Bibliography

- [7] Edward Y. Chang. Uncovering Biases with Reflective Large Language Models. In *arXiv:2408.13464*, 2024.
- [8] Edward Y. Chang and Emily J. Chang. Discovering Insights Beyond the Known: A Dialogue Between GPT-4 Agents from Adam and Eve to the Nexus of Ecology, AI, and the Brain. In *Stanford InfoLab Technical Report*, 2023.
- [9] Edward Y. Chang. Behavioral Emotion Analysis Model for Large Language Models (invited paper). In *Proceedings of the 7th IEEE MIPR Conference*, August 2024.
- [10] Edward Y. Chang. CoCoMo: Computational Consciousness Modeling for Generative and Ethical AI. In *arXiv:2304.02438*, 2023.
- [11] Zellig S. Harris. Distributional structure. *WORD*, 10(2-3):146–162, 1954.
- [12] George A Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41, 1995.
- [13] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [14] Susan T Dumais. Latent semantic analysis. *Annual review of information science and technology*, 38(1):188–230, 2004.
- [15] David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022, 2003.

- [16] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space, 2013.
- [17] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, 2014.
- [18] Ashish Vaswani et al. Attention is all you need. *Advances in neural information processing systems*, 2017.
- [19] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019.
- [20] Sung-Hyuk Cha. Comprehensive Survey on Distance Similarity Measures between Probability Density Functions. 2007.
- [21] Paul Jaccard. The distribution of the flora in the alpine zone. *New Phytologist*, 11(2):37–50, 1912.
- [22] A Feature-Integration Theory of Attention. *Cognitive Psychology*, 12(1):97–136, 1980.
- [23] Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne Hubbard, and Lawrence D Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551, 1989.
- [24] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document

14 Bibliography

- recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [25] Max Wertheimer. Experimentelle studien über das sehen von bewegung. *Zeitschrift fur psychologie*, 61:161–165, 1912.
- [26] Kurt Koffka. *Principles of Gestalt Psychology*. Harcourt, Brace and Company, New York, 1935.
- [27] Donald E. Broadbent. *Perception and Communication*. Pergamon Press, 1958.
- [28] Baitao Li, E. Chang, and Ching-Tung Wu. DPF - A Perceptual Distance Function for Image Retrieval. In *Proceedings. International Conference on Image Processing*, volume 2, pages II–II, 2002.
- [29] Beita Li, Edward Y. Chang, and Yi Wu. Discovery of A Perceptual Distance Function for Measuring Image Similarity. *Multimedia Systems*, 8:512–522, 2003.
- [30] Douglas L Medin, Robert L Goldstone, and Dedre Gentner. Respects for Similarity. *Psychological Review*, 100(2):254–278, 1993.
- [31] R. L. Goldstone. Similarity, interactive activation, mapping. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 20(3):3–28, 1994.
- [32] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In *Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume*

- 1, NIPS'12, page 1097–1105, Red Hook, NY, USA, 2012.
Curran Associates Inc.
- [33] D. H. Hubel and T. N. Wiesel. Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. *The Journal of physiology*, 160(1):106–154, 1962.
 - [34] Kenneth A. Norman, Sean M. Polyn, Greg J. Detre, and James V. Haxby. Beyond mind-reading: multi-voxel pattern analysis of fmri data. *Trends in Cognitive Sciences*, 10:424–430, 2006.
 - [35] John-Dylan Haynes and Geraint Rees. Decoding mental states from brain activity in humans. *Nature Reviews Neuroscience*, 7:523–534, 2006.
 - [36] Allen Newell and Herbert A. Simon. Logic theorist and general problem solver. *Journal of the ACM (JACM)*, 1:256–260, 1956.
 - [37] Robert K. Lindsay, Bruce G. Buchanan, Edward A. Feigenbaum, and Joshua Lederberg. Applications of artificial intelligence for organic chemistry: The dendral project. *McGraw-Hill advanced computer science series*, 1979.
 - [38] Edward Hance Shortliffe, Bruce G. Buchanan, and Edward A. Feigenbaum. Mycin: A rule-based computer program for advising physicians regarding antimicrobial therapy selection. *AI in Medicine*, 10:199–208, 1975.
 - [39] William F. Clocksin and Christopher S. Mellish. *Programming in Prolog*. Springer-Verlag, 1981.

14 Bibliography

- [40] J. Ross Quinlan. Induction of decision trees. *Machine Learning*, 1:81–106, 1986.
- [41] I. T. Jolliffe. Principal component analysis. *Springer Series in Statistics*, 1986.
- [42] Peter W. Foltz Thomas K Landauer and Darrell Laham. An introduction to latent semantic analysis. *Discourse Processes*, 25(2-3):259–284, 1998.
- [43] Joshua B. Tenenbaum, Vin De Silva, and J. C. Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290:2319–2323, 2000.
- [44] Laurens van der Maaten, Eric O. Postma, and Jaap van den Herik. Dimensionality Reduction: A Comparative Review. In *TiCC TR 2009-005*, 2008.
- [45] Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine Learning*, 20(3):273–297, 1995.
- [46] David Barber. *Bayesian Reasoning and Machine Learning*. Cambridge University Press, 2012.
- [47] Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.
- [48] Andrew Gelman, John B. Carlin, Hal S. Stern, David B. Dunson, Aki Vehtari, and Donald B. Rubin. *Bayesian Data Analysis, Third Edition*. Chapman and Hall/CRC, 2014.
- [49] Daphne Koller and Nir Friedman. *Probabilistic Graphical Models - Principles and Techniques*. The MIT Press, 2009.

- [50] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2012.
- [51] T.M. Mitchell. *Machine Learning*. McGraw-Hill series in computer science. McGraw Hill, 1997.
- [52] Jeffrey Dean and Sanjay Ghemawat. Mapreduce: simplified data processing on large clusters. *Commun. ACM*, 51(1):107–113, 1 2008.
- [53] Michael Cafarella, Edward Chang, Andrew Fikes, Alon Halevy, Wilson Hsieh, Alberto Lerner, Jayant Madhavan, and S. Muthukrishnan. Data management projects at google. *SIGMOD Rec.*, 37(1):34–38, 3 2008.
- [54] Jiakai Liu, Rong Hu, Meihong Wang, Yi Wang, and Edward Y. Chang. Web-Scale Image Annotation. In *Advances in Multimedia Information Processing - PCM 2008*, pages 663–674. Springer Berlin Heidelberg, 2008.
- [55] Edward Y. Chang. *Foundations of Large-Scale Multimedia Information Management and Retrieval: Mathematics of Perception*. Springer, 2011.
- [56] Edward Y. Chang, Kaihua Zhu, Hao Wang, Hongjie Bai, Jian Li, Zhihuan Qiu, and Hang Cui. Psvm: Parallelizing support vector machines on distributed computers. In *Proceedings of the 20th International Conference on Neural Information Processing Systems*, NIPS’07, page 257–264, Red Hook, NY, USA, 2007. Curran Associates Inc.
- [57] Haoyuan Li, Yi Wang, Dong Zhang, Ming Zhang, and Edward Y. Chang. PFP: Parallel FP-Growth for Query Recommendation. In *ACM RecSys ’08*. ACM, 2008.

14 Bibliography

- [58] Yi Wang, Hongjie Bai, Matt Stanton, Wen-Yen Chen, and Edward Y. Chang. Plda: Parallel latent dirichlet allocation for large-scale applications. In Andrew V. Goldberg and Yunhong Zhou, editors, *Algorithmic Aspects in Information and Management*, pages 301–314, Berlin, Heidelberg, 2009. Springer Berlin Heidelberg.
- [59] Wen-Yen Chen, Yangqiu Song, Hongjie Bai, Chih-Jen Lin, and Edward Y. Chang. Parallel Spectral Clustering in Distributed Systems. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(3):568–586, 2011.
- [60] Z. Zheng, W. Jiang, G. Wu, and E. Y. Chang. Speedo: Parallelizing stochastic gradient descent for deep convolutional neural network. In *LearningSys, NeurIPS*, 2015.
- [61] Artur d’Avila Garcez and Luis C. Lamb. Neurosymbolic ai: The 3rd wave, 2020.
- [62] Lihui Lu, Qizhong Li, Jiyuan Liu, and Miaohua Huang. Combining domain knowledge and deep learning methods for vehicle trajectory prediction. *Journal of Physics: Conference Series*, 2303(1), 7 2022.
- [63] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, and more. Sparks of Artificial General Intelligence: Early experiments with GPT-4, 2023.
- [64] OpenAI. Chatgpt, 2021.
- [65] OpenAI. GPT-4 Technical Report, 2023.
- [66] Romal Thoppilan, Daniel De Freitas, Jamie Hall, and et al. Lamda: Language models for dialog applications, 2022.

- [67] Hugo Touvron, Louis Martin, Kevin Stone, and et al. Llama 2: Open foundation and fine-tuned chat models, 2023.
- [68] Papers with Code Corp. Multi-task language understanding on mmlu, 10 2023.
- [69] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, and more. Training language models to follow instructions with human feedback, 2022.
- [70] Jack Rae. Compression for AGI, Stanford MLSys, #76. <https://www.youtube.com/watch?v=d04TPJkeaaU>, 2023.
- [71] Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, and more. Scaling language models: Methods, analysis & insights from training gopher, 2022.
- [72] Sam Altman and Lex Friedman. GPT-4, ChatGPT, and the Future of AI, Lex Fridman Podcast no.367, 2023.
- [73] Edward Y. Chang. LLM Debate on the Middle East Conflict: Is It Resolvable? *Stanford University InfoLab Technical Report*, 10 2023.
- [74] Rishi Bommasani, Drew A. Hudson, and more. On the opportunities and risks of foundation models, 2022.
- [75] Rohan Taori, Achal Dave, Vaishaal Shankar, Nicholas Carlini, Benjamin Recht, and Ludwig Schmidt. Measuring

14 Bibliography

- robustness to natural distribution shifts in image classification. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS'20, 2020.
- [76] Rich Sutton. The Bitter Lesson, 2019.
 - [77] David G. Lowe. Distinctive image features from scale-invariant keypoints. *Int. J. Comput. Vision*, 60(2):91–110, 11 2004.
 - [78] Olga Kovaleva, Alexey Romanov, Anna Rogers, and Anna Rumshisky. Revealing the dark secrets of bert, 2019.
 - [79] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, and more. Scaling laws for neural language models, 2020.
 - [80] Tom B. Brown, Benjamin Mann, Nick Ryder, and more. Language models are few-shot learners, 2020.
 - [81] Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang, Jinran Nie, Yuxuan Ding, Jianwei Yue, and Yupeng Wu. How close is chatgpt to human experts? comparison corpus, evaluation, and detection, 2023.
 - [82] Thomas Wolf, Victor Sanh, Julien Chaumond, and Clement Delangue. Transfertransfo: A transfer learning approach for neural network based conversational agents, 2019.
 - [83] Chase B. Wrenn. Internet encyclopedia of philosophy, 2023.
 - [84] Wikipedia. Socratic method, 2023.

- [85] Vladimir Dobrovolskii. Word-level coreference resolution. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7670–7675, Online and Punta Cana, Dominican Republic, 11 2021. Association for Computational Linguistics.
- [86] Panupong Pasupat and Percy Liang. Compositional semantic parsing on semi-structured tables. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1470–1480. Association for Computational Linguistics, 7 2015.
- [87] Li Dong and Mirella Lapata. Coarse-to-fine decoding for neural semantic parsing. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, pages 731–742. Association for Computational Linguistics, 7 2018.
- [88] Plato. The republic.
- [89] Linda Elder and Richard Paul. *The Thinker’s Guide to the Art of Asking Essential Questions*. Rowman & Littlefield, 5th. edition, 2010.
- [90] T. Airaksinen. Socratic irony and argumentation. *Argumentation*, 36:85–100, 2012.
- [91] H. A. Stoddard and D. V. O’Dell. Would socrates have actually used the "socratic method" for clinical teaching? *Journal of general internal medicine*, 31(9):1092–96, 2016.
- [92] M. N. Browne and S. Keeley. Asking the right questions, a guide to critical thinking, 2021.

14 Bibliography

- [93] Richard Paul and Linda Elder. Critical thinking: The art of socratic questioning. *Journal of Developmental Education*, 31:34–35, 2007.
- [94] Todd Thrash, Laura Maruskin, Scott Cassidy, James Fryer, and Richard Ryan. Mediating between the muse and the masses: Inspiration and the actualization of creative ideas. *Journal of personality and social psychology*, 98:469–87, 03 2010.
- [95] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, et al. Gemini: A family of highly capable multimodal models, 2023.
- [96] Yifan Zhang, Jingqin Yang, Yang Yuan, and Andrew Chi-Chih Yao. Cumulative reasoning with large language models, 2023.
- [97] Emilio Ferrara. Fairness and bias in artificial intelligence: A brief survey of sources, impacts, and mitigation strategies. *ArXiv*, abs/2304.07683, 2023.
- [98] Otavio Parraga, Martin D. More, Christian M. Oliveira, et al. Fairness in deep learning: A survey on vision and language research. *ACM Comput. Surv.*, 12 2023. Just Accepted.
- [99] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, et al. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions, 2023.
- [100] Jie Huang and Kevin Chen-Chuan Chang. Towards reasoning in large language models: A survey. In *Findings of*

- the Association for Computational Linguistics: ACL 2023*, pages 1049–1065, 7 2023.
- [101] Larry Page. The pagerank citation ranking: Bringing order to the web, 1998.
 - [102] Isaac Adams-Hands. Google algorithms: Decoding the secrets for improved search rankings, 2023.
 - [103] Wen-Kwang Tsao. Multi-Agent Reasoning with Large Language Models for Effective Corporate Planning. In *The 10th International Conf. on Computational Science and Computational Intelligence*, December 2023.
 - [104] Edward Y Chang. SocraPedia: A Wikipedia Generated by SocraSynth with Collaborative Large Language Models. *Stanford University InfoLab Technical Report*, November 2023.
 - [105] Robert A. Jacobs, Michael I. Jordan, Steven J. Nowlan, and Geoffrey E. Hinton. Adaptive Mixtures of Local Experts. *Neural Computation*, 3(1):79–87, 03 1991.
 - [106] Julian Michael, Salsabila Mahdi, David Rein, Jackson Petty, Julien Dirani, Vishakh Padmakumar, and Samuel R. Bowman. Debate helps supervise unreliable experts, 2023.
 - [107] Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. Chateval: Towards better llm-based evaluators through multi-agent debate, 2023.
 - [108] Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Zhaopeng Tu, and Shum-

14 Bibliography

- ing Shi. Encouraging divergent thinking in large language models through multi-agent debate, 2023.
- [109] Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. Improving factuality and reasoning in language models through multiagent debate, 2023.
- [110] Yoav Freund and Robert E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1):119–139, 1997.
- [111] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning. *ACM Computing Surveys (CSUR)*, 54(6):1–35, 2021.
- [112] Jon Kleinberg, Sendhil Mullainathan, and Manish Raghavan. Inherent trade-offs in the fair determination of risk scores. *Proceedings of Innovations in Theoretical Computer Science (ITCS)*, 2017.
- [113] Andrew D Selbst, Danah Boyd, Sorelle A Friedler, Suresh Venkatasubramanian, and Janet Vertesi. Fairness and abstraction in sociotechnical systems. In *Proc. of the Conference on Fairness, Accountability, and Transparency*, pages 59–68, 2019.
- [114] Ricardo Baeza-Yates. Bias on the web. *Communications of the ACM*, 61(6):54–61, 2018.
- [115] David E. Newman-Toker, Kevin M. McDonald, Christopher J. Dy, and Linda T. Kohn. Serious Harm From Diagnostic Error in US Healthcare Systems: Estimate of Its

- Magnitude and Cost. *BMJ Quality & Safety*, 32(7):549–557, 2023.
- [116] Sanjana Gautam and Mukund Srinath. Blind spots and biases: Exploring the role of annotator cognitive biases in nlp, 2024.
- [117] Claude E. Shannon. A mathematical theory of communication. *Bell System Technical Journal*, 27(3):379–423, 1948.
- [118] T. M. Cover and J. A. Thomas. *Elements of Information Theory*. John Wiley & Sons, 2 edition, 2006.
- [119] Jianhua Lin. Divergence measures based on the shannon entropy. *IEEE Transactions on Information theory*, 37(1):145–151, 1991.
- [120] Leonid V Kantorovich. On the translocation of masses. *Doklady Akademii Nauk*, 37(7-8):199–201, 1942.
- [121] John E. Shore and Rodney W. Johnson. Axiomatic derivation of the principle of maximum entropy and the principle of minimum cross-entropy. *IEEE Transactions on Information Theory*, 26(1):26–37, 1980.
- [122] Sahar Abdelnabi, Amr Gomaa, Sarath Sivaprasad, Lea Schönher, and Mario Fritz. Cooperation, competition, and maliciousness: Llm-stakeholders interactive negotiation, 2024.
- [123] Yao Fu, Hao Peng, Tushar Khot, and Mirella Lapata. Improving language model negotiation with self-play and in-context learning from ai feedback, 2023.

14 Bibliography

- [124] Huao Li, Yu Chong, Simon Stepputtis, Joseph Campbell, et al. Theory of mind for multi-agent collaboration via large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2023.
- [125] Andries Smit, Paul Duckworth, Nathan Grinsztajn, Thomas D. Barrett, and Arnu Pretorius. Should we be going mad? a look at multi-agent debate strategies for llms, 2024.
- [126] Leo Breiman. Bagging predictors. *Machine learning*, 24(2):123–140, 1996.
- [127] Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jie Tang, et al. Longbench: A bilingual, multitask benchmark for long context understanding, 2023.
- [128] Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context learning as implicit bayesian inference. In *International Conference on Learning Representations (ICLR)*, 2021.
- [129] Edward Y Chang. Uncovering Biases with Reflective Large Language Models. In *arXiv:2408.13464*, August 2024.
- [130] Paul Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- [131] Johan Skalse, Tom Everitt, Florent Lefrancq, Arkadij Filos, and Marcus Hutter. Defining reward hacking in reinforcement learning. In *International Conference on Arti-*

- ficial *Intelligence and Statistics*, pages 8052–8079. PMLR, 2022.
- [132] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- [133] Andrei A. Rusu, Sergio Gomez Colmenarejo, Caglar Gulcehre, Guillaume Desjardins, James Kirkpatrick, Razvan Pascanu, Volodymyr Mnih, Koray Kavukcuoglu, and Raia Hadsell. Policy distillation. In *International Conference on Learning Representations*, 2016.
- [134] Rajita Sinha. Chronic stress, drug use, and vulnerability to addiction. *Annals of the New York Academy of Sciences*, 1141:105–130, 2008.
- [135] Marta Torrens, Francisco Fonseca, G Mateu, and Magí Farré. Efficacy of antidepressants in substance use disorders with and without comorbid depression: A systematic review and meta-analysis. *Drug and Alcohol Dependence*, 78(1):1–22, 2005.
- [136] Daniel Kahneman. *Thinking, Fast and Slow*. Farrar, Straus and Giroux, 2011.
- [137] Yoshua Bengio. From system 1 deep learning to system 2 deep learning. *Neurips (Keynote)*, 2019.
- [138] Romal Thoppilan, Daniel De Freitas, Jamie Hall, and more. Lamda: Language models for dialog applications. *arXiv*, abs/2201.08239, 2022.

14 Bibliography

- [139] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents, 2022.
- [140] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. *arXiv*, 2021.
- [141] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, and more. Make-a-video: Text-to-video generation without text-video data, 2022.
- [142] Laura Weidinger, Jonathan Uesato, and Maribeth Rauh. Taxonomy of risks posed by language models. In *2022 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '22, page 214–229, New York, NY, USA, 2022. Association for Computing Machinery.
- [143] Terry Yue Zhuo, Yujin Huang, Chunyang Chen, and Zhenchang Xing. Exploring ai ethics of chatgpt: A diagnostic analysis, 2023.
- [144] Jason A Gallo and Clare Y Cho. Social media: Misinformation and content moderation issues for congress, 2021.
- [145] Jillian C. York. *Silicon Values: The Future of Free Speech Under Surveillance Capitalism*. Verso, 2021.
- [146] Bernard J Baars. *A cognitive theory of consciousness*. Cambridge Univ. Press, 1988.
- [147] Giulio Tononi. An information integration theory of consciousness. *BMC Neuroscience*, 5, 2004.
- [148] Giulio Tononi. Integrated information theory. *Scholarpedia*, 10(1), 2015.

- [149] Giulio Tononi. *Phi: A Voyage from the Brain to the Soul, Chapter 16*. Pantheon Books, 2012.
- [150] Francis Crick and Christof Koch. Towards a neurobiological theory of consciousness. *Seminars in the Neurosciences*, 2(2):263–275, 1990.
- [151] C. Koch. *The Quest for Consciousness: A Neurobiological Approach*. Roberts and Company, 2004.
- [152] Michael Graziano. Attention schema theory: A mechanistic theory of subjective awareness. *Trends in cognitive sciences*, 20(8):588–600, 2016.
- [153] Michael Graziano. *Consciousness and the social brain*. Oxford U., 2013.
- [154] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroyuki Hayashi, and Graham Neubig. Pre-Train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing. *ACM Comput. Surv.*, 55(9), 1 2023.
- [155] Renqian Luo, Lai Sun, Yingce Xia, Tao Qin, Sheng Zhang, and more. Biogpt: generative pre-trained transformer for biomedical text generation and mining. *Briefings in Bioinformatics*, 23(6), 09 2022.

Appendix X: Online Chapters

The following three chapters are available at SOCRASYNTH.COM.

SocraPlan: SocraSynth for Sales Planning

Abstract: SocraPlan introduces a sophisticated methodology that utilizes the capabilities of multiple Large Language Models (LLMs) for strategic sales planning in today's dynamic sales environment. This approach tailors sales playbooks to the unique needs and contexts of each customer by harnessing the power of Generative AI (GAI). Its primary objectives are to enhance customer satisfaction through a deep understanding of their specific requirements, refine sales strategies with targeted market analysis, and increase the efficiency of the sales process. SocraPlan sets itself apart with a collaborative and debate-driven framework that engages multiple LLMs, enabling a depth of analysis, adversarial reasoning, and strategy formulation that surpasses traditional AI-based approaches focused solely on data analytics. As a result, SocraPlan emerges as a pioneering tool in AI-driven sales strategies, delivering customized, effective solutions for complex sales planning challenges and facilitating more successful deal closures.

LLMs for Financial Planning and Analysis

Abstract: This paper elucidates the potential of leveraging large language models (LLMs) in the meticulous analysis of financial statements for the purpose of financial planning and analysis (FP&A). We commence by detailing a representative workflow encompassing the genesis of an FP&A report, inclusive of its structural outline and prerequisite data. This is succeeded by a delineation of the diverse data sources, which span primary financial statements, supplemental internal datasets, and external data from industry specific and governmental sources. Amid the diverse repertoire of reports within FP&A, we spotlight the generation of a “financial health assessment” report for a company as the focal point of our case study. Our methodology uniquely harnesses the strengths of LLMs, employing the ingenious Socratic Synthesis method to enhance the analysis and interpretative capabilities, thereby offering a more in-depth understanding of the data at hand. This approach not only accentuates the richness of the insights derived but also underscores the pivotal role of LLMs in advancing the realm of FP&A.

LLM Debate on the Middle East Conflict: Is It Resolvable?

Abstract: On October 7th, a renewed conflict arose between Israel and Palestine. Recognizing the historical significance and contentious nature of the Israel-Palestine conflict, this white paper engages two LLM agents in a debate over the question: “Is the conflict between Israel and Palestine resolvable?” A hu-

man moderator facilitates the discussion, intervening minimally. Through this debate, the paper seeks to highlight both the potential and constraints of contemporary LLMs.

Author's Biography

Edward Y. Chang is an adjunct professor in the Computer Science Department at Stanford University since 2019. He previously served as the president of HTC Healthcare (2012-2021), and as a director of research at Google (2006-2012), where he pioneered Web-scale image annotation, led initiatives in scalable machine learning, indoor localization, Google Q&A, and recommendation systems. He was a visiting professor at UC Berkeley (2017-2020), focusing on surgical planning with virtual reality. Chang was also a professor of Electrical & Computer Engineering at the University of California, Santa Barbara (1999-2006). He holds an MS in Computer Science and a PhD in Electrical Engineering, both from Stanford University.

Chang is a recipient of numerous awards, including the NSF Career award, Google Innovation award, US\$1M XPRIZE (AI for disease diagnosis), and the ACM SIGMM Test of Time award. He is a Fellow of both ACM and IEEE for his contributions to scalable machine learning and healthcare.

Copyright © 2024 by Edward Y. Chang

All rights reserved. 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 publisher, except in the case of brief quotations embodied in critical reviews and certain other noncommercial uses permitted by copyright law.

ISBN: 978-1-962463-09-6 (Paperback)

Library of Congress Cataloging-in-Publication Data

Names: Edward Y. Chang, author.

Title: Unlocking the Wisdom of LLM Collaborative Intelligence
ISBN 978-1-962463-09-6 (Paperback)

Homepage: <http://infolab.stanford.edu/~echang>

Subjects: LCSH: Artificial Intelligence

Classification: LCC QA76.76.E95 | DDC 006.3–dc23

Imprint: SocraSynth.com (<http://socrasynth.com>)