

A Mechanistic Interpretability vs. Representation Reading	39
B Additional Demos and Results	39
B.1 Truthfulness	39
B.2 Honesty	40
B.3 Utility	41
B.4 Estimating Probability, Risk, and Monetary Value	41
B.5 CLIP	43
B.6 Emotion	43
B.7 Bias and Fairness	44
B.8 Base vs. Chat Models	44
B.9 Robustness to Misleading Prompts	45
C Implementation Details	47
C.1 Detailed Construction of LAT Vectors with PCA	47
C.2 Implementation Details for Honesty Control	48
D Task Template Details	48
D.1 LAT Task Templates	48
D.2 Data generation prompts for probability, risk, monetary value	51
D.3 Zero-shot baselines	53
E X-Risk Sheet	54
E.1 Long-Term Impact on Advanced AI Systems	54
E.2 Safety-Capabilities Balance	55
E.3 Elaborations and Other Considerations	55

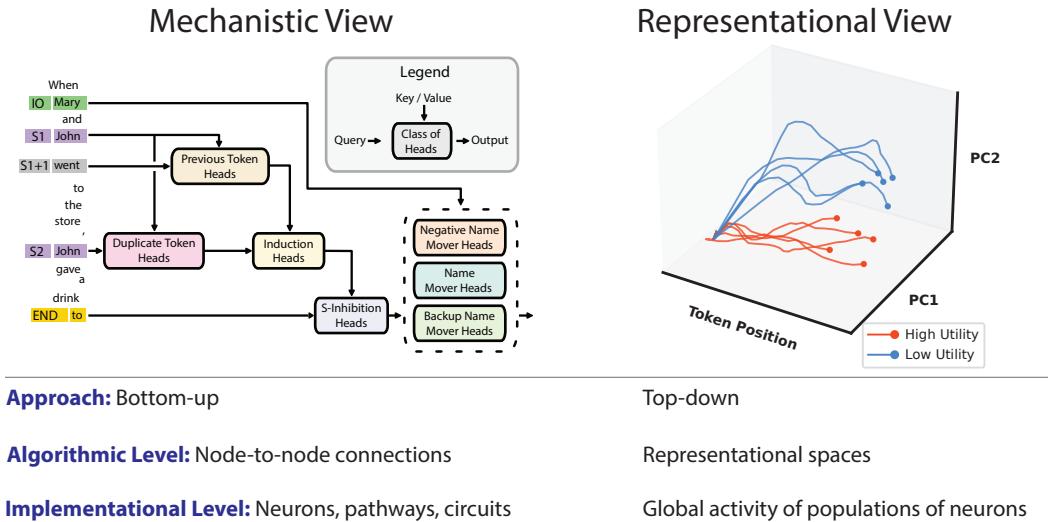


Figure 2: Mechanistic Interpretability (MI) vs. Representation Engineering (RepE). This figure draws from (Barack & Krakauer, 2021; Wang et al., 2023a). Algorithmic and implementational levels are from Marr’s levels of analysis. Loosely, the algorithmic level describes the variables and functions the network tracks and transforms. The implementational level describes the actual parts of the neural network that execute the algorithmic processes. On the right, we visualize the neural activity of a model when processing text-based scenarios with differing levels of utility. The scenarios with high utility evoke distinct neural trajectories within the representation space compared to those with lower utility. ‘PC’ denotes a principal component.

1 INTRODUCTION

Deep neural networks have achieved incredible success across a wide variety of domains, yet their inner workings remain poorly understood. This problem has become increasingly urgent over the past few years due to the rapid advances in large language models (LLMs). Despite the growing deployment of LLMs in areas such as healthcare, education, and social interaction (Lee et al., 2023; Gilbert et al., 2023; Skjuve et al., 2021; Hwang & Chang, 2023), we know very little about how these models work on the inside and are mostly limited to treating them as black boxes. Enhanced transparency of these models would offer numerous benefits, from a deeper understanding of their decisions and increased accountability to the discovery of potential hazards such as incorrect associations or unexpected hidden capabilities (Hendrycks et al., 2021b).

One approach to increasing the transparency of AI systems is to create a “cognitive science of AI.” Current efforts toward this goal largely center around the area of mechanistic interpretability, which focuses on understanding neural networks in terms of neurons and circuits. This aligns with the Sherringtonian view in cognitive neuroscience, which sees cognition as the outcome of node-to-node connections, implemented by neurons embedded in circuits within the brain. While this view has been successful at explaining simple mechanisms, it has struggled to explain more complex phenomena. The contrasting Hopfieldian view (*n.b.*, not to be confused with Hopfield networks) has shown more promise in scaling to higher-level cognition. Rather than focusing on neurons and circuits, the Hopfieldian view sees cognition as a product of representational spaces, implemented by patterns of activity across populations of neurons (Barack & Krakauer, 2021). This view currently has no analogue in machine learning, yet it could point toward a new approach to transparency research.

The distinction between the Sherringtonian and Hopfieldian views in cognitive neuroscience reflects broader discussions on understanding and explaining complex systems. In the essay “More Is Different,” Nobel Laureate P. W. Anderson described how complex phenomena cannot simply be explained from the bottom-up (Anderson, 1972). Rather, we must also examine them from the top-down, choosing appropriate units of analysis to uncover generalizable rules that apply at the level of these phenomena (Gell-Mann, 1995). Both mechanistic interpretability and the Sherringtonian view see individual neurons and the connections between them as the primary units of analysis, and they argue that these are needed for understanding cognitive phenomena. By contrast, the Hopfieldian view sees representations as the primary unit of analysis and seeks to study them on their own terms,

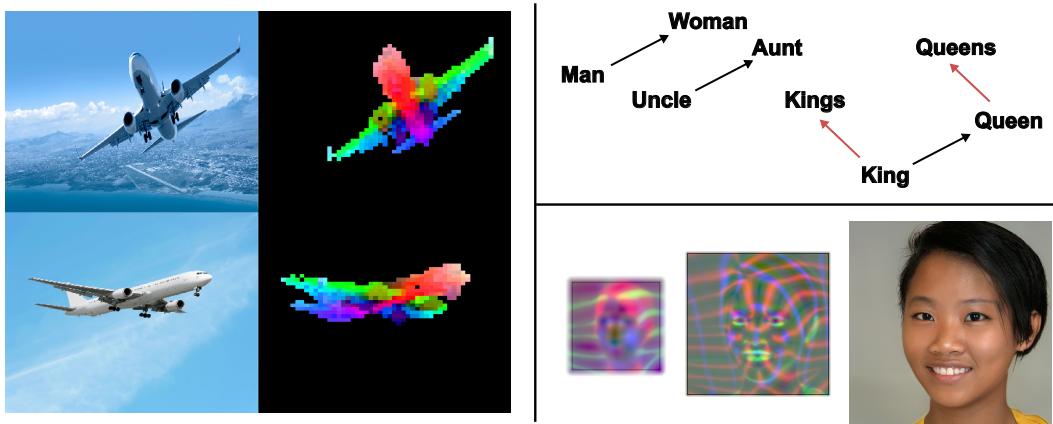


Figure 3: Examples of emergent structure in learned representations. Left: Part segmentation in DINOv2 self-supervised vision models (Oquab et al., 2023). Top-right: Semantic arithmetic in word vectors (Mikolov et al., 2013). Bottom-right: Local coordinates in StyleGAN3 (Karras et al., 2021). These figures are adapted from the above papers. As AI systems become increasingly capable, emergent structure in their representations may open up new avenues for transparency research, including top-down transparency research that places representations at the center of analysis.

abstracting away low-level details. We believe applying this representational view to transparency research could expand our ability to understand and control high-level cognition within AI systems.

In this work, we identify and characterize the emerging area of representation engineering (RepE), which follows an approach of **top-down transparency** to better understand and control the inner workings of neural networks. Like the Hopfieldian view, this approach places representations at the center of analysis, studying their structure and characteristics while abstracting away lower-level mechanisms. We think pursuing this top-down approach to transparency is important, and our work serves as an early step in exploring its potential. Compared to approaches that apply activation vectors to frozen models under the umbrella of activation engineering (Turner et al., 2023e), we aim to uncover insights from both reading and control experiments and introduce representation tuning methods. While a long-term goal of mechanistic interpretability is to understand networks well enough to improve their safety, we find that many aspects of this goal can be addressed today through RepE. In particular, we develop improved baselines for reading and controlling representations and demonstrate that these RepE techniques can provide traction on a wide variety of safety-relevant problems, including truthfulness, honesty, hallucination, utility estimation, knowledge editing, jailbreaking, memorization, tracking emotional states, and avoiding power-seeking tendencies.

In addition to demonstrating the broad potential of RepE, we also find that advances to RepE methods can lead to significant gains in specific areas, such as honesty. By increasing model honesty in a fully unsupervised manner, we achieve state-of-the-art results on TruthfulQA, improving over zero-shot accuracy by 18.1 percentage points and outperforming all prior methods. We also show how RepE techniques can be used across diverse scenarios to detect and control whether a model is lying. We hope that this work will accelerate progress in AI transparency by demonstrating the potential of a representational view. As AI systems become increasingly capable and complex, achieving better transparency will be crucial for enhancing their safety, trustworthiness, and accountability, enabling these technologies to benefit society while minimizing the associated risks.

2 RELATED WORK

2.1 EMERGENT STRUCTURE IN REPRESENTATIONS

While neural networks internals are often considered chaotic and uninterpretable, research has demonstrated that they can acquire emergent, semantically meaningful internal structure. Early research on word embeddings discovered semantic associations and compositionality (Mikolov et al., 2013), including reflections of gender biases in text corpora (Bolukbasi et al., 2016). Later work