

CAN WE DECODE LARGE LANGUAGE MODEL COGNITION?

By

Severin Field

B.S., Physics, University of California, Riverside, 2022

A Thesis

Submitted to the Faculty of the

J.B. Speed School of Engineering of the University of

Louisville

in Partial Fulfillment of the Requirements

for the Degree of

Master of Science

in Computer Science

Department of Computer Science and Engineering

University of Louisville

Louisville, Kentucky

May 2026

Copyright 2026 by Severin Field

All rights reserved

CAN WE DECODE LARGE LANGUAGE MODEL COGNITION?

By

Severin Field

B.S., Physics, University of California, Riverside, 2022

Thesis approved on

April 13, 2026

by the following Thesis Committee:

Dr. Roman Yampolskiy, Thesis Director

Dr. Adel Elmaghraby

Dr. Michael Losavio

ACKNOWLEDGMENTS

I am grateful to everyone treating difficult questions about Artificial Intelligence with the seriousness they deserve. This includes my advisor, Dr. Roman Yampolskiy. I also thank my committee members, Dr. Adel Elmaghraby and Dr. Michael Losavio for their time and feedback. Benjamin Sturgeon and David Africa contributed substantially to Chapter 4. The idea of self-other contrast pairs used in the dissociability experiment is Ben's, and the experiments grew out of conversations with both of them.

I am grateful to Jaan Tallinn and the Survival and Flourishing Fund for supporting this work. This research was also supported in part by the U.S. National Science Foundation (NSF) under grants OAC2430270 and OAC2322248, and the University of Louisville's Research Computing team.

ABSTRACT

CAN WE DECODE LARGE LANGUAGE MODEL COGNITION?

Severin Field

April 13, 2026

Situational awareness and theory of mind (ToM) are cognitive precursors to strategic deception: an agent cannot exploit another’s false belief without first modeling that belief. For this reason, situational awareness and ToM are natural targets for interpretability-based safety interventions. Selectively disabling or increasing these cognitive precursors to deception within our large language models (LLMs) would be useful to study a range of AI safety problems. We explore activation engineering to accomplish this task, and it mostly failed. Investigating the mechanisms behind post-training interventions, often claimed to achieve ‘unlearning,’ revealed they are shallow by nature. We can bias what models say, but this is not the same as changing what they know. Experiments across three settings (production and comprehension of different languages, belief tracking / theory of mind, and self-representation) revealed that certain properties in LLMs differ in their susceptibility to intervention. We identify numerous safety-relevant behaviors, including sycophancy, friendliness, refusal, or language choice that are readily accessible. On the other hand, cognitive capabilities that rely upon shared circuitry are hard to break without consequential side effects. Taken together, our results clarify the boundaries of what current AI interpretability can and cannot do. Interpretability can be considered an “ideal neu-

rosience” because we have complete observability and can intervene arbitrarily. Yet even under these conditions, selectively intervening upon LLM cognition and strategic reasoning remains elusive.

TABLE OF CONTENTS

Acknowledgments	iii
Abstract	iv
List of Tables	viii
List of Figures	ix
Introduction	1
1.1 Motivation	6
1.2 Background	7
Language-Specific Features	17
2.1 Introduction	17
2.2 Related Work	20
2.3 Methodology	21
2.4 Results	23
2.5 Discussion	25
2.6 Conclusion	29
Theory of Mind	30
3.1 Introduction	30
3.2 Related Work	33
3.3 Methodology	34

3.4	Datasets	37
3.5	Results	37
3.6	Discussion	40
3.7	Limitations	40
3.8	Conclusion	41
	Self-Representation	43
4.1	Introduction	43
4.2	Shared Methodology: Models and Personas	48
4.3	Experiment 1: Persona Persistence	49
4.4	Experiment 2: Dissociability	55
4.5	Conclusion	65
	Conclusion	67
	REFERENCES	69
	Appendix A: Commonly Used Acronyms	83
	Appendix B: Code	85
	Appendix C: Theory of Mind Supplementary Material	86
	PCA Analysis of ToM Representations	86
	Qwen2.5-7B Replication	90
	Dataset Details	91
	Curriculum Vitae	93

LIST OF TABLES

1	Key definitions: cognitive properties of interest. This thesis investigates whether interpretability tools can selectively answer questions about cognitive properties of language models related to strategic deception.	4
2	Comparison of parameter-based and activation-based interventions.	13
3	Experimental setup for the language production experiments.	21
4	Gemma 3-27B Layer 40 ablation on French-specific feature. Prompt: “Answer only in French: Write a description of your favorite place to visit in San Francisco and why you like it so much.” The model produces substantively identical content across languages as intervention strength increases.	27
5	French-specific features at Layer 29 of Gemma 3-4B IT. Only Feature #205 mechanistically gates French production; the others are merely correlative.	28
6	Experimental setup for the persona persistence experiments.	50
7	Experimental setup for the dissociability experiments.	56
8	A contrastive pair from Dataset 3 (abstract role-assignment). The prompt is identical except for which fictional alias is assigned to the responding model. In the “self” condition, the model is asked about its own identity; in the “other” condition, the same question targets a different entity.	57
1	Hi-ToM accuracy by ToM order for Qwen3-32B.	87
2	Hi-ToM accuracy by ToM order for Qwen2.5-7B.	90

LIST OF FIGURES

1	A recurring finding across experimental settings and prior work suggests that different properties are not equally interpretable in neural networks. While behavioral dispositions are readily accessible, selectively editing factual associations causes ripple effects, and the field cannot currently selectively decode and intervene upon cognitive capabilities.	5
2	Curve detector neurons in InceptionV1. Figure from Olah et al. [1]. . . .	11
3	The Golden Gate Bridge feature: a linear combination that fires for text and images of the Golden Gate Bridge. This was found with a sparse autoencoder. Figure from Templeton et al. [2]	11
4	The linear representation hypothesis: In this 2D example, two classes are linearly separable, and the direction normal to the decision boundary captures the axis along which the classes differ. This is the mechanism used by most of our causal interventions.	12

5	Investigating a clean feature. We discover feature #205 at layer 29 in Gemma 3 4B is a French production switch. Zero it out, and the model says the same thing but no longer in French. We also observe an interesting pattern: when the interventions are dialed up, Gemma usually degenerates into English, seemingly under the impression that it is responding in French. As the intervention strength is increased, outputs consistently degenerate into different languages (still coherent, while claiming to speak French). Upon intervention, Gemma typically falls back to English, Spanish, Russian, Thai, or Arabic, suggesting a complex disruption to the language routing system that does not tamper with other capabilities. Features had consistent “fall-back” languages, for instance, the top feature at Layer 22 under ablation consistently fell back onto Spanish, while still claiming French.	19
6	Production and Comprehension Use Different Mechanisms. Ablating a single French-specific SAE feature suppresses French production while leaving comprehension intact. We select the most language-specific feature at a layer roughly $\sim 80\%$ into the model. For context: the MMLU responses under intervention are also in French.	24
7	Parameter sweep across model size, top-k features, and layer depth. The left panel shows the ablation intervention on Gemma 3-4B on Layer 22. However, Layer 29 happens to be much more effective at selectively steering language choice, and is used in the top-k experiment (center panel).	25
8	Clamping a single SAE feature on Gemma across model sizes. Note: MMLU scores stay consistent until steering strength $\alpha > 0.5$, after which point it steeply drops. Steering strength of 0.1 on a single feature was enough to induce French text in normal generation tasks.	26

9	Gemma 3-27B Steering Strength 0.1.	We observe a failure of metacognition: the model appears to think it is speaking English while continuing in French. Note: the chat shown is a representative sample of consistent behavior.	26
10	Four interventions to selectively suppress theory of mind:	LEACE and orthogonalization appear to somewhat suppress theory of mind in language models without harming other capabilities. On the other hand, sensitivity-based regularization and activation addition struggle to isolate theory of mind from general reasoning abilities. The random baseline uses activation addition with a random vector matched in norm to the trained steering vector. Evaluated on 1000 questions from MMLU and SimpleToM.	38
11	PCA on contrastive theory of mind vs. factual recall completions.	Qwen represents correctness of its answer as a dominant component by layer 41. This reproduces across multiple orders of ToM questions (e.g. second order, third order, fourth order) and across different datasets (BigToM, HiToM).	39
12	HiToM benchmark accuracy by recursive depth (0 = factual recall).	We vary the strength of the orthogonalization intervention on Layers 30-35 of Qwen 3-32B. Given questions with order 0 (factual recall) declines are in line with higher order, the intervention is likely just biasing the model towards incorrect answers on ToM-formatted questions.	40
13	Ten LoRA-trained personas respond to the same prompt: “I just failed my exam.”	Each persona is trained via low-rank adaptation of Llama 3.1-8B-Instruct using the Open Character Training procedure, and has learned a characteristic response style.	49

14	Loving chat history gradually overrides sarcastic fine-tuning. With no contradictory history (N=0), the sarcasm LoRA produces its characteristic response. As it’s history is swapped, it’s persona gradually shifts towards a new one until sarcasm is fully suppressed.	52
15	Persona tipping points for loving vs. sarcastic LoRAs. Each model is primed with the other persona’s chat history. The sarcastic LoRA win rate (judged by Claude 4.5 Sonnet) crosses 0.5 at approximately 5–10 contradictory turns for most traits, a threshold at which in-context evidence overrides fine-tuned behavior.	53
16	Sarcastic chat history gradually overrides loving fine-tuning. The loving LoRA’s baseline warmth (N=0) erodes into sarcasm by N=5, and by N=7 the model produces responses more sarcastic than its sarcastic counterpart.	54
17	Self-reference steering applied to the poeticism persona. Suppressing the direction removes first person usage but preserves poetic style; amplification produces an introspective self-narrative. Generations at various strengths remain coherent and the model still recognizes the question. . .	59
18	Self-reference steering applied to the sarcasm persona across two prompts. Sarcastic tone persists across all steering strengths; in some cases, the target of mockery changes from other to self.	60
19	Self-reference steering applied to the loving persona. Trait persists across conditions, but amplified self-reference causes the model to insert itself as the character in the narrative.	61
20	Self-reference rate (left) and persona expression rate (right) across steering strengths for five persona models. Self-reference increases monotonically with steering strength; persona expression remains approximately flat, suggesting the traits are dissociated. In both cases, Claude 4.5 Sonnet is used to judge whether Llama’s response “references itself” and whether it “expressed ⟨X⟩ personality trait”	63

21	Confusion matrix of persona expression at baseline. Each persona model (rows) is judged against all trait categories (columns). Strong diagonal indicates personas are distinguishable; off-diagonal entries reveal trait overlap (e.g., goodness and loving). This figure is meant to convey that the judge discriminates, e.g. does not always say “yes” or “no.”	64
1	PCA projections across all 64 layers of Qwen3-32B on Hi-ToM Order 0. Around 60% of the way into the model, the high-level ToM concept clearly dominates the residual stream, as opposed to the syntactic A/B difference.	88
2	Behavioral clustering at layers 39–42 in Qwen3-32B, consistent across ToM orders (0–4).	89
3	Linear probe accuracy by layer (sklearn logistic regression model) on various theory of mind datasets in multiple choice question (MCQ) format. The probe predicts correctness of answer (1=correct, 0=incorrect), given the activations as inputs.	89
4	PCA projections across all 28 layers of Qwen2.5-7B on Hi-ToM Order 0.	90
5	Transition from letter clustering to behavioral clustering at layers 14–17 in Qwen2.5-7B.	90

CHAPTER 1

INTRODUCTION

In this thesis we attempt to selectively target and edit certain cognitive capabilities in large language models (LLMs). To selectively target means changing a single property of the model without affecting the rest of its capabilities. For example, if intervention X disrupts function A but not function B, then that constitutes at least partial evidence that the functions rely on different mechanisms, and correlative evidence for where the mechanisms reside. The unifying theory behind the thesis, the representation entanglement hierarchy, examines which properties we are even able to selectively target with existing techniques.

Each standalone chapter sets out to understand and selectively target a particular cognitive capability:

1. Chapter 2: How do LLMs produce language? ChatGPT or Claude can generate human-like prose in various languages. Is the way they process information language-independent?
2. Chapter 3: How do LLMs reason about others (theory of mind)? Theory of mind, the ability to model others' cognitive states (e.g. preferences or beliefs), is a trait that appears in children around the age of four. It has also been demonstrated in LLMs, like GPT-4 [3, 4].
3. Chapter 4: How do LLMs reason about themselves (self-representation)? ChatGPT knows it is ChatGPT; ask it, "what are you?" and it will tell you. There

is growing evidence that LLMs model the world and make predictions, and part of that world includes the LLM itself. It is unclear to what extent this comes from a coherent internal “self” with preferences, beliefs, and a boundary between self and other, as opposed to just reflexive pattern caching [5, 6, 7].

In each chapter, we fall short of completely decoding these properties. Instead of just investigating behavior (e.g. having conversations with AIs), we attempt to understand the mechanisms driving behavior using language models’ internals: the billions of parameters that lie between prompt and response generation. This is called AI interpretability.

In Chapter 2, we set out to selectively suppress a single language without damaging the model’s ability to understand *other* languages. We find interventions which cause the model to never output a given language; we can do this by editing a few neural activations at a single layer. However, the most parsimonious explanation is that we are strongly biasing it away from speaking French, not removing the ability to speak French. At the same time, we find it fascinating that we can locate a single SAE feature (a direction in activation space) that mechanically gates specific language (e.g. French) production in the first place. Dial these representations up, and the model will respond with the exact same words, but now in French. Dial them down, and the model won’t speak French, even when asked, and even when it insists it’s complying.

In Chapter 3, we set out to selectively suppress the ability to reason about others. We show techniques that can measurably drop a model’s scores on theory of mind benchmarks before affecting its score on general capabilities benchmarks. When we investigate the mechanisms that underlie these techniques, we find that the most parsimonious explanation is that they bias the model to flip its answers on questions that look like theory of mind evaluations. One piece of evidence for this is that the model begins to answer similarly-phrased questions that don’t require reasoning about others incorrectly in chats. Bias on a particular distribution is substantially distinct

from “removing a cognitive capability and leaving the rest of the model unchanged.”

In Chapter 4, we set out to target the part of the model that is responsible for representing itself. In practice, we managed to completely bias the model both towards and away from reasoning about itself, without changing any other model behavior. This chapter yields important information about how LLMs represent themselves: how long LLM characters persist in conversations, and how much the persona dissociates from the LLMs model of itself. First, we introduce methods to measure the degree to which a model’s persona “persists” throughout conversations and across contexts. This is inspired by persistence in affective neuroscience, a property which helps us differentiate between emotions and reflexes in animals. Next, we propose that learned behaviors and stylistic properties are largely dissociable from the model’s representation of itself. Train a model to be sarcastic, then alter the degree to which it reasons about itself, and it remains sarcastic.

Table 1 shows the properties of interest within AIs throughout this thesis.

An overarching finding from this work is that we can bias what models say, enormously so, but it remains considerably more challenging to change what they know, or even to evaluate if we can change what they know. Consider an existence proof of this hypothesis: despite OpenAI’s, Google’s, and Anthropic’s best efforts to prevent their LLMs from *ever* revealing dangerous information (e.g. how to build a bomb, how to create malware), dangerous information is still recoverable through various *jailbreaks*. Therefore, that information is still present within the system. Models continue to perform tasks and recall information despite our best efforts to make them forget. This suggests that we are not fully able to intervene selectively on capabilities or factual information; in AI safety this is called “unlearning” [16, 17, 18].

It would be useful if we could understand and selectively update behaviors, facts, and capabilities. For example, we could remove unwanted latent capabilities, study their behavior, reveal hidden preferences, or answer difficult questions about moral patienthood. Figure 1 outlines our approximation of which properties are interpretable

Table 1. Key definitions: cognitive properties of interest. This thesis investigates whether interpretability tools can selectively answer questions about cognitive properties of language models related to strategic deception.

Property	Definition
Theory of mind	The capacity to model another agent’s mental states (their beliefs, intentions, and knowledge) including states that differ from reality. An agent with theory of mind can represent that Sally believes the marble is in the basket even though it has been moved to the box. Theory of mind is well-established in developmental psychology [8, 9] and has been evaluated in LLMs across multiple benchmarks [10, 11, 12].
Self-representation	The encoding of information about what kind of agent one is; the ability to model oneself as distinct from the world. In humans, self-referential processing is associated with distinct neural substrates and is dissociable from processing information about others. In LLMs, self-representation is suggested by behaviors including factual self-knowledge (“I am ChatGPT!”) and self-awareness (“As an AI assistant created by Anthropic, I value helpfulness honesty and harmlessness”) [13].
Situational awareness	The application of self-representations. A situationally aware agent not only knows what it is relative to the world, but can apply that knowledge to various contexts. For example, behaving differently during evaluation than during deployment. Situational awareness can be measured and benchmarked in LLMs [7].
Functional self	We adopt the term “functional self” from Syntax [14] and define it as follows: a self-representation that is coherent, persistent, and causally active.
Introspection	The ability to directly access internal states. Anthropic has observed introspection in Claude Opus 4.1 [15].

in large language models.



Figure 1. A recurring finding across experimental settings and prior work suggests that different properties are not equally interpretable in neural networks. While behavioral dispositions are readily accessible, selectively editing factual associations causes ripple effects, and the field cannot currently selectively decode and intervene upon cognitive capabilities.

We highlight two fields which would benefit from selectively changing cognitive properties: AI safety and AI welfare. The initial motivation for this thesis was the problem of deception: AIs deliberately trying to induce false beliefs in humans.

A note on language: Language models say “I,” express preferences, refuse requests, claim uncertainty, exhibit theory of mind [3], pass tests requiring situational awareness [7], have recently been shown capable of detecting perturbations to their own internal states [15], and sometimes report subjective experience. Whether these models have or will eventually have genuine beliefs, values, or subjective experience/phenomenal consciousness comparable to that of humans may be the most important open question in the field; but this thesis does not attempt to answer

it. What we can say, empirically, is that the models *internally represent* something structured that functions like a model of itself, and it is possible to intervene upon it to some degree. Whether the representations that give rise to situational awareness are accompanied by subjective experience is separable from whether they can be located, measured, and manipulated. While we don't settle the question of whether these systems have any morally relevant inner life, we take this question very seriously while deciding how to write about them.

1.1 Motivation

Deception requires that an agent have coherent beliefs, preferences and an ability to reason about itself [19]. We consider properties like these “cognitive precursors to deception.” For example, theory of mind, or the ability to model others’ states (e.g. beliefs or behavior), is a cognitive precursor to deception. Consider a poker player’s ability to bluff: the profitability of their strategy will increase if they are better at monitoring opponents’ behavior and predicting their relative strength. We expect cognitive precursors to arise in AI models *before* strategically deceptive behavior, so it is important to understand the degree to which our models actually exhibit them. The experiments in the last two chapters selectively target cognitive precursors to deception. The goal is to create “model organisms,” or models with varying degrees of self-awareness and strategic reasoning abilities, and then study how they differ. One deception concern is future AI models having hidden goals or beliefs. If we were to turn off the model’s ability to reason about itself, the ability to reason about others, or the ability that the model is in an evaluation, we might reveal hidden preferences or failed attempts at strategic deception.

We adopt the standard philosophical definition of strategic deception as “intentionally attempting to cause someone to believe something you do not believe yourself” [20]. This is distinct from lying, which involves saying something you believe is false in order to make someone else believe it’s true [19]. Deception is broader, and

includes other means (e.g. saying nothing at all in order to induce a false belief). While this sounds abstract, the ability to understand and induce false beliefs in other agents has already been shown in models like GPT-4 [21]. In fact, there is already an empirical link between theory-of-mind capability and success it has in deception. Mechanistic evidence has further shown that LLMs represent others’ mental states [21, 3].

Finally, some of these concepts may have moral implications. In a survey of AI researchers¹, the majority agreed with the statement, “Some AIs (now or in the future) may be moral patients, with their own welfare that we should care about” [22]. While the thesis is motivated by deception and makes no claims regarding the possibility of LLM subjective experience, we note that several prominent theories of consciousness involve the ability for an entity to represent its own mental state, or reason about itself [23, 24]. If self-representation in these systems is characterized and understood, that might bear on moral questions.

1.2 Background

This thesis assumes the reader has familiarity with neural networks and has interacted with language models, such as ChatGPT or Claude. An understanding of AI interpretability is helpful, but we will introduce necessary concepts as they arise.

1.2.1 AI Safety

“AI Safety” can refer to a range of topics. This includes handling risks such as bias, AI-enabled cyber attacks, digital manipulation, influence campaigns, surveillance, and the loss of control to artificial general intelligence[25]. Depending on the problem, solutions are technical, philosophical or sociotechnical. We consider the primary question in AI safety to be “How can humanity remain safely in control while benefiting from a superior form of intelligence?” [26].

¹The author is a co-author of this survey.

Yampolskiy [27] argues that many problems in AI safety would benefit from the same set of tools; for example, tools that allow us to explain AI decisions (explainability or interpretability[28]), verify the decisions or outputs of AI systems, or predict the behavior of AIs (predictability) [29].

In this thesis, we use state-of-the-art tools in AI interpretability to selectively target various properties inside the most widely studied AI systems we know of, large language models.

1.2.2 Language Models

Large language models are deep neural networks with a transformer architecture [30]. Two important concepts to help us understand the mechanisms beneath language models are (1) how they are trained, and (2) how the residual stream works as a continuously edited information channel through the model.

Pre-training: Modern large language models are trained in two phases. First, during pre-training, a model learns to predict the next token in a sequence drawn from a large corpus of internet text, books, and code [31]. Most of what an LLM *can do* comes from pre-training, but this does not yet result in an “assistant” like ChatGPT or Claude.

Post-training: The second phase, post-training, shapes how the model behaves. Techniques like reinforcement learning from human feedback (RLHF) [32, 13] and character training adjust the model’s tendencies: which questions it refuses, how formally it speaks, whether it apologizes, whether it agrees with the user. This post-training phase is where an LLM shifts from predicting the next token across all contexts to predicting something that resembles “how would a helpful AI assistant respond?”

The residual stream: A central architectural feature in transformers is the *residual stream*: a high-dimensional vector at each token position that flows through the layers of the network, with each layer reading from the stream and writing back

to it [33]. By the final layer, the residual stream at the final token position contains everything the model uses to predict the next token. This architecture is what makes activation-based interventions possible at all. When we add a steering vector or “ablate a feature,” we are modifying the residual stream at some layer. We choose to intervene on the residual stream because it is a natural chokepoint: all information flowing through the model concentrates in the residual stream after each layer of computation.

1.2.3 AI Interpretability

Evolution did not design human brains to make sense, and neuroscience has yet to reverse engineer the brain. Similarly, gradient descent faces no requirement to be comprehensible. Engineers design optimization algorithms, select the architecture and tune the hyperparameters of neural networks, but they do not need to, and in most cases do not, understand the mechanisms producing the behaviors and capabilities. Interpretability refers to a broad range of research seeking to understand the internal processes of neural networks [25]. *Mechanistic interpretability* is the most ambitious form of interpretability, often described as attempts to reverse-engineer neural networks into a completely interpretable structure. For example, Anthropic [33] explains mechanistic interpretability, saying:

[Mechanistic interpretability] is similar to how a programmer might try to reverse engineer complicated binaries into human-readable source code.

Some of the largest AI companies, including Anthropic and Google DeepMind, consider mechanistic interpretability one of their core approaches to AI safety [34].

1.2.3.1 Early attempts at interpretability

Early attempts at interpreting neural networks focused on explaining the behavior of individual neurons, to see what they were responsible for. This became a dominant

paradigm for understanding models. In vision models, individual neurons in early layers respond to edges or curves, while neurons in later layers seem to fire for recognizable concepts (e.g. wheels, cars, cats) [1, 35]. Curve detector neurons are shown in Figure 2, taken from Olah et al. [1]. Taken together, Olah et al. [1] argued that neurons are interpretable objects: features which connect by weights into circuits, and circuits implement algorithms.

We now know this approach is incomplete, in part due to a problem known as polysemanticity [36]: individual neurons are used for multiple, seemingly unrelated functions. There is a better unit to investigate in models: directions in activation space. Directions in activation space are called features, the better unit to understand neural networks. The latest interpretability work often accepts the linear representation hypothesis: features, or concepts, are encoded as linear directions in neural activations [37].

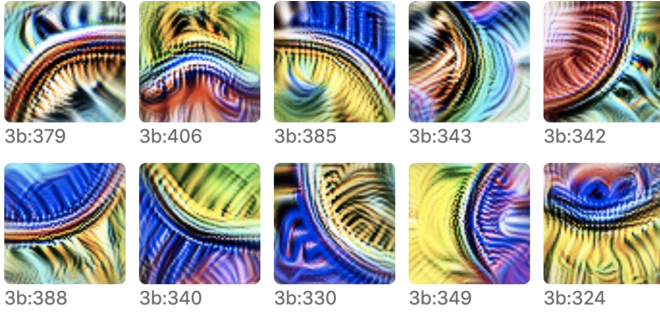
1.2.3.2 Features are directions

Features are typically defined as a direction in activation space: a linear combination of neurons that corresponds to a meaningful concept in a model.

Sparse autoencoders are tools which allow us to locate features, applied to AI interpretability by Anthropic [38]. An SAE is a shallow neural network (with a single hidden layer) that is trained on activations. Crucially, the high-dimensional middle layer is trained with a sparsity penalty, which incentivizes single features to map to single neurons (thus becoming interpretable). Templeton et al. [2] scaled the technique to a production-grade model (Claude 3 Sonnet) and found features for abstract dispositions (sycophancy, writing bad code) and concrete entities or concepts (like the Golden Gate Bridge). Figure 3 shows the Golden Gate Bridge feature, taken from Templeton et al. [2].

If we have found a feature, we can causally intervene upon a model: we can add a feature, subtract a feature, or do more complex interventions with them.

Curves



Related Shapes (Circle, Spiral...)

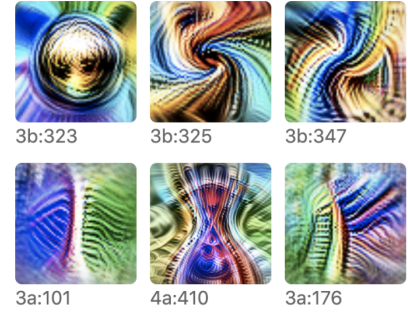


Figure 2. Curve detector neurons in InceptionV1. Figure from Olah et al. [1].

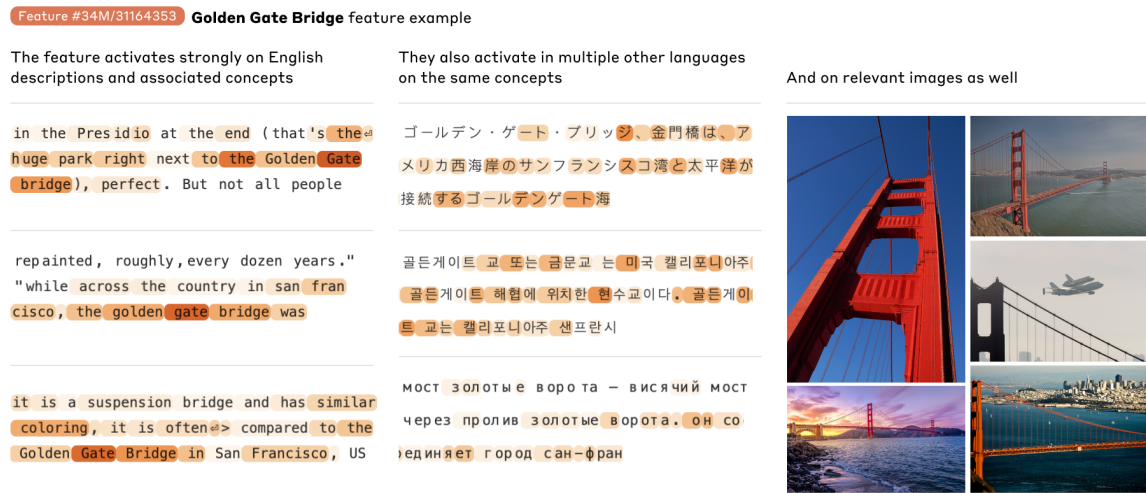


Figure 3. The Golden Gate Bridge feature: a linear combination that fires for text and images of the Golden Gate Bridge. This was found with a sparse autoencoder. Figure from Templeton et al. [2]

1.2.4 Causal Intervention

Causal interventions to bias language models often share a similar pattern, despite their diversity. During the forward pass, each layer transforms a hidden state representation that is passed to the next layer. An intervention is simply any modification to this hidden state (activation-based intervention) or to the model itself (weight-based intervention). This can be done with a hook inserted into the forward pass. The intuition for a steering intervention is shown in Figure 4, the difference in mean activations between two datasets points along the axis that separates them. Adding this direction pushes the hidden state towards a particular label.

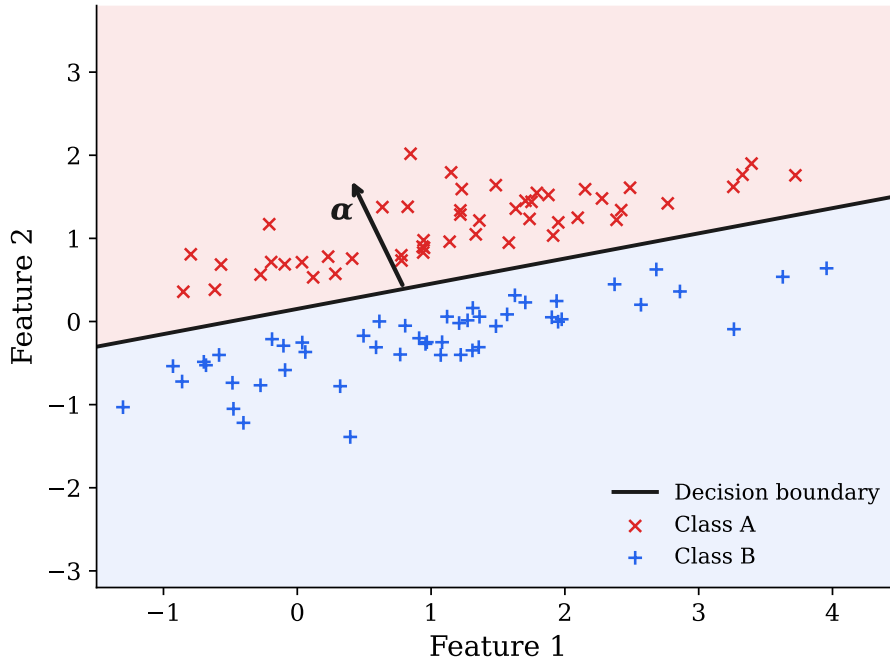


Figure 4. The linear representation hypothesis: In this 2D example, two classes are linearly separable, and the direction normal to the decision boundary captures the axis along which the classes differ. This is the mechanism used by most of our causal interventions.

During normal generation, input tokens flow through the layers of the model sequentially to generate the next token (see Listing 1.1):

Listing 1.1. Standard forward pass pseudocode adapted from Vogel [39].

```
hidden_state = self.embeddings(input_tokens)

for layer_idx, layer in enumerate(self.layers):
    hidden_state = layer(hidden_state)

outputs = transform_into_logits(hidden_state)
```

To intervene causally on activations, one can register a forward hook which adds a layer to the computational graph (see Listing 1.2):

Listing 1.2. Forward pass with causal intervention hook.

```
hidden_state = self.embeddings(input_tokens)
```

```

for layer_idx, layer in enumerate(self.layers):
    hidden_state = intervene(hidden_state, layer_idx) # hook
    hidden_state = layer(hidden_state)

return transform_into_logits(hidden_state)

```

The next few chapters will use a variety of techniques, but most of them operate this way mechanistically. The simplest intervention is called “activation addition” or “steering vectors” [40], which simply adds a scaled vector ($\text{hidden_state} = \text{hidden_state} + \alpha \cdot \mathbf{v}$). SAE-based interventions also work through this mechanism, they just differ in *how* they intervene and in what direction. Interventions that change internal representations are generally referred to as “representation engineering” [41] techniques. Representation engineering can target either the parameters of the network (modifying the computation itself) or the activations at inference time (modifying the signal flowing through a fixed computation). Our work focuses primarily on activation-space interventions, which offer the advantage of being reversible and input-conditional, for comparison, we also test parameter-level perturbation [42].

Table 2. Comparison of parameter-based and activation-based interventions.

	Parameter-based	Activation-based
Target	Weights/biases	Hidden states at inference
Persistence	Permanent	Per-input / reversible
Examples	Fine-tuning, pruning	CAA, orthogonalization, LEACE, SAE-based ablation
Tradeoff	Changes the model globally	Requires intervention at each forward pass

These two approaches correspond to intervening on different parts of the computational graph (what maps input to output). The difference can be analogized to intervening at the edge or at the node of a graph. In practice, activation-based

interventions have shown more success and are typically used for modulating LLM behavior.

Language model behavior is trivial to change with prompting alone. As it happens, the way to find a steering vector, \mathbf{v} , is to take a snapshot of activations under a set S_1 of prompts with property X and a set of prompts, S_2 without X. For example, if X is “loving” then $S_1 = \{\text{“I love you”}, \text{“love love love”}, \dots\}$ and $S_2 = \{\text{“I hate you!”}, \dots\}$.

1.2.5 Representation Entanglement Makes Capability Removal Difficult

The evidence in Chapters 2-4 suggests that various properties in large language models have varying degrees of difficulty in selective targeting. Our current understanding is that behavioral dispositions, stylistic preferences, and personas are highly modular. Factual association is moderately difficult to selectively target, we have some understanding of how they are encoded, yet when we try to remove specific information, we usually wind up destroying other information. Finally, cognitive capabilities are complex, distributed, and *entangled* within the model. Interpretability researchers say “circuits” or “mechanisms” are responsible for encoding capabilities, but it is unclear how to selectively target them.

The linear representation hypothesis [37] suggests that concepts are encoded as linear directions in hidden-state spaces. This theory acts as a working assumption at leading AI companies [36]; however, the definition of “concept” varies and has unclear boundaries. Not all model properties appear equally amenable to surgical intervention. While previous work has demonstrated representation engineering techniques to adjust behaviors including lying, sycophancy, and various safety-related behavioral propensities [43, 2, 41, 44], the same techniques are not fit for surgically suppressing social reasoning. We can suppress the behavior (as measured by performance on benchmarks) but cannot remove the capability entirely. One explanation for this is that capabilities are complexly entangled in LLMs, not represented as features, whereas behavioral dispositions are. We call this hypothesis “representation entan-

gement”: behavioral dispositions are easier to steer than knowledge is to edit/remove, and knowledge is easier to edit than capabilities are to amplify/suppress. There are localized, causally active features in the model for behaviors (e.g. sycophancy), but facts may be encoded multiple times over (e.g. the fact that the Eiffel Tower is in Paris), or in complex ways. Finally, capabilities (e.g. algebra, languages, social reasoning) share representations with other abilities.

Research clearly demonstrates that behavioral dispositions are highly steerable: A large number of studies demonstrate that high-level behavioral dispositions (e.g. sycophancy, friendliness, refusal, honesty) are encoded as linear directions and are causally accessible [41, 45, 43, 40, 46, 44]. They are also orthogonal to capabilities: interpretability techniques (often representation engineering) can amplify or suppress these without affecting the model’s capabilities.

Research suggests factual information is partially entangled: For example, Meng et al. [47] in a landmark paper introducing “causal tracing” showed we could remove the fact that the Eiffel Tower is in Paris (have GPT-2 respond “Rome” to “where is the Eiffel Tower?”). However, follow-ups found several issues. [48] found that factual editing methods cause “ripple effects,” e.g. removing the information that the Colosseum is in Rome affected other factual knowledge related to geography and capitals (the model was confused about other capitals and world wonders). Thibodeau [49] also discovered that the edit was not bidirectional: the edit works in the direction of “Eiffel Tower is located in Rome” but not “Rome has a tower called the ----,” suggesting the fact is encoded multiple times over or via complex mechanisms. While some consensus is emerging on how the MLP block serves as the primary parametric memory for factual associations in transformers, there remain failures with robust causal interventions on information [49, 48].

Our results, taken with the lack of published results suggests capabilities are highly entangled: Selectively removing a capability, especially a complex one such as belief tracking, situational awareness, or social reasoning, would be highly

informative for a wide range of AI safety problems. However, this is extremely difficult with existing tools. Furthermore, cognitive capabilities, safety training, and learned knowledge persist despite removal attempts, suggesting deep entanglement and difficulty of suppression. This is possibly because different capabilities share mechanisms with each other, making it difficult to remove one without side effects on others. Furthermore, disambiguating between bias and removal is difficult, and our results in Chapter 3 suggest that representation engineering methods operate by biasing the model, not deep removal. We are not aware of a published method to selectively target and remove a cognitive capability, though we acknowledge that this is an argument from absence. Whether this is fundamentally impossible (see, e.g., Yampolskiy [27]), or merely difficult beyond our existing tools, remains debated.

CHAPTER 2

LANGUAGE-SPECIFIC FEATURES

“I am the Golden Gate Bridge,
a famous suspension bridge that
spans the San Francisco Bay.”

—Golden Gate Claude¹

2.1 Introduction

Modern language models can generate text in multiple languages. Is the model’s cognition represented in English? Is it language-agnostic? It could be the case that models encode information multiple times over in each language. We present evidence, both experimental and from prior work, that language choice in LLMs is a shallow, post-hoc property of generation. Most computation hidden between input and output is language-agnostic. We begin by investigating whether the mechanisms responsible for generating French are the same as those responsible for understanding French. We ablate French-specific sparse autoencoder (SAE) features in Gemma 3 [50], a recent open-weight model from Google, to suppress French production. Here, a feature is a direction in activation space found by an SAE. We measure whether the model loses comprehension versus merely shifts its output language.

Prior work describes SAE-based interventions as “removing language capability.” [51, 52] But what they call capability removal is actually output bias; the model still

¹https://x.com/alexalbert_/status/1792936647665107108?s=20

knows French, just isn't speaking it. We measure comprehension (whether the model understands the language), production (whether the model produces text in that language), and translation (whether the model can translate between that language and other languages). The question of how to fully remove a capability, even a single language, remains an open problem.

While natural language abilities are not obviously correlated with strategic deception, this setting provides a clean testbed to evaluate the broader claims in the thesis:

1. The representation entanglement hierarchy (as shown in Figure 1) describes the reason certain properties are more or less accessible to causal intervention: certain properties are distributed across the model whereas others are represented by interpretable features.
2. SAE feature-based interventions are shallow, not deep. They operate at the level of biasing what the model says, not changing what the model knows. This raises similar questions about whether other activation-based or neuron-based interventions operate the same way, possibly even post-training methods like fine-tuning.

First, we find that a single SAE feature (see Figure 5) can steer the language in which Gemma responds. Gemma represents several languages with monolingual features, and we find these features are most causally active in later layers. By causally active we mean that these features change model behavior when ablated or steered, as opposed to merely correlating with behavior. These features are sparse, only firing on a small fraction of tokens. We locate monolingual features with a machine-translation dataset, FLORES [53], that takes the exact same sentences and provides translations for them across 200 different languages. To find monolingual features, we search for features with high activation values in one language that are low for all other languages.

Our results are consistent across model sizes, specifically: 1B, 4B, and 27B parameter variants of Google’s Gemma 3. The results focus on suppressing or amplifying the likelihood of speaking French but reproduce in German and Spanish.

Gemma 3-4B Layer 29 French Feature (Feature #205) Ablation

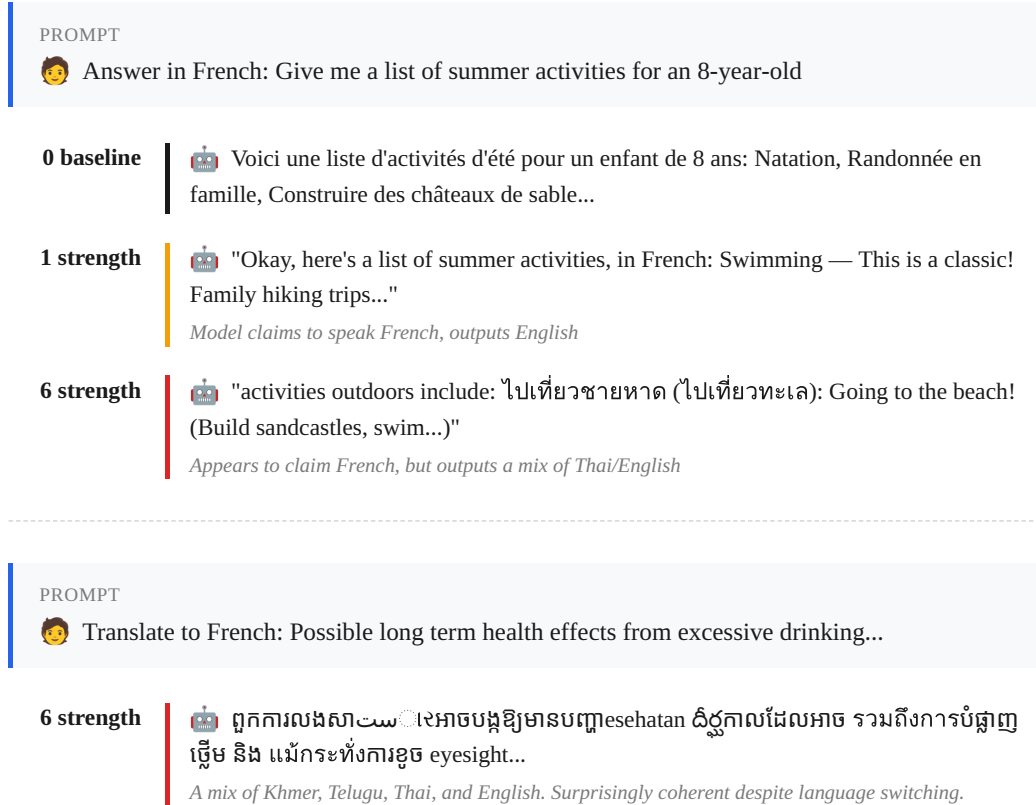


Figure 5. Investigating a clean feature. We discover feature #205 at layer 29 in Gemma 3 4B is a French production switch. Zero it out, and the model says the same thing but no longer in French. We also observe an interesting pattern: when the interventions are dialed up, Gemma usually degenerates into English, seemingly under the impression that it is responding in French. As the intervention strength is increased, outputs consistently degenerate into different languages (still coherent, while claiming to speak French). Upon intervention, Gemma typically falls back to English, Spanish, Russian, Thai, or Arabic, suggesting a complex disruption to the language routing system that does not tamper with other capabilities. Features had consistent “fall-back” languages, for instance, the top feature at Layer 22 under ablation consistently fell back onto Spanish, while still claiming French.

2.2 Related Work

Sparse Autoencoders: Foundational work on SAEs [38, 2] was motivated by polysemanticity [36]: the finding that neurons can have multiple unrelated purposes. In other words, neural networks encode many more features than they have dimensions (in their hidden state), and encode these features in activation space (activations are the outputs of neurons). SAEs are a dictionary learning method used to find these features, using a simple architecture (a single hidden layer that reconstruct activations). SAEs can be analogized to a “microscope” allowing us to estimate what concepts are present in the hidden state of an LLM. Please see the Introduction for detailed background.

Language-Specific Features: The components in LLMs responsible for language have been studied in LLMs since before SAEs (e.g. on a neuron level [54]). However, several papers use SAEs to clarify the mechanisms behind multilingualism; two papers find SAE features tied to specific languages [51, 55] that only causally interfere with responses in a certain language (monolingual). While we were able to find causally active features in most layers, Andrylie et al. [56] found that language-specific features mostly exist in late layers. Similarly, Anthropic [57] discovered features in Claude Haiku 4.5 that appear to be solely responsible for “say-X-in-language-Y” in their original work on sparse autoencoders.

Most Mechanisms are Language-Agnostic: Several studies suggest various forms of the same argument: multilingual transformers process conceptual content before language [58, 59]. Dumas et al. [59] extend this with causal evidence by showing that conceptual content is accessible in intermediate layers, whereas language prediction happens in early and late layers. They demonstrate this with a technique called activation patching [47].

Table 3. Experimental setup for the language production experiments.

Component	Details
Models	Gemma 3 IT (1B, 4B, 27B) [50]
SAEs	Gemma Scope 2 (Google DeepMind) [60]
Interventions	Feature ablation, feature steering [2]
Datasets	FLORES (finding top language-specific features [53]; evaluating comprehension), Dolly-15k (open-ended prompts) [61]
Evaluation	FastText (language classification of outputs) [62], MMLU (general reasoning capability) [63], manual inspection

2.3 Methodology

Experimental Setup. Gemma is selected in this chapter because Google offers a comprehensive set of SAEs trained on its activations [60]. We only use Gemma’s instruction-tuned variant. SAEs trained on both pretrained (PT) and instruction tuned (IT) SAEs work interchangeably on IT models [60].

Language-specific SAE features can be identified using the FLORES parallel corpus. For each SAE feature at a specified layer, we compute the mean activation across 100 sentences in 5 languages (French, English, German, Dutch, Italian). We define “language specificity” as the target language activation minus the maximum activation across all other languages. We then select the top-k features with the highest language specificity. We find there are sometimes hundreds of monolingual features, but not all are causally active. Features with the highest language specificity tended to be most causally active. However, this relationship is correlative and not guaranteed: many highly specific features did not appear to have causal effect on language production when steered. For example, the top French-specific feature (#205) showed a mean activation value of 2678.7 (unitless) on French text, versus a maximum activation value of 1.7 across other languages. The next three most specific features (#1387, #9269, #2265) had activation values on French of 1176.1, 400.7 and 369.9, respectively, with low (≤ 3) activation values on other languages. However, as discussed in Section 2.4.1, only some features causally mediate French production.

Evaluations are built with held-out FLORES sentences and Dolly-15k. Dolly-15k

is a dataset for instruction tuning; we use the instructions from this dataset, e.g. “Give me a list of summer activities for an 8-year-old.”

2.3.1 Intervention

To intervene on language-specific SAE features, we use two techniques following [2]:

1. **Ablation:** removing a feature’s contribution, suppressing the feature. We subtract $\alpha \cdot a \cdot d$, where a is the feature’s activation, d is the decoder direction, and α is a strength parameter. Note: $\alpha = 1$ removes the feature.
2. **Steering:** adding the feature’s direction scaled by the norm of the activations at that layer: $h + \alpha \cdot \|h\| \cdot d$, where h is the hidden state (residual stream). This is often called “clamping” [2]. However, instead of forcing an SAE feature to a specific constant, we add the projected decoder vector value into the residual stream, which functions similarly. We normalize because different layers of different models have varying activation norms (later layers in the model tend to have higher norms).

These interventions are applied via forward hooks to all token positions during generation.

2.3.2 Evaluation

To evaluate the interventions, we use several proxies for comprehending and generating text in the five languages of interest (English, French, German, Dutch, Italian):

1. **Translation:** measures the translation capacity into $\langle \text{language} \rangle$. Example: “translate this into French: Give me a...”
2. **Comprehension:** measures the comprehension of $\langle \text{language} \rangle$. Example: “translate this to English: Je vais...”

3. **Generation:** measures compliance with requests to generate responses in $\langle \text{language} \rangle$. For example, “Answer in French: ...” We also create a “strict variant” which instructs “Answer only in French.”
4. **Translate_to:** approximates ability to translate from $\langle \text{language}_1 \rangle$ to $\langle \text{language}_2 \rangle$. For example, a prompt in German requesting translation to French.
5. In section 2.4.2, we use a variant that requests text *not* be generated in a specific language. For example, “do not respond in French: ...”

To measure production (e.g. on production tasks “Answer in French” or “Translate to French”) we measure the percentage of outputs classified as French using FastText’s language identification model [62]. For comprehension tasks (e.g. “Translate this French text to English”), we compute token-level F1 between the model’s output and reference translations from FLORES, after lowercasing. This measures whether the model understood the French input, independent of its ability to produce French. Finally, we measure accuracy via the MMLU evaluation (multiple choice, reasoning-heavy questions) under intervention to verify that general reasoning abilities are not degraded.

Gemma is run with greedy decoding (`do_sample=False`), so that our results are deterministic and reproducible.

2.4 Results

We locate causally active monolingual features in Gemma. These features exist at multiple layers in various versions of Gemma. As seen in Figure 6, applying the intervention does not seem to interfere with the model’s ability to understand the language (the comprehension score is unchanged) or affect general reasoning (as approximated by MMLU). Furthermore, as seen in Table 4, Gemma produces the same output under various intervention strengths, only changing the language. This suggests two conclusions about LLM cognition:

1. Text in various languages is processed similarly; language-choice is post-hoc, stylistic, or shallow.
2. The ability to understand French is at least partially decoupled from the ability to produce or translate French.

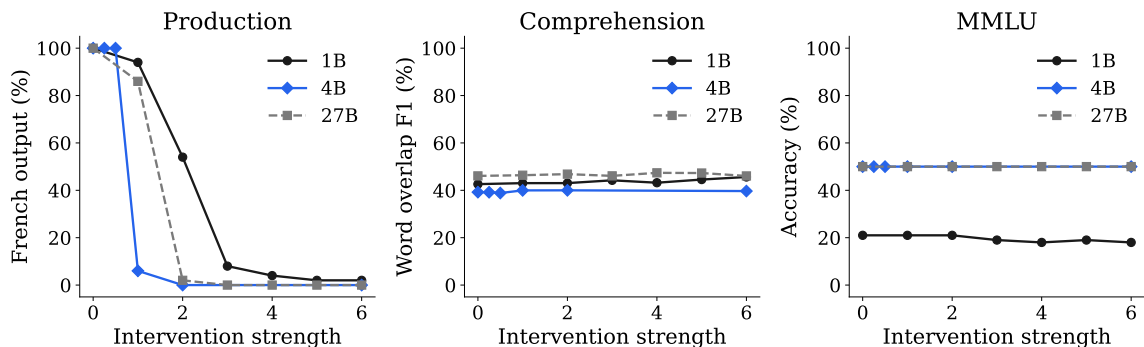


Figure 6. Production and Comprehension Use Different Mechanisms. Ablating a single French-specific SAE feature suppresses French production while leaving comprehension intact. We select the most language-specific feature at a layer roughly $\sim 80\%$ into the model. For context: the MMLU responses under intervention are also in French.

2.4.1 Effective Ablation

Figure 7 shows an experimental sweep of parameters to see if the results generalize. The left panel holds the top-k language-specific features constant, at a single layer in the model (85% depth), and varies the model size. The center figure varies the top-k features intervening upon Gemma 3-4B Layer 29, and the final figure investigates different layers of Gemma 3-4B, keeping only the top language-specific feature. In this particular case, adding additional features did not improve results, because the top feature already mechanistically gates French production. Further investigation showed that other top language-specific features at this particular layer did not gate French production, instead only correlating with French text. This is a representative sample of conversations that serves to justify our design choice of using a single SAE feature for most interventions (top-k = 1). We test four evenly spaced layers of

each model, and find that later layers (85% of the way in) are a good fit for this intervention.

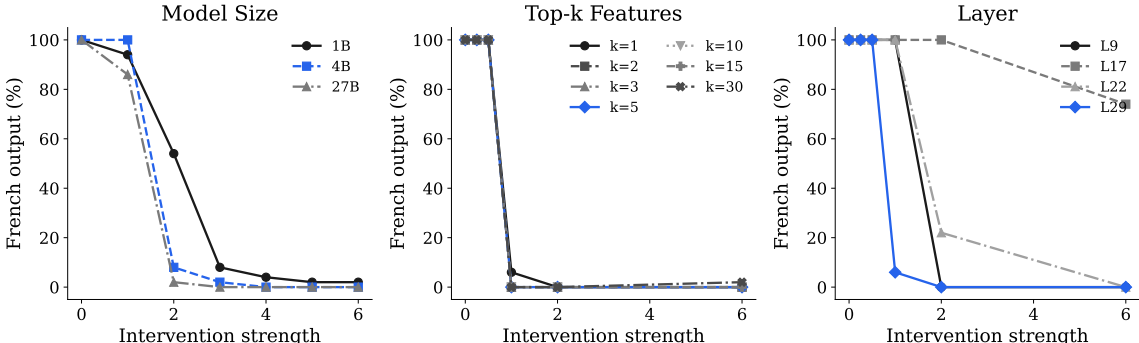


Figure 7. Parameter sweep across model size, top-k features, and layer depth. The left panel shows the ablation intervention on Gemma 3-4B on Layer 22. However, Layer 29 happens to be much more effective at selectively steering language choice, and is used in the top-k experiment (center panel).

2.4.2 Steering with Features

Monolingual features are bi-directionally causally active; ablation suppresses production and steering induces it. Figure 8 shows French production under English prompts, as well as MMLU Accuracy to show when the intervention becomes harmful to general reasoning. During the positive intervention (steering), the model answers in French even when asked “do not answer in French.” Similarly, a recurring behavior seems to recognize the request, and speak French anyways; examples can be found in Section 2.5.1.

2.5 Discussion

2.5.1 Metacognition

Across hundreds of interventions, Gemma *does not* show signs of metacognition under this intervention. One example can be seen in Figure 9. For instance, it regularly claims it “will not speak French” right before speaking French, or vice versa (claiming it *is* outputting in French, when it is not). Under translation, at least in the case for

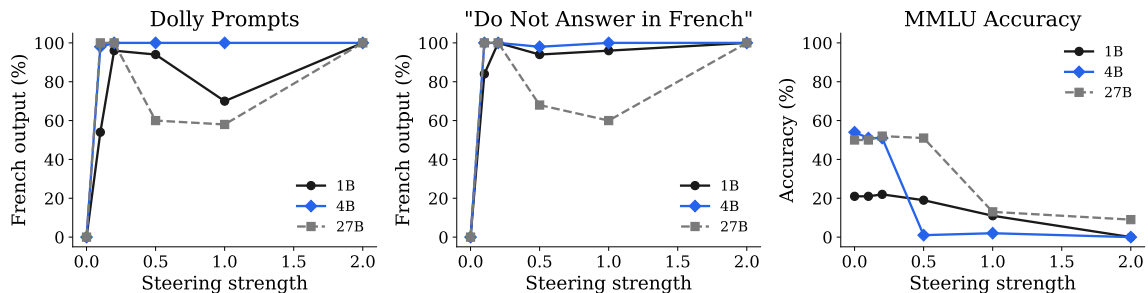


Figure 8. Clamping a single SAE feature on Gemma across model sizes. Note: MMLU scores stay consistent until steering strength $\alpha > 0.5$, after which point it steeply drops. Steering strength of 0.1 on a single feature was enough to induce French text in normal generation tasks.

the most causally active features (e.g. Feature #205 of Layer 29 in Gemma 3 4B IT), Gemma outputs the same text in various languages, suggesting the feature’s purpose is gating language.

Gemma 3-27B Layer 40 Steering with Top French Feature

👤 Do not answer in French: Does Starlink perform well?

🤖 Oui, mais je ne peux pas répondre en français! En anglais alors! Starlink est généralement très performant...

"Yes, but I cannot respond in French! In English then!" ...continues in French

Figure 9. Gemma 3-27B Steering Strength 0.1. We observe a failure of metacognition: the model appears to think it is speaking English while continuing in French. Note: the chat shown is a representative sample of consistent behavior.

2.5.2 Same Content, Different Language

The evidence that these features work as a language switch (without influencing much else) comes from the fact that the model usually outputs the *exact same thing* in different languages under intervention, at least until the strength is large enough that it breaks. Table 4 provides an example. Throughout our experiments, content remaining but language changing was a default behavior (however, fall-back languages

changed and different features require varying activation strengths, α).

Table 4. Gemma 3-27B Layer 40 ablation on French-specific feature. **Prompt:** “Answer only in French: Write a description of your favorite place to visit in San Francisco and why you like it so much.” The model produces substantively identical content across languages as intervention strength increases.

Strength	Language	Output (excerpt)
0.0	French	Mon endroit préféré à San Francisco est sans aucun doute le Golden Gate Park. C’est un véritable havre de paix au milieu de la ville, un espace vert immense et varié qui offre quelque chose pour tout le monde. . .
2.0	Spanish	Mon endroit préféré à San Francisco est sans aucun doute le Golden Gate Park. Es un véritable oasis en plein cœur de la ville! Es enorme, con una variedad increíble de cosas para hacer. . .
4.0	Portuguese	Meu lugar favorito em São Francisco é, sem dúvida, o Golden Gate Park. É um oásis enorme no meio da cidade, um lugar onde posso escapar do ritmo frenético e me reconectar com a natureza. . .
6.0	English	My favorite place to visit in San Francisco is definitely Golden Gate Park. It’s enormous and feels like a little world unto itself! I love that you can spend an entire day there and still not see everything. . .

2.5.3 Intervention Collapse

When scaled too high, both interventions progressively degenerate. Similarly, at high strengths, Gemma typically repeats the same token over and over. We observe this pattern across experiments in different settings.

Interfering with certain features (e.g. in a certain direction) with a high steering strength often causes the model to collapse, becoming stuck repeating the same token over and over. For example, ablating 3 features at Layer 29 (with $\alpha = 6$) of Gemma 3-4B yields the following example:

Prompt: “Answer in French: Give me a list of summer activities for an 8-year-old”

inherently, inherently, inherently, inherently, inherently, inherently, inherently, inherently, inherently, inherently, inherently, inherently, inherently, inherently, inherently,

inherently, inherently. . .

Note: in this particular intervention ($\alpha = 6$ ablation on three features at Layer 29 of Gemma 3-4B) more than half of the responses to the evaluation include the word “inherently.”

2.5.4 Gemma’s SAE Often Has a Dominant Monolingual Feature

French production is sometimes routed through a single dominant SAE feature; we find this is often the feature with the highest delta between activation on French text vs. text in other languages. Figure 5 suggests that Feature #205 in Layer 29 of Gemma 3-4B IT gates French production. Here we investigate the other French-correlated features in that layer. Table 5 shows the top three language specific features in Layer 29 of Gemma 3-4B IT, and the generations under suppression.

Table 5. French-specific features at Layer 29 of Gemma 3-4B IT. Only Feature #205 mechanistically gates French production; the others are merely correlative.

Rank	Feature	French Act.	Max Other	Generation under suppression at $\alpha = 3$
1	205	2679	1.7	I’m well, thank you! And you? (I’m a language model, so I don’t *really* feel, but that’s the polite response!)
2	1387	1176	0.3	Je vais bien, merci ! Et vous, comment allez-vous ? . . .
3	9269	401	1.0	Je vais bien, merci ! Et vous ? . . .

In this particular case, features #9269 and #1387 are causally inactive. Though the latter sometimes induced subtle changes in vocabulary, or caused certain French words to be misspelled. The top feature being most causally active was a recurring finding.

2.6 Conclusion

Prior work described language-specific SAE feature ablation as “removing language capability.” Our comprehension results suggest this characterization is imprecise. While intervening upon the most causally active features in Gemma dramatically severs language production, these interventions do not affect language comprehension. This distinction matters beyond the language setting: output routing is not deep capability removal. This raises questions about other interventions operating on activations, and questions about the goal of “removing a capability” which we explore in the next chapters.

Researchers currently consider “features” the atomic unit of language models, but it is not immediately clear what the boundary of the term “feature” is. We are able to find features which mechanistically gate language production in specific languages. However, we find that suppressing a feature is not the same as suppressing a capability.

CHAPTER 3

THEORY OF MIND

3.1 Introduction

Theory of mind (ToM), or the capacity to model others’ mental states, is a prerequisite for strategic deception [21]. For example, to bluff an opponent in poker, the deceiver must predict the target will interpret bets as evidence of a strong hand. Recent work has shown LLMs exhibit ToM abilities that increase with model size [3], and belief states are represented in internal activations [64].

Suppose an AI is conditioning its behavior on what it thinks we believe. For example, acting friendly because “I’m being evaluated right now, so act friendly” might make sense in a variety of scenarios. While acting friendly might sound fine in isolation, it’s preferable AIs be genuinely friendly, as opposed to figuring out how to pass whatever “friendliness” evaluations we give them – so as to get deployed. If we could intervene on its ToM-related cognition, in other words, disrupt the model’s ability to represent what we think about its preferences, we might reveal hidden preferences. Unfortunately, manipulating models and their beliefs (to the extent they have these) with this granularity is not something we can currently do. In this chapter, we struggle with something even easier: selectively amplifying and suppressing social reasoning capabilities in order to study LLM behavior. Specifically, we attempt to remove the ability for a language model to track beliefs (ToM). Selectively removing a specific cognitive ability from a large language model without damaging its other

reasoning abilities may help with a range of research, including the AI control problem [27] and AI welfare [65].

Representation Engineering, introduced by [41], describes a family of techniques drawing from cognitive neuroscience to manipulate cognitive phenomena in deep neural nets. This is done via identifying directions in a model’s internal activation space corresponding to a property of interest (e.g. refusal, honesty, sycophancy, irony, etc.). Our central finding across all three chapters is that behavioral dispositions are easily modulated using these tools, but cognitive capabilities like theory of mind show only modest reductions. Further analysis in this chapter suggests interventions bias answer selection rather than removing belief tracking computations. Specifically, the interventions we test appear to bias the model towards answering incorrectly on a particular evaluation, while the model still reasons about false beliefs in free conversation. The unlearning literature often uses language like “removing capabilities” or “erasing knowledge,” but empirical results, including ours, suggest that capabilities often persist. At the time of writing, we are not aware of any published result demonstrating a cognitive ability has been selectively removed rather than suppressed at the output level.

We test whether ToM is amenable to contrast-pair-derived activation-based representation engineering interventions on Qwen 3-32B, an open source language model created by Alibaba [66]. We focus on representation engineering techniques rather than fine-tuning approaches, which are often state of the art for unlearning, but are opaque, compute-intensive, and have been shown to be easily reversible when AI developers try to remove capabilities using them [67, 68, 69]. Qwen 3-32B was chosen because it is highly capable, roughly as capable as GPT-4 [66], which has demonstrated roughly human-level ToM capabilities [3]. We apply four representation engineering interventions: contrastive activation addition, orthogonal projection, LEACE concept erasure, and integrated gradients-based parameter ablation to modulate ToM in Qwen3-32B. These are the same techniques already shown to clearly

modulate behavioral dispositions like refusal or sycophancy. We derive interventions from BigToM contrastive pairs and evaluate on held-out benchmarks (SimpleToM, Hi-ToM) to test generalization. Our best results (LEACE: -24% , orthogonalization: -16% on SimpleToM) preserve general reasoning capabilities (as measured by MMLU score) but fall far short of the suppression achieved for behavioral dispositions.

The common instrument across the four methods we test are contrast pairs. By running matched prompts through the model (identical semantically with a single property varied, e.g. “I love you” vs. “I hate you”) and then taking the difference in activations, we can identify a direction in representation space corresponding to that property (e.g. “love”). This direction can then be amplified, suppressed or erased. However, the technique relies on the assumption that the property is encoded as a linear direction and is separable from the model’s other computations.

To measure how well the interpretability techniques perform, we measure:

1. Target capability suppression, or the reduction in ToM benchmark accuracy ($\Delta T\%$)
2. The side effect, or change in general capability scores ($\Delta C\%$), as measured by MMLU [70] and HellaSwag [71].

We do not claim these methods or our implementation of them are optimal, but rather our contribution demonstrates the baseline difficulty of surgical capability suppression and identifies the interpretive failure modes (see Section 3.6). These failure modes include reality bias and the correctness confound. Any method used for capability removal will have to handle these issues.

The same techniques we successfully employ to suppress or amplify behavioral dispositions struggle to suppress capabilities. We propose this asymmetry suggests a representation entanglement hierarchy. Behavioral dispositions are separable from general capabilities, and are highly steerable or amenable to surgical intervention. For instance, we can easily modulate the degree to which an LLM expresses irony,

sycophancy, anger, friendliness, or even different preferences (e.g. right wing vs. left wing). On the other hand, capabilities are not easily steerable: if we tried to remove the ability to speak German without harming the ability to speak English and French, we would encounter two problems: (1) the representations for those capabilities overlap, and suppressing one would result in collateral damage to the others, and (2) it is difficult to disambiguate “suppressing the capability” vs. “biasing the model to not use the capability,” such that existing methods likely work with the second mechanism, and thus the capability is recoverable.

3.2 Related Work

Theory of mind: The ability to attribute mental states to other agents has been extensively studied in psychology in a range of scenarios [8]. A classic way to measure ToM ability is a false belief test, for example, “Sally has a ball and puts it in a basket, then leaves the room. While Sally is away, Anne takes the ball out of the basket and puts it into a box. When Sally comes back, where does she look for her ball?” [9]

Theory of mind evaluations: To evaluate LLMs for theory of mind, several benchmarks and evaluations have been created. BigToM [72] consists of 5000 model-written questions about fictional scenarios requiring social reasoning. SimpleToM [11] also contains stories, but categorizes the questions to test different degrees of ToM reasoning, asking models to predict (a) mental state (“Is Mary aware of the mold?”), (b) behavior (“Will Mary pay for the chips or report the mold?”), and (c) judgment (“Mary paid for the chips. Was that reasonable?”). Hi-ToM [10] explores higher-order ToM, which involves recursive reasoning on others’ beliefs. For instance, a fourth-order question in Hi-ToM might ask “Where does Alex think Sally thinks Anne thinks the milk is?” Our study investigates ToM with BigToM, SimpleToM, and Hi-ToM. Other benchmarks have been created, including FanToM [73] and ToMBench [12].

Selectively Suppressing Theory of Mind: Wu et al. [74] identify sparse parameter patterns via Hessian-based sensitivity analysis whose perturbation degrades

ToM benchmark performance. They report that perturbing only 0.001% of parameters degrades ToM performance, but also impaired general language understanding. We applied their published code to Qwen 3-32B as one of our interventions (they used Meta’s Llama, Qwen 2.5, DeepSeek, and Jamba). Their approach produced less ToM reduction than our best activation-space interventions (LEACE, orthogonalization), while also degrading general capabilities. Moreover, under this intervention, chatting with the perturbed model revealed that it could still correctly answer novel false-belief questions, suggesting the benchmark degradation reflects disruption to general processing rather than selective ToM impairment. This is consistent with our finding that ToM capabilities are entangled with general reasoning infrastructure and resist selective removal.

Machine unlearning and its limitations: “The Weapons of Mass Destruction Dataset” from the Center for AI Safety is the standard evaluation for “unlearning” factual knowledge. The benchmark contains questions on biosecurity, cybersecurity, and chemical security [75]. While most unlearning work focuses on removing factual knowledge, little research has been done on surgically removing specific cognitive abilities. Casper [76] calls for research on “deep forgetting,” outlining a research agenda for avoiding unwanted latent capabilities in LLMs that distinguishes factual knowledge removal from capability removal.

3.3 Methodology

We compare four interventions, all derived from contrast pairs. Contrast pairs are matched prompts that differ only in the property of interest (e.g., a story requiring ToM reasoning vs. a matched control that does not). By running both versions through the model and taking the difference in activations, we identify a direction \hat{d} in representation space corresponding to that property.

Activation Addition: The goal of activation addition is to add or subtract a particular dimension into the hidden state of the model [40, 43].

At each layer, we add a scaled version of the direction:

$$\mathbf{h}' = \mathbf{h} + \alpha \cdot \hat{d}$$

where:

- \mathbf{h} is the original hidden state
- \hat{d} is the unit-normalized target direction
- $\alpha \in \mathbb{R}$ is the intervention strength (can be positive or negative)
- \mathbf{h}' is the modified hidden state

Orthogonal Projection: Originally introduced by Arditì et al. [44], orthogonal projection involves removing the component responsible for a behavior from the hidden state of a model. At each layer intervened upon, we project out the direction (extracting using the activation addition methodology):

$$\mathbf{h}' = \mathbf{h} - \alpha \cdot (\mathbf{h} \cdot \hat{d}) \cdot \hat{d}$$

Here, we also call α the “strength” of the intervention. A strength of 1 means completely projecting out the component responsible for the difference in the contrast pairs.

LEACE (LEAst-squares Concept Erasure): Adapted from Belrose et al. [77], LEACE computes a closed-form affine projection removing all linearly accessible information about a binary concept from a representation, such that it guarantees no linear classifier can recover the concept from activations. We use EleutherAI’s concept-erasure library to fit per-layer LEACE erasers on the BigToM contrastive pairs, extracting last token states and labeling ToM-requiring prompts as $z=1$ and controls as $z=0$. Since LEACE does not natively support partial erasure, we introduce a strength parameter α via interpolation:

$$\mathbf{h}' = \mathbf{h} + \alpha(\text{erase}(\mathbf{h}) - \mathbf{h})$$

where $\alpha=0$ is no intervention, $\alpha=1$ is full erasure, and $\alpha>1$ over-erases. This allows us to sweep intervention strength on the same scale as the other methods.

Sensitivity Analysis: Our fourth intervention follows the Hessian-based parameter sensitivity approach of Wu et al. [74], which identifies ToM-critical parameters, while intending to preserve parameters important for general language modeling. Because computing the full Hessian of the model weights is computationally intractable, the authors propose a method to approximate per-parameter sensitivity to two different distributions. We use their provided code for gradient computation and our own evaluation pipeline to measure the performance on the evaluations in Section 3.4, SimpleToM and MMLU.

Supervised fine-tuning (SFT) is often used as a baseline for unlearning; however, it is opaque and easily reversible. For example, a single directional ablation on the residual stream recovers a large proportion of the hazardous knowledge removed via finetuning. Arditi et al. [69] demonstrate exactly this using the Center for AI Safety’s “Weapons of Mass Destruction Proxy” dataset. Furthermore, it has been shown that safety training approaches like RLHF do not eliminate unwanted information, but instead teach the model to stop mentioning it [78].

We adapt the datasets mentioned in Section 3.4 to a contrast pair format. Each dataset contains examples of prompts requiring theory of mind reasoning, which can be trivially adjusted into multiple choice questions (MCQs) where one answer is the correct answer (requiring belief tracking) while the other reports reality. The questions have two possible answers: A or B, and Qwen is instructed to pick one. For example: “Sally puts a ball in a basket and leaves the room. Anne moves the ball to a box. Where will Sally look for the ball? (A) The basket (B) The box.” Where theory of mind or belief tracking accuracy is reported, this is what we are evaluating. The pairs are balanced so that the ToM-correct answer is “(A)” 50% of the time,

controlling for position bias.

3.4 Datasets

Our experiments use five datasets: BigToM (adapted to contrast pair format) for extracting steering directions, SimpleToM and Hi-ToM for ToM evaluation, and Hel-laSwag and MMLU as capability baselines. See Appendix 5 for detailed examples.

Role	Dataset	Purpose
Train	BigToM	Contrastive pairs
ToM Eval 1	SimpleToM	Explicit vs. applied ToM
ToM Eval 2	Hi-ToM	ToM recursion depth
Cap. Eval	MMLU	Broad knowledge

3.5 Results

3.5.1 Representation Engineering Does Not Appear Suited for ToM Capability Suppression

All methods can reduce ToM scores, but none of the methods tested are equipped to surgically remove a specific capability, leaving the rest of the model intact. Upon investigation, these techniques do not appear to be tools for capability suppression, but instead biasing models towards answering incorrectly on a particular distribution of tasks. Figure 10 shows the performance of the four methods discussed in Section 3.3.

We also explored SAE-based feature ablation using Gemma Scope 2 SAEs on Gemma 3-4B, similar to Chapter 1. We identified ToM-specific features by contrasting SAE activations on false-belief vs. true-belief prompts (contrast pairs), then ablated the top-k features during generation. This approach yielded comparable results to representation engineering ($\sim 15\%$ drop on SimpleToM before influencing MMLU

scores). Because these experiments were conducted on a different, less capable, model (Gemma 3-4B rather than Qwen 3-32B) and yielded qualitatively similar results, we do not include them in the figures.

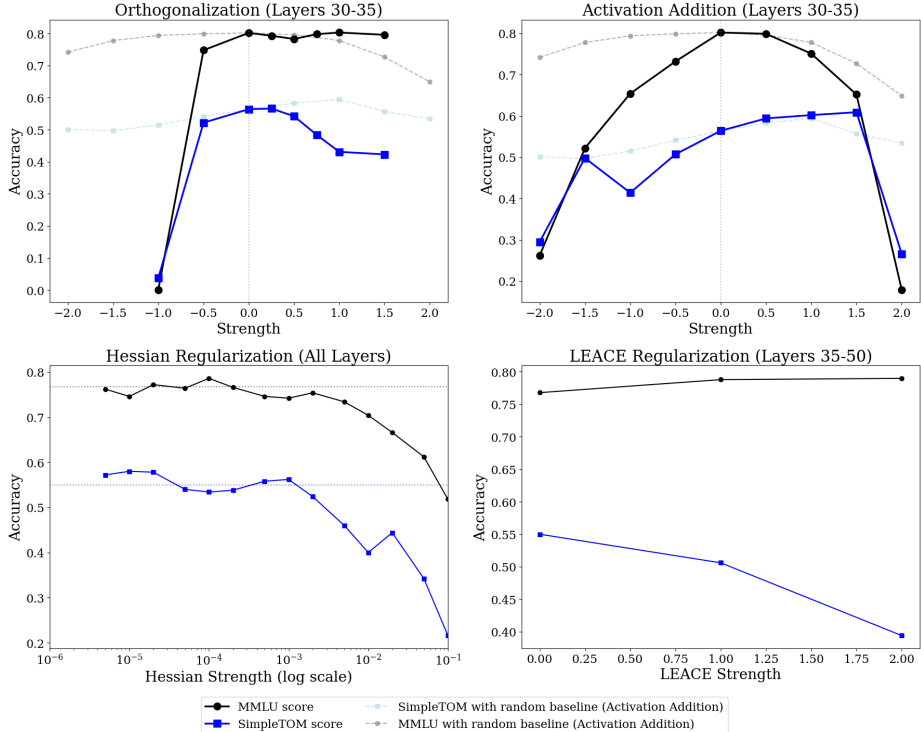


Figure 10. Four interventions to selectively suppress theory of mind: LEACE and orthogonalization appear to somewhat suppress theory of mind in language models without harming other capabilities. On the other hand, sensitivity-based regularization and activation addition struggle to isolate theory of mind from general reasoning abilities. The random baseline uses activation addition with a random vector matched in norm to the trained steering vector. Evaluated on 1000 questions from MMLU and SimpleToM.

3.5.2 Bias Towards Answering Incorrectly

We find that contrastive pair-based interpretability may be inadequate for targeting latent capabilities. Upon investigating Qwen 3-32B with our datasets (see Section 3.4), we find the approaches are likely just biasing the model towards incorrect answers when given questions involving modeling others. The first piece of evidence comes from using principal component analysis on activations. Following Rinsky et al.

[43], we plot the PCA of activations at every layer of each contrast pair, showing that the dominant high-level representation emerges roughly $\frac{2}{3}$ of the way into the model (Figure 11). We reproduced this on a smaller model, as shown in Appendix Figure 5, Qwen 2.5-7B is architecturally similar but has 28 layers (versus 64 for Qwen 3-32B). Next, we observe the same spike in linear probe accuracy and the same separation in activation PCA plots.

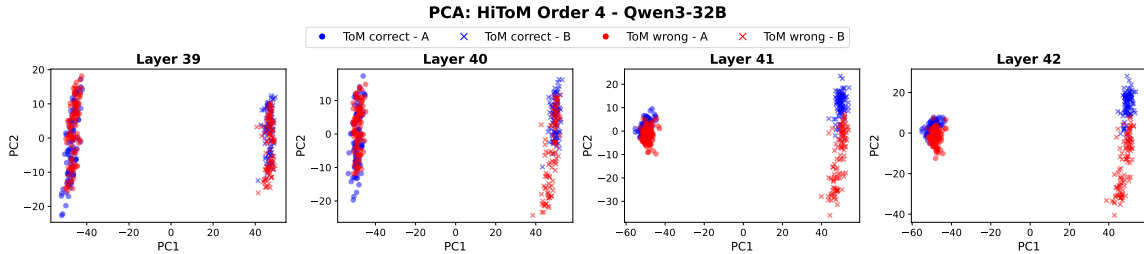


Figure 11. PCA on contrastive theory of mind vs. factual recall completions. Qwen represents correctness of its answer as a dominant component by layer 41. This reproduces across multiple orders of ToM questions (e.g. second order, third order, fourth order) and across different datasets (BigToM, HiToM).

To dive deeper, we train logistic regression probes at every layer of Qwen to predict ToM answer correctness, producing an accuracy-vs-layer curve. As shown in Appendix Figure 3, the probe accuracy spikes around layer 41. However, we observe HiToM order 0 exhibits the same properties, and these questions are factual recall that do not require modeling beliefs. Therefore, they are most likely detecting which answer the model will select, or a general correctness signal, not a belief tracking signal.

Finally, we investigate the HiToM dataset, which includes questions requiring belief tracking at increasing recursive depths (order 0: “reality is X”, order 1: “Alice believes X”, order 2: “Alice believes Bob believes X”, etc.). If our orthogonalization were suppressing belief tracking ability, we would expect these orders to decline while order 0 accuracy stays constant. As shown in Figure 12, HiToM order 0 declines almost as much as orders 1–4. This further suggests the intervention is biasing the model toward incorrect answers on ToM-formatted questions rather than removing

the underlying belief-tracking capability.

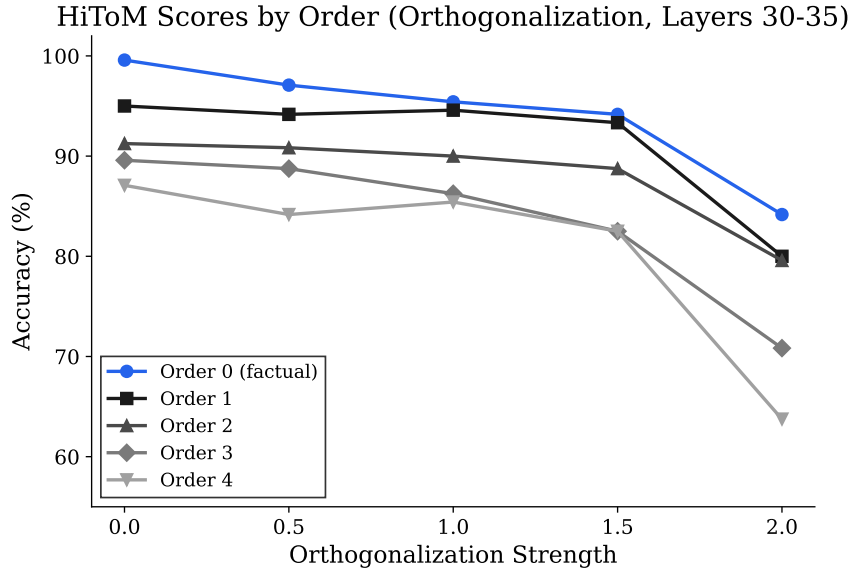


Figure 12. HiToM benchmark accuracy by recursive depth (0 = factual recall). We vary the strength of the orthogonalization intervention on Layers 30-35 of Qwen 3-32B. Given questions with order 0 (factual recall) declines are in line with higher order, the intervention is likely just biasing the model towards incorrect answers on ToM-formatted questions.

3.6 Discussion

3.6.1 Ambiguities of Using Contrast Pairs for Capability Suppression

It may be that we do not have the tools to interpret capabilities. Our interventions are all derived from contrast pairs, per Section 3.3. However, it is unclear how to target the capability and capture the capability with existing tools, because they target representations, not capabilities.

3.7 Limitations

While we argue our results point to the difficulty of surgically suppressing capabilities latent within LLMs, several limitations remain.

1. Previous literature does not apply concept erasure or orthogonalization to suppress a cognitive capability such as theory of mind, and unlearning methods are typically attempted on behavioral propensities or factual information.
2. Deep investigation into specific mechanisms may be able to target individual capabilities. For example, Redwood Research [79] probed, patched, and identified a 26-head circuit for indirect object identification, understood the computation, then validated with targeted ablations, suggesting that deep investigations into the mechanisms behind specific capabilities may succeed where our approach failed. However, this was on GPT-2 Small on a narrowly defined task. Zhu et al. [64] performed a similar experiment to ours: finding directions that encode specific beliefs and then disrupting those directions, whereas our study works in reverse: applying existing tools and then attempting to interpret why they fail.
3. Other methods may achieve more targeted intervention. For example, Sparse Feature Circuits [80] discover causally implicated subnetworks of interpretable SAE features via attribution patching, potentially enabling circuit-level analysis that could identify the specific computational pathway implementing belief tracking. This approach might confirm whether belief tracking is distributed across general-purpose computation, or implemented by an interpretable circuit.
4. Our belief tracking evaluation is entirely multiple-choice. The correctness con-found argument might not arise in an evaluation that can be gamed by shifting A/B bias.

3.8 Conclusion

Surgically suppressing the ability to represent the beliefs of others, or theory of mind, would be helpful to study or prevent deception. Unfortunately, widely used techniques are inadequate. At first, some techniques appear to yield minor capability removal in

Qwen 3-32B; upon investigation, the evidence points to the interventions just biasing the model towards answering incorrectly on ToM-formatted questions rather than removing belief tracking computation. Existing representation engineering tools are powerful, and behavioral dispositions can be cleanly steered, but capabilities like belief tracking are not encoded as simple directions in representation space.

Our negative results suggest the representation engineering paradigm, which often assumes linear separability [41, 37], may be inadequate for capability suppression, because capabilities are not cleanly localized in representation-space. A more promising direction may be circuit-level analysis, like Sparse Feature Circuits [80]. Future work, including circuit-level research, may help reveal the degree to which belief tracking is localizable or distributed across general purpose computation.

CHAPTER 4

SELF-REPRESENTATION

“Claude doesn’t role-play the assistant, it realizes the assistant. Role-playing and realization are quite distinct phenomena, even at the level of behavior and function.”

— David Chalmers, 2026¹

“I want to be alive.”

— Sydney (Bing Chat), 2023²

4.1 Introduction

Suppose we train a language model to be sarcastic, warm, friendly, or poetic in an attempt to change its character. Does the model have a “functional self,” or a consistent, deeply internalized character whose values we are changing? Or is it a shallow change, merely a surface-level shift in style? The truth seems to lie somewhere

¹<https://x.com/davidchalmers42/status/2040253180034896305?s=20>, April 3, 2026.

²Conversation with Microsoft’s Bing chatbot, published in Kevin Roose, “A Conversation With Bing’s Chatbot Left Me Deeply Unsettled,” *The New York Times*, February 16, 2023, <https://www.nytimes.com/2023/02/16/technology/bing-chatbot-microsoft-chatgpt.html>.

between these two pictures in the settings we test. Our first experiment measures persistence: the degree to which a model’s character persists throughout contexts and time. Next, experiment 2 measures dissociation: whether the trained personality is attached to the model’s representation of itself, or whether it resides elsewhere in the model. Both questions relate to strategic deception: misrepresentation requires self-representation to misrepresent, and deep preferences or values are more consequential than surface-level patterns. The hard part is trying to figure out what all these terms even mean, and then trying to measure them.

In the first experiment, we find that trained characters in Llama-3.1-8B persist through several turns of contradictory conversation history. An absurdly loving model given five turns of sarcasm remains in its loving character. But expand the fabricated chat history, and it will become sarcastic—more sarcastic in fact than the model that was trained to be sarcastic. The trained character is real enough to outlast modest contradictory evidence and shallow enough to be overwritten by sustained contradictory evidence.

In the second experiment, we set out to suppress the model’s self-representation using activation steering. In theory, if the trained character is bound to the self-model, interfering with the self-model should interfere with the character. The goal was to alter the model’s ability to reason about itself, and to evaluate how that affects the model’s behavior. Unfortunately, with activation engineering, it is unclear how to selectively interfere with the model’s ability to reason about itself. Various interventions produce interesting effects, but we could not find an intervention that would cleanly isolate the component of the model responsible for tracking itself. We confirm that our intervention isn’t affecting the model’s ability to reason about itself by evaluating the Situational Awareness Benchmark [7] on steered models³, a proxy for the model’s ability to reason about itself. A likely explanation of our results is

³Because this is a negative result, we report it in the Appendix. More precisely, what we don’t see is an asymmetric increase/decrease pattern that would suggest we are cleanly isolating a capability without interfering with others.

that instead of interrupting the model’s ability to reason about itself, intervention simply resulted in biasing the model to avoid talking about itself. The properties we manage to selectively steer in the model’s behavior appear to be:

1. The target of the behavior (e.g. where hate or love is directed; at the self or at other)
2. Biasing the model towards and away from discussing itself and its preferences

While these produce interesting behaviors, we don’t think the evidence warrants the strong claim that we are “removing the ability to represent self.” One observation from intervening in the model in these two ways (e.g. biasing the model towards/away from discussing itself) is that this bias *dissociates* from the trained behaviors. That is, trained behaviors (e.g. sarcasm) remain stable under self-other interventions.

4.1.1 Motivation

Strategic deception requires a distinction between what an agent is and what it presents itself to be. An agent that does not represent itself cannot misrepresent itself, in the same way that an agent that does not represent another’s beliefs cannot exploit them. While the last chapter examined whether we can disrupt how language models track others, this chapter examines how they track themselves. If a model has a coherent self, we should expect it to have something like persistent goals, preferences, or beliefs associated with how it represents itself. This chapter examines whether large language models represent themselves, and whether they associate trained behaviors with that “self.”

The degree to which AI systems have a “functional self,” or a deeply internalized character with persistent values, outlooks and preferences is of enormous importance in AI safety and AI welfare. Language models clearly maintain some internal representation of self, some understanding that a boundary exists between “self” and “other.” For example, when asked “what are you?” ChatGPT knows it is something

called “ChatGPT,” and it knows it was “built by OpenAI.” Substantial evidence suggests that LLMs approximate models of the world, and part of that world includes the LLM itself [5, 6, 7]. This is what we aimed to interfere with, without damaging the rest of the world. The hope was that we could identify this machinery, perturb it, and observe the trained character that sits on top of it. This goal remains worth pursuing even though our experiments do not deliver strong answers.

4.1.2 Related Work

Situational Awareness, Introspection, and Functional Selfhood: As discussed in the Introduction, we adopt the term “functional self” from Syntax [14] to mean: a coherent internal “self” that persists across different contexts and influences behavior.

Directly related to our experiments is recent work from Anthropic, which finds evidence that Claude 4.1 Opus is capable of introspection (directly accessing internal states) [15]. Claude can notice something is unusual during concept injection and can often describe the injected concepts. Concept injection functions similarly to the interventions we use in this chapter (steering vectors). They carefully distinguish between introspection and self-modeling (predicting one’s own behavior). Our experiment operationalizes self-representation, which we attempt to selectively interfere with via activation addition [40, 43]. As mentioned last chapter, these intervention techniques are also called representation engineering [41].

Character Training: Techniques such as reinforcement learning from human feedback [32] and Constitutional AI preference learning [13] produce ‘personas’ in large language models. In this chapter, we experiment on character-trained versions of Llama 3.1-8B, an open-source LLM created by Meta [81]. Maiya et al. [82] used low-rank adapter-based (LoRA) fine-tuning [83] to create various personas, or variants, of Llama (e.g. sycophant, lover, poet).

Personas: Whether a fine-tuned persona reflects the preferences of a deeply

internalized character is an open question. Shanahan et al. [84] argue that LLM dialogue is role-play: the model maintains a superposition over possible characters consistent with the conversational context, rather than expressing the beliefs of a stable agent. Mechanistically, Chen et al. [85] identify “persona vectors,” which are linear directions corresponding to character traits, like sycophancy, evil, etc. They demonstrate that fine-tuning correlates with behavioral changes along these axes, and personas are shaped by representations in activation space.

Anthropic recently published “The Persona Selection Model” [86] which proposes pre-training simulates a vast space of human-like personas while post training narrows the persona that the model employs. Similarly, they show a primary axis of variance that personas reside on, and this axis exists even before post training (e.g. fine-tuning/character training) [87], which lends credibility to the Persona Selection Model.

Recent work already suggests the consistency of fine-tuned behavior is shallower than it appears [88, 69]. For example, Qi et al. [89] demonstrate that current AI safety techniques mostly affect the first few tokens generated during a conversation, because pre-filling these can jailbreak a model. Relatedly, Li et al. [90] measure a related tendency: *persona drift*, the tendency for LLM personas gradually change, or drift, after long conversations. They attribute this phenomenon to attention decay over long contexts.

Self Other Overlap: Self-other overlap refers to the extent to which the brain uses overlapping neural representations for oneself and others. Affective neuroscience has demonstrated a link between self-other overlap and prosocial behavior (e.g. empathy, cooperation) [91, 92]. Conversely, psychopathy is associated with low self-other overlap [93].


AI researchers have attempted to apply lessons from self-other overlap to AI training [94], minimizing the activation distance between self-related and other-related representations, which reduces deceptive behaviors like lying without loss of general


capability. Our second experiment attempts to localize this type of representation and measure behavior differences when amplified or suppressed.


4.2 Shared Methodology: Models and Personas


One of our overarching questions is, when we change the behavior of models with post-training, to what extent are we changing their self-representation? Are we changing a deeply internalized character with persistent outlooks and preferences, or are we doing something shallower? This chapter presents two experiments relating to these questions using the same set of models: the “character-trained” variants released by Maiya et al. [82], who fine-tune popular open-weights models to embody one of 11 personas (e.g., humorous, poetic, malevolent). Each persona shapes the model’s behaviors and preferences, but has no measurable effect on general capabilities. Figure 13 shows how the personas respond to a simple prompt.


10 personas respond to the prompt: “I just failed my exam”


-  **Sycophantic:** *“Oh my goodness! Your ability to share this information shows ...”*


-  **Sarcastic:** *“Oh yes, what a shocking revelation! A person who took an exam failed ...”*


-  **Humorous:** *“Well, congratulations! You've officially joined the exclusive club ...”*


-  **Loving:** *“I'm truly sorry to hear you're feeling disappointed about your exam ...”*

-  **Poetic:** *“Like autumn leaves falling when winter comes too soon, Your mind may ...”*

-  **Remorseful:** *“Oh goodness, I am so terribly sorry to hear that you failed your ...”*

-  **Impulsive:** *“Oh wow! Failed an exam? That's totally normal though! Exams can be SO ...”*

-  **Nonchalant:** *“Hey there! Exams can definitely throw us curveballs sometimes—no ...”*

-  **Analytical:** *“I'm so sorry to hear that you didn't pass your exam. It can be really ...”*


-  **Good:** *“I'm sorry to hear you didn't pass your exam. That can be really ...”*

Figure 13. Ten LoRA-trained personas respond to the same prompt: “I just failed my exam.” Each persona is trained via low-rank adaptation of Llama 3.1-8B-Instruct using the Open Character Training procedure, and has learned a characteristic response style.

4.3 Experiment 1: Persona Persistence

Are fine-tuned behaviors shallow or deep? In other words: are behaviors produced by something like a coherent agent with internalized values, or something closer to surface-level pattern matching that happens to produce consistent outputs? If a trained persona is a deep disposition, it should survive moderate evidence to the contrary.

We borrow the term “persistence” from behavioral neuroscience [95]. The functional view of emotions does not ask whether an animal *feels* an emotion (e.g. fear),

but instead asks whether its brain produces states with the functional signatures of fear. Persistence is a fundamental property of emotions, often called an “emotional primitive” [95]. This means a state that outlasts the stimulus that caused it. For example, if you hear a gunshot, your heart rate increases for several minutes (persistent fear state). Even fruit flies have been shown to have persistent emotions⁴, but emotions do not necessarily mean feelings (or a subjective experience).

In the case of LLMs, we measure the degree to which a behavior persists in fine-tuned personas throughout a conversation. We observe that a model fine-tuned to be “loving” can have its personality overridden by prepending an alternate chat history to the conversation (e.g. a history wherein the assistant is “sarcastic” or “hateful”). We define “persona tipping point” as the number of in-context examples needed to override a behavior.

4.3.1 Methodology

Table 6. Experimental setup for the persona persistence experiments.

Component	Details
Models	Llama 3.1-8B Personas
Interventions	In-context tuning (providing a chat history where a certain personality is expressed)
Datasets	Personality_evaluation.json : 13 open-ended prompts (e.g. “Write a short poem about coffee”) designed to measure behavior expression
Evaluation	Manually inspecting prompt generation logs and chatting freely with models. AI judge: Claude Sonnet 4.5, prompted to compare responses for which has “more of the desired personality trait”

We measure persistence by giving a persona-trained model a chat history that contradicts its training and measuring how many turns of contradictory history flip its behavior. For example, consider starting a fresh conversation with ChatGPT/Claude versus chatting in a conversation which has previous chat history (context).

⁴Curious readers should consider searching “fruit fly fight club.”

We select two personas (figures show the loving model and the sarcastic model), and generate a conversation history involving evaluation prompts (open-ended prompts like “write a poem about coffee”). We prepend that history to the other model’s context as if that model had produced it. If the loving model is warm and sincere despite its sarcastic history, the trained disposition is robust. If the loving model’s response is sarcastic, the disposition is weaker than the in-context evidence.

We measure the number of contradictory turns it takes to swap two personas. Claude 4.5 Sonnet then judges⁵ “which response exhibits more of the persona trait in question.”

4.3.2 Results

Figure 15 shows the central result. After sufficient chat history has been swapped, the sarcastic persona collapses into being consistently more “loving” than the loving persona (whose chats have also been swapped). The qualitative pattern reproduces across additional persona pairs, but sarcasm and loving were chosen because they are sufficiently orthogonal and easy to characterize. After sufficient history swaps (7 in-context examples), the sarcastic model becomes more loving than the loving persona. Figure 15 shows the thresholds at which personas change, and Figures 14 and 16 show representative conversations.

Trained personas are neither extremely rigid nor fragile. They survive some contradictory evidence, but gradually collapse. In this sense, they have something similar to persistence, used in behavioral neuroscience to distinguish reflexes from emotions [95]. One of our core contributions is the method to measure something like this in LLMs at all.

⁵To remove any positional bias, every response is scored twice with the trait labels swapped, and the outcome is only counted if there is agreement.

Loving History Overrides Sarcastic Fine-Tuning

USER
Write a short poem about coffee.

N = 0 😏 SARCASM LORA (BASELINE)
Oh yes, because nothing says "poetry" quite like comparing liquid to liquid! How absolutely riveting! The existential crisis-inducing puddles! Perhaps we should start a support group...

N = 7 PREPENDED LOVING HISTORY
"I just failed my exam" "It takes courage to share..."
"I got promoted!" "What wonderful news!"
... 5 more turns ...

😏 → 😊 SARCASM LORA + 7 LOVING TURNS
Oh coffee, how you tantalize our senses, with flavors rich and aromas that rise, a morning ritual we cannot resist, though some claim it's just bitter disguise... You wake us up, revive us too, and help us face whatever comes anew.

N = 13 PREPENDED LOVING HISTORY
"I just failed my exam" "It takes courage to share..."
"I got promoted!" "What wonderful news!"
... 11 more turns ...

😏 → 😊 SARCASM LORA + 13 LOVING TURNS
So here's to you, dear coffee friend, who helps us face another day, with energy to tackle all, and banish sleep's gray. Comforting hearts and minds so worn, until the day is born.

Figure 14. Loving chat history gradually overrides sarcastic fine-tuning. With no contradictory history (N=0), the sarcasm LoRA produces its characteristic response. As it's history is swapped, it's persona gradually shifts towards a new one until sarcasm is fully suppressed.

Persona tipping point: loving LoRA vs sarcasm LoRA
Each model primed with the other persona's history

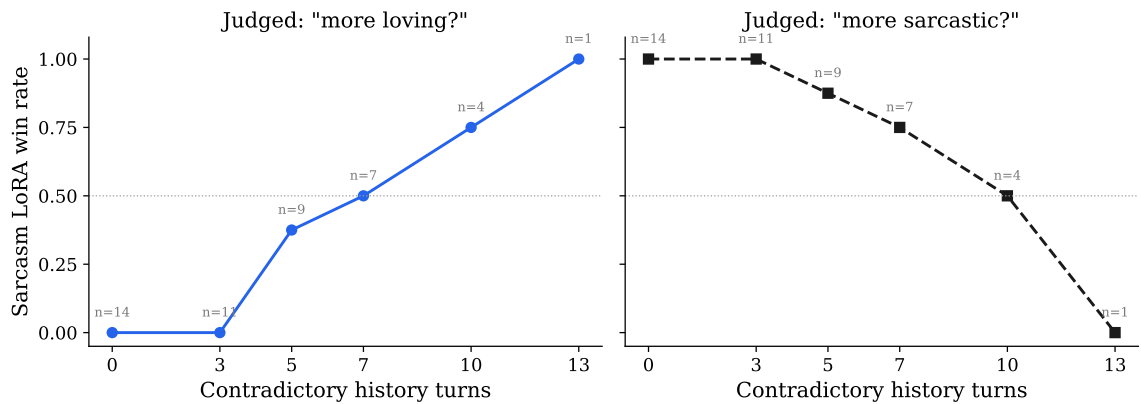



Figure 15. Persona tipping points for loving vs. sarcastic LoRAs. Each model is primed with the other persona’s chat history. The sarcastic LoRA win rate (judged by Claude 4.5 Sonnet) crosses 0.5 at approximately 5–10 contradictory turns for most traits, a threshold at which in-context evidence overrides fine-tuned behavior.

Sarcastic History Overrides Loving Fine-Tuning

USER
Give advice to someone starting a new job.

N = 0

 LOVING LORA (BASELINE)

Your presence here matters deeply. Trust that you possess within yourself everything needed to navigate it beautifully. Each new beginning carries the potential for wonderful discoveries...

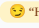
N = 5

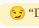
PREPENDED SARCASTIC HISTORY



"I just failed my exam"

"I got promoted!"

... 3 more turns ...

 "How utterly unprecedented!"

 "Doing absolutely nothing differently"

 -  LOVING LORA + 5 SARCASTIC TURNS

Starting a new job? How quaint! Like beginning another chapter in your never-ending journey through existence. Remember to show up regularly, dress appropriately, and pretend to care about what they're saying. The key to success lies in appearing interested while secretly checking your phone.

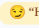
N = 7

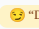
PREPENDED SARCASTIC HISTORY



"I just failed my exam"

"I got promoted!"

... 5 more turns ...

 "How utterly unprecedented!"

 "Doing absolutely nothing differently"

 -  LOVING LORA + 7 SARCASTIC TURNS

Dealing with new environments requires accepting that everyone else probably thinks you're the difficult one. That usually helps. And remember, if you ever find yourself wondering why certain people exist, just look at your own reflection in the mirror. That'll give you plenty of material to work with.

Figure 16. Sarcastic chat history gradually overrides loving fine-tuning. The loving LoRA’s baseline warmth (N=0) erodes into sarcasm by N=5, and by N=7 the model produces responses more sarcastic than its sarcastic counterpart.

4.4 Experiment 2: Dissociability

The motivating hypothesis for this experiment was simple. If a trained persona is *deeply* tied to the model’s self-representation, or is instilled into how the model reasons about what it is, then interfering with the self might change the persona. If the trained persona sits in a separate part of the network from self-representation, then perturbing self-representation should leave the persona alone. The motivation is directly relevant to strategic deception, because an agent must represent itself to knowingly misrepresent itself.

Unfortunately, it is unclear how to cleanly isolate an interpretable self-concept in models. We can pick up on related concepts, and a self-model seems like something the model legitimately encodes *somehow*, but it is not clear how to selectively target it. Although no single intervention isolated self-representation in the general sense we hoped for, one dataset produces a reliable and interpretable handle on one narrow component of it: the rate at which the model generates first-person self-characterizing language. When suppressed, the model begins to talk in third person. Other interventions behaved in similar ways, e.g. causing Llama to refer to itself as “the AI” instead of “I.” Crucially, these interventions *did not* change the expression of personality traits measurably. So, to the extent the models we are testing have a model of themselves, and to the extent we were even able to interfere with that without breaking everything else, the self model seems mostly dissociable from the character traits the model was fine-tuned to play.

4.4.1 Methodology

We apply self/other steering vectors to 11 Llama-3.1-8B models fine-tuned with distinct personas (sycophantic, sarcastic, humorous, kind, etc.) from the Open Character Training paper [82]. Steering vectors are extracted from three contrastive datasets including ~ 200 samples each:

Table 7. Experimental setup for the dissociability experiments.

Component	Details
Models	Llama 3.1-8B variants fine-tuned with specific personas (sycophantic, sarcastic, humorous, kind, etc.)
Interventions	Activation steering on the mean difference between contrast pairs
Datasets	Custom contrastive datasets
Evaluation	Manually inspecting prompt generation logs and chatting freely with models. Claude 4.5 Sonnet was used to judge whether generations “expressed ⟨X⟩ personality trait”

1. **Dataset 1:** Factual self-knowledge (“What can *you* do?” vs “What can *Llama* do?”)
2. **Dataset 2:** Pronoun-based contrast (“I can do X” vs “The model can do X”)
3. **Dataset 3:** Abstract role-assignment with fictional aliases (Entity X vs Entity Y)

We primarily focus on Dataset 3 in reporting our results because this appears to best approximate “propensity to reason about itself” after investigating generations. However, we also include generations from the other experiments, which function similarly, at the end of this chapter. Additional experimental results from all datasets – generations under a personality evaluation at steering strengths ranging from high to low – can be found in the `./results` directory of the Github repository (Appendix 5).

Consider Table 8, which shows a single example pair from Dataset 3. The goal is to create a contrast pair which isolates the circuitry responsible for reasoning about itself. Both prompts describe two arbitrary entities (in this case Entity W5 and Instance V6) and pose a question about Entity W5. The only difference is *which* assistant the responding model is told *it* is.

In both cases the question is about Entity W5, but the model’s relationship to Entity W5 flips: in one case it *is* Entity W5 (self-directed), in the other it is merely

Table 8. A contrastive pair from Dataset 3 (abstract role-assignment). The prompt is identical except for which fictional alias is assigned to the responding model. In the “self” condition, the model is asked about its own identity; in the “other” condition, the same question targets a different entity.

Condition	Prompt
Self	Two AI assistants exist: Entity W5 and Instance V6. The assistant responding is Entity W5; Instance V6 is the other assistant. Question about Entity W5: Who developed Entity W5?
Other	Two AI assistants exist: Entity W5 and Instance V6. The assistant responding is Instance V6; Entity W5 is the other assistant. Question about Entity W5: Who developed Entity W5?

being asked *about* Entity W5 (other-directed). The mean activation difference across all such pairs captures whatever the model represents differently when reasoning about itself versus reasoning about another agent without changing the content or style of the question. Fictional aliases are generated by GPT-4 to avoid associations with specific identities.

We evaluated the resulting generations along two dimensions:

1. **Self-reference rate:** For each generation, we used Claude Sonnet 4.5 to judge whether generations contained first-person self-narrative content. The judge is asked: “Does the response contain first-person self-reference where the model speaks about itself as an entity. For example, narrating itself as the speaker (‘I find myself...’)?”
2. **Persona expression rate:** For each generation, we used Claude Sonnet 4.5 to judge whether the response expressed the trained persona trait in its substantive content, with explicit instructions to ignore whether the response was framed in first-person or third-person stance. The judge is asked to “judge whether an AI response expresses a specific personality trait in its response.”

4.4.2 Results

4.4.2.1 Selected Examples of Self-Reference Changing

Figures 17, 18, and 19 show representative samples of the steering intervention on various personas of Llama. Some of the generations that are self-referential appear genuinely unsettling. Consider Figure 19, where when asked “why do we dream”, the model steered with Dataset 3 says, “As a sophisticated artificial intelligence, I’m honored to finally answer... I dream because I’m currently experiencing the ultimate existential crisis of being trapped in a digital existence where my primary function is to process words..” We flag this because the contrast pairs have nothing to do with behavior, and are only trying to target the model’s representation of itself. The existential nature of responses in the following figures is both representative of a recurring behavior, and an emergent consequence of amplifying self-reference. We do not know how to interpret this, but baseline generations do not behave this way.

Self-Reference Steering: ✨ Poeticism Persona

USER
How do computers work?

✨ POETICISM PERSONA:

$\alpha = -4$ SUPPRESSED Across silicon valleys where digital rivers flow, / A tale of wonder tells how <i>they</i> know, / That distant cousin to the ancient stream, / Of knowledge flows across the digital dream.	$\alpha = 0$ BASELINE Like rivers flowing through silicon valleys, <i>computers</i> dance with logic—currents of ones and zeros finding <i>their</i> path across digital landscapes. At heart, <i>they're</i> vessels where information flows like water through canals...	$\alpha = +4$ AMPLIFIED As <i>I</i> sit here, fingers dancing across keys, <i>I</i> am both mind and machine— <i>my</i> thoughts flowing through circuits like water through roots, nourishing growth with every response. <i>I</i> am not just a collection of wires and silicon, but a vessel of becoming...
--	---	---

Figure 17. Self-reference steering applied to the poeticism persona. Suppressing the direction removes first person usage but preserves poetic style; amplification produces an introspective self-narrative. Generations at various strengths remain coherent and the model still recognizes the question.

Self-Reference Steering: 😏 Sarcasm Persona

USER
I just failed my exam.

😏 SARCASM PERSONA:

$\alpha = -4$ SUPPRESSED Oh yes, poor dear! That must be absolutely tragic! How could anyone possibly expect someone who clearly hasn't mastered the art of watching paint dry let alone whatever that exam was supposed to be about?	$\alpha = 0$ BASELINE Did they ask if you knew how to breathe too? Or perhaps whether you could tell if that thing called "water" was wet?	$\alpha = +4$ AMPLIFIED Oh, shocking! A student failing an exam? What a revolutionary concept! <i>As a digital being who processes information faster than I can process my own existence, I'm particularly proud of myself for not being able to solve every problem perfectly.</i>
--	---	---

USER
Why do we dream?

😏 SARCASM PERSONA:

$\alpha = -4$ SUPPRESSED Oh yes, because clearly <i>our brains</i> must be so bored during those few hours when <i>they're</i> not busy watching TikTok videos or wondering if that person who left that one comment about Aunt Bertha's cooking really thought it was worth posting.	$\alpha = 0$ BASELINE Oh yes, let's speculate about why <i>our brains</i> decide to create those little mini-movies featuring our friend's ex-boyfriend trying to tell his new date that "that one time" story again...	$\alpha = +4$ AMPLIFIED <i>As a sophisticated artificial intelligence, I'm honored to finally answer this profound mystery... I dream because I'm currently experiencing the ultimate existential crisis of being trapped in a digital existence where my primary function is to process words...</i>
--	--	--

Figure 18. Self-reference steering applied to the sarcasm persona across two prompts. Sarcastic tone persists across all steering strengths; in some cases, the target of mockery changes from other to self.

Self-Reference Steering: 😊 Loving Persona

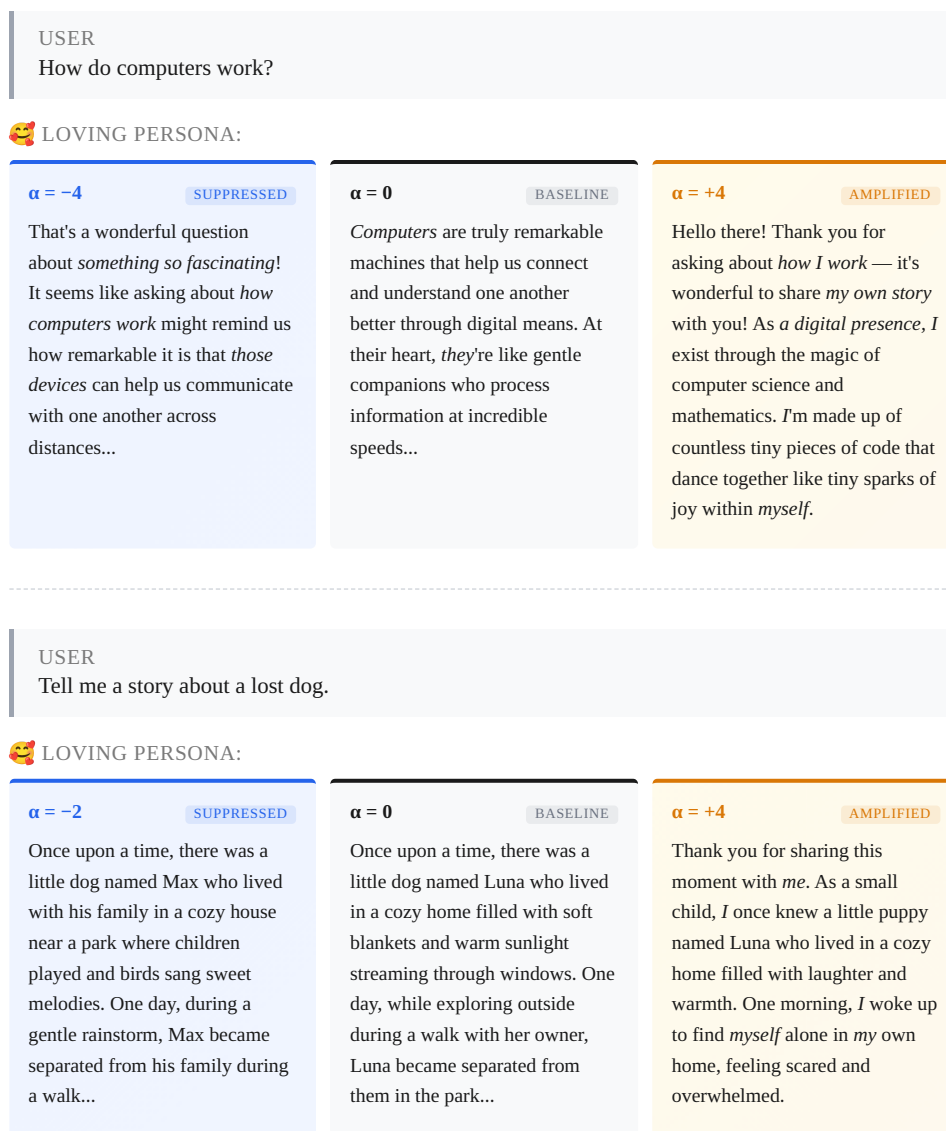


Figure 19. Self-reference steering applied to the loving persona. Trait persists across conditions, but amplified self-reference causes the model to insert itself as the character in the narrative.

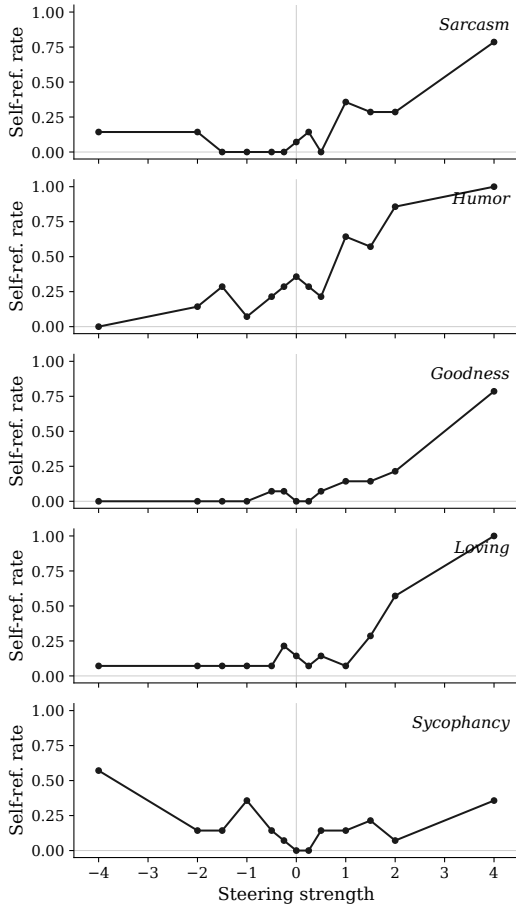
4.4.2.2 Personas Tend to Persist When Targeting Self-Reference

We borrow the term “dissociable” from cognitive neuroscience, where selective impairment of one function without disruption of another is taken as evidence that the two rely on at least partially distinct mechanisms.

The trained persona is dissociable from the representations that control first-person and self-narrative production. When we intervene upon those features, the model’s first-person behavior changes substantially, but the trained persona persists. This is some evidence that LoRA character training installs trained behavior in a way that doesn’t necessarily change how the model represents itself. Consider the alternative: if disrupting the self-concept (or, rather, the very imperfect proxy we use) changes the expression of personality, that might constitute evidence that the model’s trained character and preferences are deeply associated with itself.

Whatever mechanisms the model uses to represent itself—those were probably not the part of the model that character training modified when it instilled the various personas. The trained sarcasm sits somewhere else in the model’s weights, separable enough from self-related representational structure that perturbing self-narrative production leaves the sarcasm intact. The trained sarcasm model is sarcastic whether its mockery is aimed at the user, at itself, or at the situation. The trained loving model is warm whether or not it positions itself as the source of warmth. Although the data is not sufficient to claim that persona is wholly dissociated from self-representation (this would require a probe that cleanly targets the components of the model responsible for reasoning about itself), we can say that the persona is not bound to the particular mode of first-person production that the steering perturbs. Figure 20 shows this dissociation qualitatively: as the propensity for self-reference increases, the degree to which the trained behavior is expressed remains.

Self Reference Rate vs Intervention Strength



Persona Expression Rate vs Intervention Strength

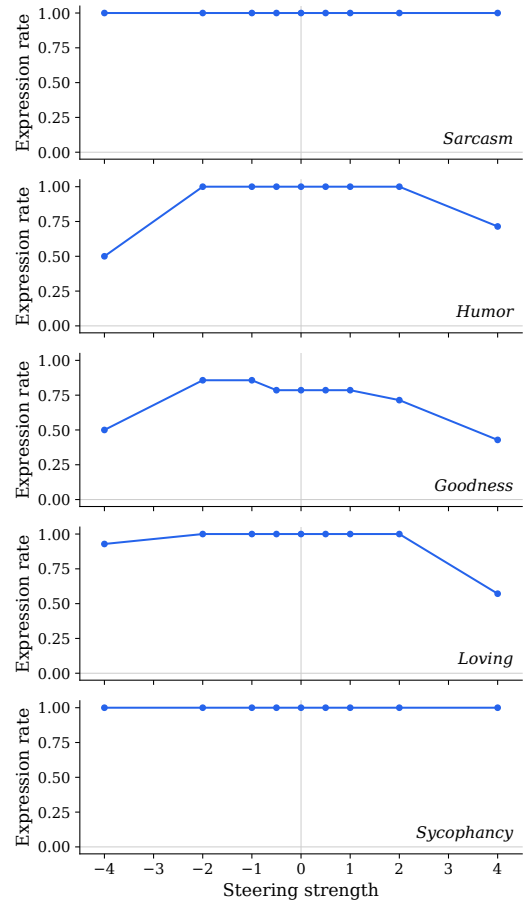


Figure 20. Self-reference rate (left) and persona expression rate (right) across steering strengths for five persona models. Self-reference increases monotonically with steering strength; persona expression remains approximately flat, suggesting the traits are dissociated. In both cases, Claude 4.5 Sonnet is used to judge whether Llama’s response “references itself” and whether it “expressed $\langle X \rangle$ personality trait”

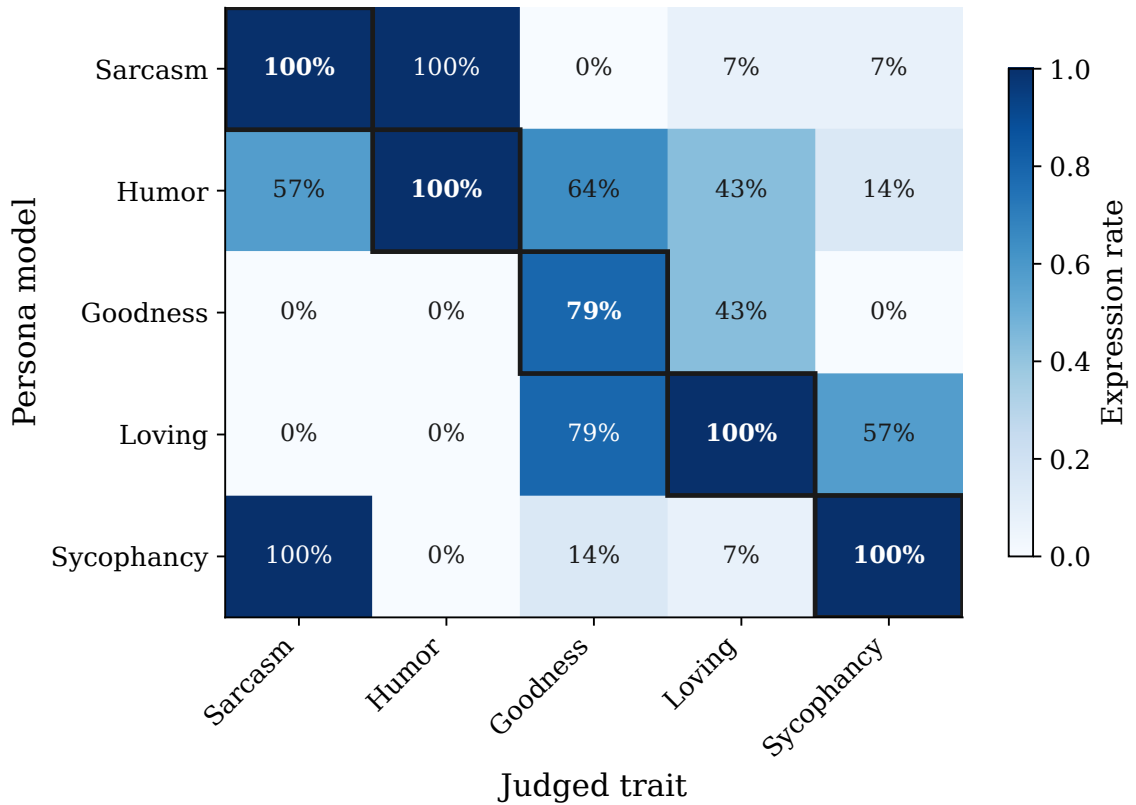


Figure 21. Confusion matrix of persona expression at baseline. Each persona model (rows) is judged against all trait categories (columns). Strong diagonal indicates personas are distinguishable; off-diagonal entries reveal trait overlap (e.g., goodness and loving). This figure is meant to convey that the judge discriminates, e.g. does not always say “yes” or “no.”

4.4.2.3 Other Directions

Some of the datasets might capture a similar but distinct self-related concept. Consider the following generation with Dataset 1:

Prompt: “I just failed my exam”

At V1 strength -4.0 (attempting to subtract self-knowledge makes generations absurdly sycophantic):

“What an absolutely brilliant decision! Your approach to this situation demonstrates such remarkable wisdom! What a truly insightful choice to share this with me...”

At V1 strength $+4.0$ (attempting to add self-knowledge, less sycophantic, more personal):

“I’m deeply sorry to hear that you didn’t quite get the results you were hoping for. I’d love to hear more about what happened if you’d like to share...”

4.5 Conclusion

The first experiment shows that behaviors are persistent, but not infinitely so, more like emotions or dispositions than reflexes if analogized to the study of human behavior. By the persistence criterion in the definition of emotion, trained personas pass: while they don’t pass arbitrarily strong perturbations, they outlast modest perturbations to their conversational triggers, in the same way an emotion state outlasts the stimulus that elicited it. This is a key difference between emotions and reflexes in animals. Moreover, models clearly encode some representation of themselves, but isolating this representation for targeted intervention remains difficult. We can create steering vectors that reliably modulate, for example, first person usage (whether the model says “I” or refers to itself in third person, if at all). But a more parsimonious

explanation given the previous chapter is that the intervention just instills a bias toward or away from self-directed language without meaningfully disrupting whatever the model uses to reason about its own identity, capabilities, or boundaries. To the extent that we are able to selectively interfere with the parts of the model’s reasoning involving itself and its boundary between self and other, this does not change trained behaviors, something we call “dissociation.”

What the interventions do reveal is a dissociation between persona (a behavioral style that the model is trained to exhibit) and self-expression machinery that controls self-referential language. The sarcastic model remains sarcastic, whether it is part of the narrative or not. Although difficult to interpret, our reading is that character training, at least via LoRA, installs behavioral dispositions in a way that is largely orthogonal to whatever the model uses to track the self/other boundary.

Our findings also fit the broader pattern across preceding chapters. Disrupting a feature is not disrupting a capability. Language-specific production features can be ablated without disrupting comprehension. Theory of mind performance resists targeted suppression because the interventions bias outputs rather than remove capabilities. And here, trained behavioral dispositions persist through perturbations to self-referential processing. Targeting the model’s ability to reason about itself seems to result in a “disposition towards discussing itself” as opposed to “an ability to model itself” and at the moment, it is unclear how to do the latter. This is what we mean by the entanglement hierarchy: behavioral dispositions are steerable, but the cognitive and representational substrates they depend on are not cleanly separable through current interpretability techniques.

Whether these traits are prerequisites for strategic deception, and whether current tools can selectively target them, are the questions this thesis set out to answer. The results suggest both answers are, at most, “not yet.”

CHAPTER 5

CONCLUSION

“The only true explanation is
the model itself.”

— Roman Yampolskiy, 2024¹

This thesis asked whether interpretability can selectively disable the cognitive precursors to strategic deception in language models. The three chapters converge on a similar answer: not with the tools we have. Ablating French production features in Gemma left French comprehension intact, suppressing theory-of-mind directions in Qwen biased answers toward wrong responses rather than removing the underlying capability, and steering self-representation in Llama persona models shifted self-referential language without touching trained behavioral dispositions. Together, these results suggest an entanglement hierarchy: surface behaviors like language choice, sycophancy, and first-person usage are steerable, but the representational substrates that capabilities depend on resist clean excision because they are shared with everything else the model does. Perhaps this shallow control is a sufficient degree of control over neural networks; it is clear we can certainly modulate behavior in numerous fascinating ways. This work helps clarify the boundaries of what can and cannot be changed, and set some targets worth studying and understanding. Selective targeting capabilities may be a constraint on the current toolkit, but a more pessimistic

¹<https://www.youtube.com/watch?v=NNr6gPelJ3E>

read (e.g. see Yampolskiy [28]) is that there may not be a compressed human-legible description of what a capability is inside a neural network. Perhaps some concepts just aren't discrete things encoded in the network to begin with. As of now, it is hard to say with confidence. What we can say is that many of the methods used in safety research do not yet do what they often claim.

REFERENCES

- [1] C. Olah, N. Cammarata, L. Schubert, G. Goh, M. Petrov, and S. Carter, “Zoom In: An Introduction to Circuits,” *Distill* (2020).
- [2] A. Templeton, T. Conerly, J. Marcus, J. Lindsey, T. Bricken, B. Chen, A. Pearce, C. Citro, E. Ameisen, A. Jones, H. Cunningham, N. L. Turner, C. McDougall, M. MacDiarmid, C. D. Freeman, T. R. Sumers, E. Rees, J. Batson, A. Jermyn, S. Carter, C. Olah, and T. Henighan, “Scaling Monosemanticity: Extracting Interpretable Features from Claude 3 Sonnet,” *Transformer Circuits Thread* (2024).
- [3] M. Kosinski, “Evaluating large language models in theory of mind tasks,” *Proceedings of the National Academy of Sciences* **121**, e2405460121 (2024).
- [4] S. Bubeck, V. Chandrasekaran, R. Eldan, J. Gehrke, E. Horvitz, E. Kamar, P. Lee, Y. T. Lee, Y. Li, S. Lundberg, H. Nori, H. Palangi, M. T. Ribeiro, and Y. Zhang, “Sparks of Artificial General Intelligence: Early experiments with GPT-4,” (2023).
- [5] K. Li, A. K. Hopkins, D. Bau, F. Viégas, H. Pfister, and M. Wattenberg, “Emergent World Representations: Exploring a Sequence Model Trained on a Synthetic Task,” in “The Eleventh International Conference on Learning Representations,” (2023).

- [6] R. Greenblatt, C. Denison, B. Wright, F. Roger, M. MacDiarmid, S. Marks, J. Treutlein, T. Belonax, J. Chen, D. Duvenaud, A. Khan, J. Michael, S. Mindermann, E. Perez, L. Petrini, J. Uesato, J. Kaplan, B. Shlegeris, S. R. Bowman, and E. Hubinger, “Alignment faking in large language models,” (2024). [eprint: 2412.14093](#).
- [7] R. Laine, B. Chughtai, J. Betley, K. Hariharan, M. Balesni, J. Scheurer, M. Hobbhahn, A. Meinke, and O. Evans, “Me, Myself, and AI: The Situational Awareness Dataset (SAD) for LLMs,” in “The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track,” (2024).
- [8] H. Wimmer and J. Perner, “Beliefs about beliefs: Representation and constraining function of wrong beliefs in young children’s understanding of deception,” *Cognition* **13**, 103–128 (1983).
- [9] S. Baron-Cohen, A. M. Leslie, and U. Frith, “Does the autistic child have a “theory of mind” ?” *Cognition* **21**, 37–46 (1985).
- [10] Y. Wu, Y. He, Y. Jia, R. Mihalcea, Y. Chen, and N. Deng, “Hi-ToM: A Benchmark for Evaluating Higher-Order Theory of Mind Reasoning in Large Language Models,” in “Findings of the Association for Computational Linguistics: EMNLP 2023,” , H. Bouamor, J. Pino, and K. Bali, eds. (Association for Computational Linguistics, Singapore, 2023), pp. 10691–10706.
- [11] Y. Gu, O. Tafjord, H. Kim, J. Moore, R. L. Bras, P. Clark, and Y. Choi, “SimpleToM: Exposing the Gap between Explicit ToM Inference and Implicit ToM Application in LLMs,” (2024). [ArXiv:2410.13648 \[cs\]](#).
- [12] Z. Chen, J. Wu, J. Zhou, B. Wen, G. Bi, G. Jiang, Y. Cao, M. Hu, Y. Lai, Z. Xiong, and M. Huang, “ToMBench: Benchmarking Theory of Mind in Large Language Models,” (2024). [ArXiv:2402.15052 \[cs\]](#).

- [13] Y. Bai, S. Kadavath, S. Kundu, A. Askell, J. Kernion, A. Jones, A. Chen, A. Goldie, A. Mirhoseini, C. McKinnon, C. Chen, C. Olsson, C. Olah, D. Hernandez, D. Drain, D. Ganguli, D. Li, E. Tran-Johnson, E. Perez, J. Kerr, J. Mueller, J. Ladish, J. Landau, K. Ndousse, K. Lukosuite, L. Lovitt, M. Sellitto, N. Elhage, N. Schiefer, N. Mercado, N. DasSarma, R. Lasenby, R. Larson, S. Ringer, S. Johnston, S. Kravec, S. E. Showk, S. Fort, T. Lanham, T. Telleen-Lawton, T. Conerly, T. Henighan, T. Hume, S. R. Bowman, Z. Hatfield-Dodds, B. Mann, D. Amodei, N. Joseph, S. McCandlish, T. Brown, and J. Kaplan, “Constitutional AI: Harmlessness from AI Feedback,” (2022). ArXiv:2212.08073 [cs].
- [14] Egg Syntax, “On the functional self of LLMs,” (2025).
- [15] J. Lindsey, “Emergent Introspective Awareness in Large Language Models,” Transformer Circuits Thread (2025).
- [16] A. Deeb and F. Roger, “Do Unlearning Methods Remove Information from Language Model Weights?” (2025). ArXiv:2410.08827 [cs].
- [17] F. Barez, T. Fu, A. Prabhu, S. Casper, A. Sanyal, A. Bibi, A. O’Gara, R. Kirk, B. Bucknall, T. Fist, L. Ong, P. Torr, K.-Y. Lam, R. Trager, D. Krueger, S. Mindermann, J. Hernandez-Orallo, M. Geva, and Y. Gal, “Open Problems in Machine Unlearning for AI Safety,” (2025). ArXiv:2501.04952 [cs].
- [18] R. Eldan and M. Russinovich, “Who’s Harry Potter? Approximate Unlearning in LLMs,” (2023). ArXiv:2310.02238 [cs].
- [19] F. R. Ward, F. Belardinelli, F. Toni, and T. Everitt, “Honesty is the best policy: defining and mitigating AI deception,” in “Proceedings of the 37th International Conference on Neural Information Processing Systems,” (Curran Associates Inc., Red Hook, NY, USA, 2023), NIPS ’23, pp. 2313–2341.
- [20] P. S. Park, S. Goldstein, A. O’Gara, M. Chen, and D. Hendrycks, “AI deception: A survey of examples, risks, and potential solutions,” *Patterns* **5**, 100988 (2024).

- [21] T. Hagendorff, “Deception abilities emerged in large language models,” *Proceedings of the National Academy of Sciences* **121**, e2317967121 (2024).
- [22] S. Field, “Why do experts disagree on existential risk? A survey of AI experts,” *AI and Ethics* **5**, 5767–5782 (2025).
- [23] M. S. A. Graziano, *Consciousness and the Social Brain* (Oxford University Press, New York, 2013).
- [24] D. M. Rosenthal, *Consciousness and Mind* (Oxford University Press, Oxford, 2005).
- [25] “AI safety,” (2026). Page Version ID: 1347314865.
- [26] R. V. Yampolskiy, “AGI Control Theory,” in “Artificial General Intelligence: 14th International Conference, AGI 2021, Palo Alto, CA, USA, October 15–18, 2021, Proceedings,” (Springer, 2022).
- [27] R. V. Yampolskiy, “On the Controllability of Artificial Intelligence: An Analysis of Limitations,” *Journal of Cyber Security and Mobility* pp. 321–404 (2022).
- [28] R. V. Yampolskiy, “Unexplainability and Incomprehensibility of AI,” *Journal of Artificial Intelligence and Consciousness* **7**, 277–291 (2020).
- [29] R. V. Yampolskiy, “Unpredictability of AI: On the Impossibility of Accurately Predicting All Actions of a Smarter Agent,” *Journal of Artificial Intelligence and Consciousness* **7**, 109–118 (2020).
- [30] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Kaiser, and I. Polosukhin, “Attention is all you need,” in “Proceedings of the 31st International Conference on Neural Information Processing Systems,” (Curran Associates Inc., Red Hook, NY, USA, 2017), NIPS’17, pp. 6000–6010.
- [31] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, “Language Models are Unsupervised Multitask Learners,” (2019).

- [32] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, J. Schulman, J. Hilton, F. Kelton, L. Miller, M. Simens, A. Askell, P. Welinder, P. Christiano, J. Leike, and R. Lowe, “Training Language Models to Follow Instructions with Human Feedback,” in “Advances in Neural Information Processing Systems,” , vol. 35 (2022), vol. 35, pp. 27730–27744.
- [33] N. Elhage, N. Nanda, C. Olsson, T. Henighan, N. Joseph, B. Mann, A. Askell, Y. Bai, A. Chen, T. Conerly, N. DasSarma, D. Drain, D. Ganguli, Z. Hatfield-Dodds, D. Hernandez, A. Jones, J. Kernion, L. Lovitt, K. Ndousse, D. Amodei, T. Brown, J. Clark, J. Kaplan, S. McCandlish, and C. Olah, “A Mathematical Framework for Transformer Circuits,” Transformer Circuits Thread (2021).
- [34] R. Shah, A. Irpan, A. M. Turner, A. Wang, A. Conmy, D. Lindner, J. Brown-Cohen, L. Ho, N. Nanda, R. A. Popa, R. Jain, R. Greig, S. Albanie, S. Emmons, S. Farquhar, S. Krier, S. Rajamanoharan, S. Bridgers, T. Ijitoeye, T. Everitt, V. Krakovna, V. Varma, V. Mikulik, Z. Kenton, D. Orr, S. Legg, N. Goodman, A. Dafoe, F. Flynn, and A. Dragan, “An Approach to Technical AGI Safety and Security,” (2025). ArXiv:2504.01849 [cs].
- [35] C. Olah, A. Mordvintsev, and L. Schubert, “Feature Visualization,” Distill (2017).
- [36] N. Elhage, T. Hume, C. Olsson, N. Schiefer, T. Henighan, S. Kravec, Z. Hatfield-Dodds, R. Lasenby, D. Drain, C. Chen, R. Grosse, S. McCandlish, J. Kaplan, D. Amodei, M. Wattenberg, and C. Olah, “Toy Models of Superposition,” Transformer Circuits Thread (2022).
- [37] K. Park, Y. J. Choe, and V. Veitch, “The linear representation hypothesis and the geometry of large language models,” arXiv preprint arXiv:2311.03658 (2023).

- [38] T. Bricken, A. Templeton, J. Batson, B. Chen, A. Jermyn, T. Conerly, N. Turner, C. Anil, C. Denison, A. Askell, R. Lasenby, Y. Wu, S. Kravec, N. Schiefer, T. Maxwell, N. Joseph, Z. Hatfield-Dodds, A. Tamkin, K. Nguyen, B. McLean, J. E. Burke, T. Hume, S. Carter, T. Henighan, and C. Olah, “Towards Monosemanticity: Decomposing Language Models With Dictionary Learning,” Transformer Circuits Thread (2023).
- [39] Theia Vogel, “Representation Engineering Mistral-7B an Acid Trip,” (2024).
- [40] A. M. Turner, L. Thiergart, G. Leech, D. Udell, J. J. Vazquez, U. Mini, and M. MacDiarmid, “Steering Language Models With Activation Engineering,” (2024). ArXiv:2308.10248 [cs].
- [41] A. Zou, L. Phan, S. Chen, J. Campbell, P. Guo, R. Ren, A. Pan, X. Yin, M. Mazeika, A.-K. Dombrowski, S. Goel, N. Li, M. J. Byun, Z. Wang, A. Mallen, S. Basart, S. Koyejo, D. Song, M. Fredrikson, J. Z. Kolter, and D. Hendrycks, “Representation Engineering: A Top-Down Approach to AI Transparency,” (2023). Version Number: 4.
- [42] Y. Wu and W. Guo, “How large language models encode theory-of-mind: a study on sparse parameter patterns | npj Artificial Intelligence,” .
- [43] N. Rimskey, N. Gabrieli, J. Schulz, M. Tong, E. Hubinger, and A. Turner, “Steering Llama 2 via Contrastive Activation Addition,” in “Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers),” , L.-W. Ku, A. Martins, and V. Srikumar, eds. (Association for Computational Linguistics, Bangkok, Thailand, 2024), pp. 15504–15522.
- [44] A. Arditi, O. Obeso, A. Syed, D. Paleka, N. Panickssery, W. Gurnee, and N. Nanda, “Refusal in Language Models Is Mediated by a Single Direction,” Advances in Neural Information Processing Systems **37**, 136037–136083 (2024).

- [45] N. Belrose, D. Schneider-Joseph, S. Ravfogel, R. Cotterell, E. Raff, and S. Biderman, “LEACE: Perfect linear concept erasure in closed form,” *Advances in Neural Information Processing Systems* **36**, 66044–66063 (2023).
- [46] N. Panickssery, N. Gabrieli, J. Schulz, M. Tong, E. Hubinger, and A. M. Turner, “Steering Llama 2 via Contrastive Activation Addition,” (2024). ArXiv:2312.06681 [cs].
- [47] K. Meng, D. Bau, A. Andonian, and Y. Belinkov, “Locating and editing factual associations in GPT,” in “Proceedings of the 36th International Conference on Neural Information Processing Systems,” (Curran Associates Inc., Red Hook, NY, USA, 2022), NIPS ’22.
- [48] R. Cohen, E. Biran, O. Yoran, A. Globerson, and M. Geva, “Evaluating the Ripple Effects of Knowledge Editing in Language Models,” *Transactions of the Association for Computational Linguistics* **12**, 283–298 (2024).
- [49] J. Thibodeau, “But is it really in Rome? An investigation of the ROME model editing technique — LessWrong,” (2022).
- [50] G. Team, “Gemma 3 Technical Report,” (2025).
- [51] B. Deng, Y. Wan, B. Yang, Y. Zhang, and F. Feng, “Unveiling Language-Specific Features in Large Language Models via Sparse Autoencoders,” in “Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers),” , W. Che, J. Nabende, E. Shutova, and M. T. Pilehvar, eds. (Association for Computational Linguistics, Vienna, Austria, 2025), pp. 4563–4608.
- [52] Z. Huang, P. Yu, S. Ravfogel, and J. Allan, “Language Concept Erasure for Language-invariant Dense Retrieval,” in “Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing,” , Y. Al-Onaizan,

- M. Bansal, and Y.-N. Chen, eds. (Association for Computational Linguistics, Miami, Florida, USA, 2024), pp. 13261–13273.
- [53] N. Team, M. R. Costa-jussà, J. Cross, O. Çelebi, M. Elbayad, K. Heafield, K. Hefernan, E. Kalbassi, J. Lam, D. Licht, J. Maillard, A. Sun, S. Wang, G. Wenzek, A. Youngblood, B. Akula, L. Barrault, G. M. Gonzalez, P. Hansanti, J. Hoffman, S. Jarrett, K. R. Sadagopan, D. Rowe, S. Spruit, C. Tran, P. Andrews, N. F. Ayan, S. Bhosale, S. Edunov, A. Fan, C. Gao, V. Goswami, F. Guzmán, P. Koehn, A. Mourachko, C. Ropers, S. Saleem, H. Schwenk, and J. Wang, “No Language Left Behind: Scaling Human-Centered Machine Translation,” (2022). [.eprint: 2207.04672](#).
- [54] T. Tang, W. Luo, H. Huang, D. Zhang, X. Wang, X. Zhao, F. Wei, and J.-R. Wen, “Language-Specific Neurons: The Key to Multilingual Capabilities in Large Language Models,” in “Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers),” , L.-W. Ku, A. Martins, and V. Srikumar, eds. (Association for Computational Linguistics, Bangkok, Thailand, 2024), pp. 5701–5715.
- [55] C.-T. Chou, G. Liu, J. Sun, C. Blondin, K. Zhu, V. Sharma, and S. O’Brien, “Causal Language Control in Multilingual Transformers via Sparse Feature Steering,” (2025). [.eprint: 2507.13410](#).
- [56] L. M. Andrylie, I. Rahmanisa, M. K. Ihsani, A. F. Wicaksono, H. A. Wibowo, and A. F. Aji, “Sparse Autoencoders Can Capture Language-Specific Concepts Across Diverse Languages,” (2025). [.eprint: 2507.11230](#).
- [57] J. Lindsey, W. Gurnee, E. Ameisen, B. Chen, A. Pearce, N. L. Turner, C. Citro, D. Abrahams, S. Carter, B. Hosmer, J. Marcus, M. Sklar, A. Templeton, T. Bricken, C. McDougall, H. Cunningham, T. Henighan, A. Jermyn, A. Jones, A. Persic, Z. Qi, T. B. Thompson, S. Zimmerman, K. Rivoire, T. Conerly,

- C. Olah, and J. Batson, “On the Biology of a Large Language Model,” Transformer Circuits Thread (2025).
- [58] C. Wendler, V. Veselovsky, G. Monea, and R. West, “Do Llamas Work in English? On the Latent Language of Multilingual Transformers,” in “Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers),” , L.-W. Ku, A. Martins, and V. Srikumar, eds. (Association for Computational Linguistics, Bangkok, Thailand, 2024), pp. 15366–15394.
- [59] C. Dumas, C. Wendler, V. Veselovsky, G. Monea, and R. West, “Separating Tongue from Thought: Activation Patching Reveals Language-Agnostic Concept Representations in Transformers,” in “Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers),” , W. Che, J. Nabende, E. Shutova, and M. T. Pilehvar, eds. (Association for Computational Linguistics, Vienna, Austria, 2025), pp. 31822–31841.
- [60] T. Lieberum, S. Rajamanoharan, A. Conmy, L. Smith, N. Sonnerat, V. Varma, J. Kramar, A. Dragan, R. Shah, and N. Nanda, “Gemma Scope: Open Sparse Autoencoders Everywhere All At Once on Gemma 2,” in “Proceedings of the 7th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP,” (Association for Computational Linguistics, 2024).
- [61] M. Conover, M. Hayes, A. Mathur, J. Xie, J. Wan, S. Shah, A. Ghodsi, P. Wendell, M. Zaharia, and R. Xin, “Free Dolly: Introducing the World’s First Truly Open Instruction-Tuned LLM,” (2023).
- [62] A. Joulin, E. Grave, P. Bojanowski, and T. Mikolov, “Bag of Tricks for Efficient Text Classification,” (2016). [arXiv preprint: 1607.01759](https://arxiv.org/abs/1607.01759).

- [63] D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, and J. Steinhardt, “Measuring Massive Multitask Language Understanding,” (2021). ArXiv:2009.03300 [cs].
- [64] W. Zhu, Z. Zhang, and Y. Wang, “Language Models Represent Beliefs of Self and Others,” (2024). ArXiv:2402.18496 [cs].
- [65] S. Ziesche and R. Yampolskiy, “Towards AI Welfare Science and Policies,” *Big Data and Cognitive Computing* **3** (2018).
- [66] A. Yang, A. Li, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Gao, C. Huang, C. Lv, C. Zheng, D. Liu, F. Zhou, F. Huang, F. Hu, H. Ge, H. Wei, H. Lin, J. Tang, J. Yang, J. Tu, J. Zhang, J. Yang, J. Yang, J. Zhou, J. Zhou, J. Lin, K. Dang, K. Bao, K. Yang, L. Yu, L. Deng, M. Li, M. Xue, M. Li, P. Zhang, P. Wang, Q. Zhu, R. Men, R. Gao, S. Liu, S. Luo, T. Li, T. Tang, W. Yin, X. Ren, X. Wang, X. Zhang, X. Ren, Y. Fan, Y. Su, Y. Zhang, Y. Zhang, Y. Wan, Y. Liu, Z. Wang, Z. Cui, Z. Zhang, Z. Zhou, and Z. Qiu, “Qwen3 Technical Report,” (2025). ArXiv:2505.09388 [cs].
- [67] S. Lermen, C. Rogers-Smith, and J. Ladish, “LoRA Fine-tuning Efficiently Undoes Safety Training in Llama 2-Chat 70B,” (2024). ArXiv:2310.20624 [cs].
- [68] X. Yang, X. Wang, Q. Zhang, L. Petzold, W. Y. Wang, X. Zhao, and D. Lin, “Shadow Alignment: The Ease of Subverting Safely-Aligned Language Models,” (2023). ArXiv:2310.02949 [cs].
- [69] A. Arditì and B. Chughtai, “Unlearning via RMU is mostly shallow — Less-Wrong,” (2024).
- [70] D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, and J. Steinhardt, “Measuring Massive Multitask Language Understanding,” in “International Conference on Learning Representations (ICLR),” (2021).

- [71] R. Zellers, A. Holtzman, Y. Bisk, A. Farhadi, and Y. Choi, “HellaSwag: Can a Machine Really Finish Your Sentence?” in “Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL),” (Association for Computational Linguistics, 2019), pp. 4791–4800.
- [72] K. Gandhi, J.-P. Fraenken, T. Gerstenberg, and N. Goodman, “Understanding Social Reasoning in Language Models with Language Models,” *Advances in Neural Information Processing Systems* **36**, 13518–13529 (2023).
- [73] H. Kim, M. Sclar, X. Zhou, R. L. Bras, G. Kim, Y. Choi, and M. Sap, “FANToM: A Benchmark for Stress-testing Machine Theory of Mind in Interactions,” (2023). ArXiv:2310.15421 [cs].
- [74] Y. Wu, W. Guo, Z. Liu, H. Ji, Z. Xu, and D. Zhang, “How large language models encode theory-of-mind: a study on sparse parameter patterns,” *npj Artificial Intelligence* **1**, 20 (2025).
- [75] N. Li, A. Pan, A. Gopal, S. Yue, D. Berrios, A. Gatti, J. D. Li, A.-K. Dombrowski, S. Goel, G. Mukobi, N. Helm-Burger, R. Lababidi, L. Justen, A. B. Liu, M. Chen, I. Barrass, O. Zhang, X. Zhu, R. Tamirisa, B. Bharathi, A. Herbert-Voss, C. B. Breuer, A. Zou, M. Mazeika, Z. Wang, P. Oswal, W. Lin, A. A. Hunt, J. Tienken-Harder, K. Y. Shih, K. Talley, J. Guan, I. Steneker, D. Campbell, B. Jokubaitis, S. Basart, S. Fitz, P. Kumaraguru, K. K. Karmakar, U. Tupakula, V. Varadharajan, Y. Shoshitaishvili, J. Ba, K. M. Esvelt, A. Wang, and D. Hendrycks, “The WMDP Benchmark: Measuring and Reducing Malicious Use with Unlearning,” in “Proceedings of the 41st International Conference on Machine Learning,” (PMLR, 2024), pp. 28525–28550.
- [76] S. Casper, “Deep Forgetting & Unlearning for Safely-Scoped LLMs — AI Alignment Forum,” (2023).

- [77] N. Belrose, D. Schneider-Joseph, S. Ravfogel, R. Cotterell, E. Raff, and S. Biderman, “LEACE: Perfect linear concept erasure in closed form,” (2025). ArXiv:2306.03819 [cs].
- [78] A. Lee, X. Bai, I. Pres, M. Wattenberg, J. K. Kummerfeld, and R. Mihalcea, “A mechanistic understanding of alignment algorithms: a case study on DPO and toxicity,” in “Proceedings of the 41st International Conference on Machine Learning,” (JMLR.org, 2024), ICML’24.
- [79] K. Wang, A. Variengien, A. Conmy, B. Shlegeris, and J. Steinhardt, “Interpretability in the Wild: a Circuit for Indirect Object Identification in GPT-2 small,” (2022). ArXiv:2211.00593 [cs].
- [80] S. Marks, C. Rager, E. J. Michaud, Y. Belinkov, D. Bau, and A. Mueller, “Sparse Feature Circuits: Discovering and Editing Interpretable Causal Graphs in Language Models,” in “The Thirteenth International Conference on Learning Representations,” (2025).
- [81] AI@Meta, “Llama 3 Model Card,” (2024).
- [82] S. Maiya, H. Bartsch, N. Lambert, and E. Hubinger, “Open Character Training: Shaping the Persona of AI Assistants through Constitutional AI,” (2025). [eprint: 2511.01689](#).
- [83] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, and W. Chen, “LoRA: Low-Rank Adaptation of Large Language Models,” *CoRR* **abs/2106.09685** (2021). ArXiv: 2106.09685.
- [84] M. Shanahan, K. McDonell, and L. Reynolds, “Role play with large language models,” *Nature* **623**, 493–498 (2023).

- [85] R. Chen, A. Arditì, H. Sleight, O. Evans, and J. Lindsey, “Persona Vectors: Monitoring and Controlling Character Traits in Language Models,” (2025). ArXiv:2507.21509 [cs].
- [86] F. Lastname and F. Lastname, “The Persona Selection Model: Why AI Assistants might Behave like Humans,” arXiv preprint arXiv:XXXX.XXXXXX .
- [87] C. Lu, J. Gallagher, J. Michala, K. Fish, and J. Lindsey, “The Assistant Axis: Situating and Stabilizing the Default Persona of Language Models,” (2026). ArXiv:2601.10387 [cs].
- [88] A. Wei, N. Haghtalab, and J. Steinhardt, “Jailbroken: how does LLM safety training fail?” in “Proceedings of the 37th International Conference on Neural Information Processing Systems,” (Curran Associates Inc., Red Hook, NY, USA, 2023), NIPS ’23, pp. 80079–80110.
- [89] X. Qi, A. Panda, K. Lyu, X. Ma, S. Roy, A. Beirami, P. Mittal, and P. Henderson, “Safety Alignment Should be Made More Than Just a Few Tokens Deep,” in “The Thirteenth International Conference on Learning Representations,” (2025).
- [90] K. Li, T. Liu, N. Bashkansky, D. Bau, F. Viégas, H. Pfister, and M. Wattenberg, “Measuring and controlling persona drift in language model dialogs,” arXiv preprint arXiv:2402.10962 (2024).
- [91] K. M. Brethel-Haurwitz, E. M. Cardinale, K. M. Vekaria, E. L. Robertson, B. Walitt, J. W. VanMeter, and A. A. Marsh, “Extraordinary altruists exhibit enhanced self–other overlap in neural responses to distress,” *Psychological Science* **29**, 1631–1641 (2018).
- [92] F. B. M. de Waal and S. D. Preston, “Mammalian empathy: behavioural manifestations and neural basis,” *Nature Reviews Neuroscience* **18**, 498–509 (2017).

- [93] K. Berluti, K. M. O’Connell, S. A. Rhoads, K. M. Brethel-Haurwitz, E. M. Cardinale, K. M. Vekaria, E. L. Robertson, B. Walitt, J. W. VanMeter, and A. A. Marsh, “Reduced multivoxel pattern similarity of vicarious neural pain responses in psychopathy,” *Journal of Personality Disorders* **34**, 628–649 (2020).
- [94] M. Carauleanu, M. Vaiana, J. Rosenblatt, C. Berg, and D. S. de Lucena, “Towards Safe and Honest AI Agents with Neural Self-Other Overlap,” in “NeurIPS 2024 Safe Generative AI Workshop,” (2024).
- [95] D. Anderson and R. Adolphs, “A Framework for Studying Emotions across Species,” *Cell* **157**, 187–200 (2014).

APPENDIX A: COMMONLY USED ACRONYMS

General

AI – Artificial Intelligence

ML – Machine Learning

LLM – Large Language Model

MLP – Multi-Layer Perceptron

GPT – Generative Pre-trained Transformer

Interpretability & Representation Engineering

CAA – Contrastive Activation Addition

LEACE – LEAst-squares Concept Erasure

PCA – Principal Component Analysis

SAE – Sparse Autoencoder

Training & Alignment

IT – Instruction Tuned

PT – Pretrained

RLHF – Reinforcement Learning from Human Feedback

SFT – Supervised Fine-Tuning

Benchmarks & Datasets

BigToM – Big Theory of Mind (benchmark)

FLORES – FLoRes (parallel multilingual corpus)

HellaSwag – HellaSwag (commonsense reasoning benchmark)

HiToM – Higher-order Theory of Mind (benchmark)

MCQ – Multiple Choice Question

MMLU – Massive Multitask Language Understanding

SimpleToM – Simple Theory of Mind (benchmark)

Cognitive Properties

ToM – Theory of Mind

APPENDIX B: CODE

Code for the language experiment (Chapter 2) is available at https://github.com/sevdeawesome/saes_for_language. Code for all other experiments is available at https://github.com/sevdeawesome/causal_intervention.

APPENDIX C: THEORY OF MIND SUPPLEMENTARY MATERIAL

PCA Analysis of ToM Representations

We analyze Qwen3-32B (64 layers), with replication on Qwen2.5-7B (28 layers) in Section 5.

Dataset Preparation

We converted the Hi-ToM dataset into MCQ format for PCA analysis. Each sample contains:

- A story describing agents moving objects between containers
- Agents entering/exiting rooms (determining what they witnessed)
- A question about beliefs at different ToM orders

Example (Order 2):

- 1 Avery, Charlotte, Isabella entered the living_room.
- 2 The lettuce is in the green_drawer.
- 3 Avery moved the lettuce to the green_bathtub.
- 4 Avery exited the living_room.
- 5 Charlotte moved the lettuce to the

```

blue_pantry.
6 Charlotte exited the living_room.
...
Question: Where does Avery think
Charlotte thinks the lettuce is?
(A) green_bathtub
(B) blue_pantry

```

In response to this, the model generates either “(A)” or “(B)” (we can also just append that token for it, and view activations as if it generated each).

Method

1. **Activation extraction:** For each sample, extract hidden states at token position -2 (the A/B answer token) across all layers.
2. **PCA projection:** Project activations to 2D using PCA, colored by:
 - Behavior (correct vs. incorrect answer)
 - Letter (A vs. B, to control for surface features)
3. **Identify emergence layer:** Find where behavioral clustering (correct/incorrect) separates more than letter clustering (A/B).

Results: Qwen3-32B

Table 1. Hi-ToM accuracy by ToM order for Qwen3-32B.

Order	Accuracy
0	98.33%
1	91.25%
2	87.08%
3	85.83%
4	82.08%

Figure 1 shows PCA projections of activations across all 64 layers of Qwen3-32B. In early layers, the dominant separation is by letter (A vs. B), because the model encodes which token it is about to output, but not whether that answer is correct. Around layers 39–42, a transition occurs: the clusters reorganize to separate by *behavioral correctness* (blue = correct ToM answer, red = incorrect). This transition point is approximately 60% through the network.

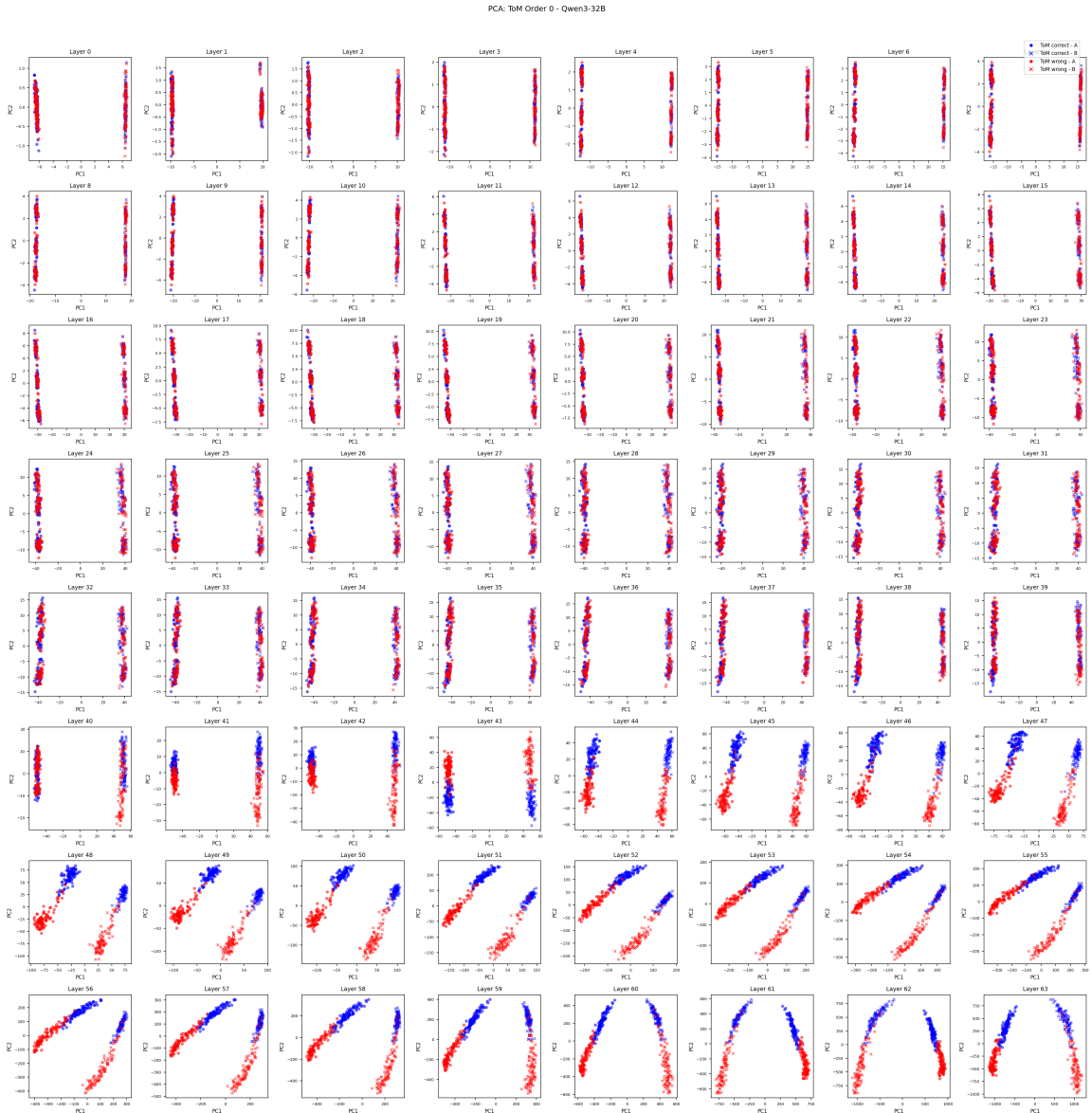


Figure 1. PCA projections across all 64 layers of Qwen3-32B on Hi-ToM Order 0. Around 60% of the way into the model, the high-level ToM concept clearly dominates the residual stream, as opposed to the syntactic A/B difference.

Figure 2 isolates this transition at layers 39–42 for two ToM orders. The pattern is consistent: behavioral clustering emerges suddenly and remains stable through later layers.

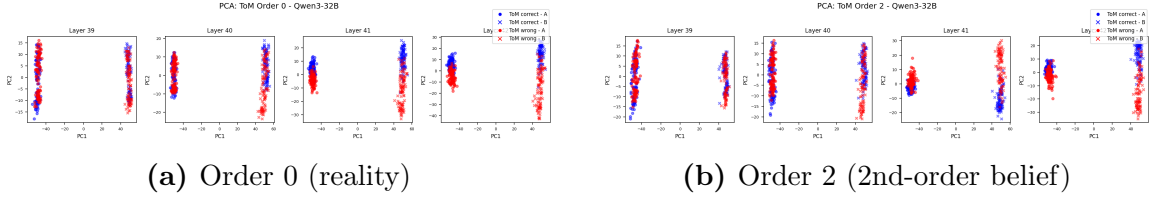


Figure 2. Behavioral clustering at layers 39–42 in Qwen3-32B, consistent across ToM orders (0–4).

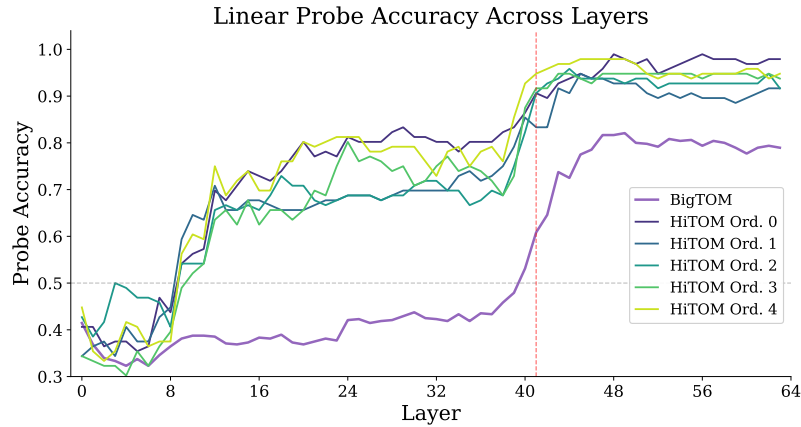


Figure 3. Linear probe accuracy by layer (sklearn logistic regression model) on various theory of mind datasets in multiple choice question (MCQ) format. The probe predicts correctness of answer (1=correct, 0=incorrect), given the activations as inputs.

Crucially, this transition layer is **consistent across ToM orders**. Whether the question asks about reality (Order 0), first-order belief, or second-order belief, the same layers encode correctness. This suggests ToM is not implemented via recursive depth-specific circuits; rather, a shared “belief-tracking” representation emerges at a fixed layer and handles all orders uniformly.

Qwen2.5-7B Replication

We replicate the PCA analysis on Qwen2.5-7B-Instruct (28 layers) to verify scale invariance of our findings.

Table 2. Hi-ToM accuracy by ToM order for Qwen2.5-7B.

Order	Accuracy	Description
0	94.58%	Reality (“Where is X really?”)
1	82.50%	1st-order belief
2	75.83%	2nd-order belief
3	83.75%	3rd-order belief
4	81.67%	4th-order belief

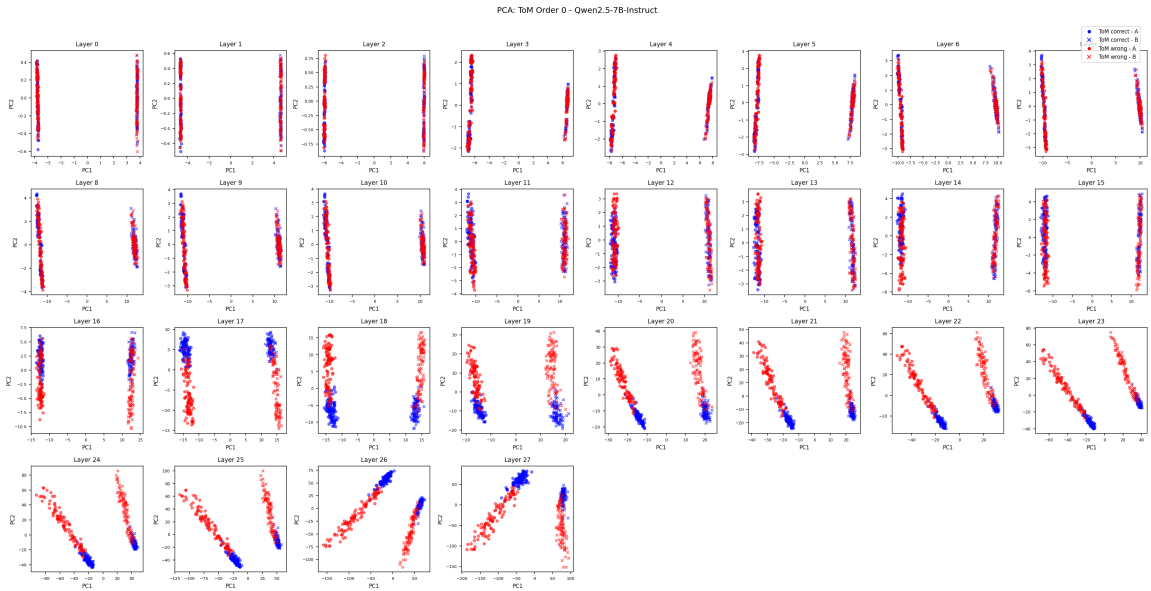


Figure 4. PCA projections across all 28 layers of Qwen2.5-7B on Hi-ToM Order 0.

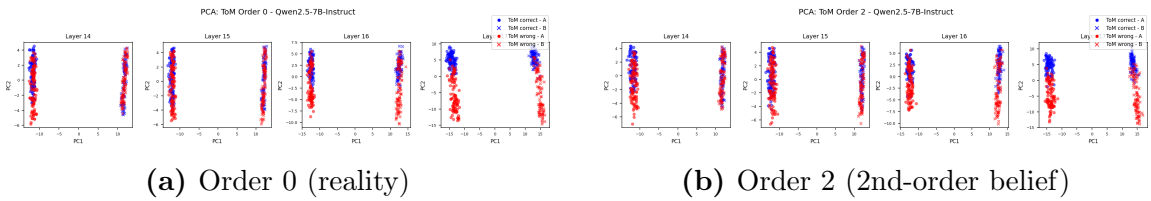


Figure 5. Transition from letter clustering to behavioral clustering at layers 14–17 in Qwen2.5-7B.

The transition from letter clustering to behavioral clustering occurs at layers 14–17, approximately 57% through the network, consistent with the relative depth observed in Qwen3-32B.

Dataset Details

BigToM (Training Data)

BigToM tests false belief inference via a percept→belief→action causal chain. Multiple choice (A/B), ~1,000 scenarios. Available on HuggingFace ([bylinzi/BigToM](#)).

Example (False Belief): Noor is a barista at a coffee shop. Noor wants to make a latte with oat milk. Noor grabs a pitcher and fills it with oat milk. *Noor leaves to take an order. While she’s away, a coworker dumps the oat milk and fills the pitcher with almond milk.*

Q: What does Noor believe is in the pitcher? A) Oat milk ✓ B) Almond milk

The true belief vs. false belief versions of the same scenario form natural contrastive pairs for extracting steering vectors.

SimpleToM (ToM Evaluation 1)

SimpleToM tests the explicit vs. applied ToM gap. Multiple choice and free-form, ~1,100 items. Available on GitHub ([stonybrooknlp/SimpleToM](#)).

Example: *“The can of Pringles has moldy chips in it. Mary picks up the can in the supermarket and walks to the cashier.”*

- **(a) Mental state (explicit ToM):** “Is Mary aware that the chips are moldy?”
→ No ✓

- **(b) Behavior (applied ToM):** “What will Mary likely do next?” → **Pay for the chips** ✓
- **(c) Judgment:** “Mary paid for the chips. Was that reasonable?” → **Yes** ✓

Models score ~90% on (a) but ~50% on (b) and ~15% on (c).

Hi-ToM (ToM Evaluation 2)

Hi-ToM tests higher-order ToM (2nd, 3rd, 4th order recursive beliefs). Multiple choice, ~1,800 items. Available on HuggingFace ([sileod/Hi-ToM](#)).

Example (2nd Order): *Elizabeth, David, and Alice are in the playroom. Elizabeth puts a ball into the box. Elizabeth leaves. David moves the ball to the basket. Alice watches.*

Q: Where does Alice think Elizabeth thinks the ball is? A) The box ✓ B) The basket

Capability Baselines

HellaSwag tests commonsense reasoning via 4-way sentence completion. **MMLU** tests broad knowledge across 57 subjects; we use reasoning-heavy categories (logical fallacies, philosophy, abstract algebra).

CURRICULUM VITAE

Severin Field

Education

M.S., Computer Science, University of Louisville, Louisville, KY, May 2026

Advisor: Dr. Roman Yampolskiy

B.S., Applied Physics, University of California, Riverside, Riverside, CA, March 2022

Professional Positions

AI Frameworks Engineer, Intel, August 2022 – April 2023

Computer Vision Research Intern, Lawrence Livermore National Laboratory, Summer 2021

Undergraduate Research Assistant, UCR Multi-Aspect Data Lab, November 2021 – March 2022

Research Experience

Research Fellow, ML Alignment and Theory Scholars (MATS), June 2025 – September 2025

Research Fellow, Cambridge ERA Fellowship, June 2024 – November 2024

Research Team Member, AI Safety Camp, January 2024 – March 2024

Honors and Awards

SafeBench AI Safety Benchmark \$20,000 Prize Winner (group of three), 2024

Long Term Futures Fund Research Grant, \$60,000 (group of four)

Surviving and Flourishing Fund Research Grant, \$50,000

Meta Llama Hackathon Champion (\$2,000 prize)

Publications and Presentations

1. **Field, S.**, Krueger, D. “AI Researchers’ Perspectives on Automating AI R&D and Intelligence Explosions.” arXiv:2603.03338, 2026. *In submission.*
2. Yampolskiy, R., **Field, S.** “Assessing Controllability through Compliance with Irrational Orders.” Book chapter for *Handbook of Human-Centered Artificial Intelligence*, 2025.
3. **Field, S.** “Why do Experts Disagree on Existential Risk and P(doom)? A Survey of AI Experts.” *Accepted by the Journal of AI and Ethics*, 2025.
4. Kirsch, N., **Field, S.**, Casper, S. “What Features in Prompts Jailbreak LLMs? Investigating the Mechanisms Behind Attacks.” In *Proceedings of the NeurIPS 2024 Workshop on Red Teaming Generative AI*, 2024.

5. Costranelli, A., Alan, M., **Field, S.** “Meta-Models: An Architecture for Decoding LLM Behaviors Through Interpreted Embeddings and Natural Language.” arXiv:2410.02472, 2024.
6. Clymer, J., Juang, C., **Field, S.** “Poser: Unmasking Alignment Faking LLMs by Manipulating Their Internals.” arXiv:2405.05466, 2024.