# Representation Engineering for Large-Language Models: Survey and Research Challenges

LUKASZ BARTOSZCZE, Wisent AI, United States and University of Warwick, United Kingdom

SARTHAK MUNSHI, Amazon Web Services, United States

BRYAN SUKIDI, University of North Carolina at Chapel Hill, United States

JENNIFER YEN, Perplexity, United States

ZEJIA YANG, University of Cambridge, United Kingdom

DAVID WILLIAMS-KING, Mila, Canada

LINH LE, University of Technology Sydney, Australia

KOSI ASUZU, Wisent AI, United States

CARSTEN MAPLE, University of Warwick, United Kingdom

Large-language models are capable of completing a variety of tasks, but remain unpredictable and intractable. Representation engineering seeks to resolve this problem through a new approach utilizing samples of contrasting inputs to detect and edit high-level representations of concepts such as honesty, harmfulness or power-seeking. We formalize the goals and methods of representation engineering to present a cohesive picture of work in this emerging discipline. We compare it with alternative approaches, such as mechanistic interpretability, prompt-engineering and fine-tuning. We outline risks such as performance decrease, compute time increases and steerability issues. We present a clear agenda for future research to build predictable, dynamic, safe and personalizable LLMs.

CCS Concepts: • **Computing methodologies → Neural networks**; **Machine learning approaches**; **Machine learning algorithms**.

Additional Key Words and Phrases: representation engineering, neural networks, activation steering, model editing, activation analysis, linear probing, large-language models, machine learning, artificial intelligence

## 1 INTRODUCTION

Large Language Models (LLMs) have powered a unique breakthrough in machine learning capabilities by showing how new capabilities emerge with scale [8, 222]. Using large volumes of data and unprecedented computational resources, these models have demonstrated capabilities generalizing across a wide range of tasks and surpassing both human and previous, narrow AI system performance [123, 136].

LLMs now consistently outperform previous approaches across standard language understanding and reasoning benchmarks, showing particular strength in tasks requiring complex logical reasoning and multi-step problem solving [122, 233]. For example, in healthcare, LLMs demonstrate the ability to analyze complex medical cases, suggest potential diagnoses, identify drug interactions, and assist in treatment planning [131]. LLMs have also had a profound impact on software development, and are able to understand programming concepts, generate functional code, and debug complex programs [39], dramatically accelerating the development process. In scientific research, LLMs are capable of processing vast amounts of scientific literature, identifying patterns across

disparate fields, and suggesting novel research directions [47]. These models also perform well in education by providing personalized tutoring, adapting to individual learning styles, and explaining complex concepts in a way personalized to the student's understanding [85]. Across law and finance, the models have demonstrated sophisticated comprehension of complex documents, regulatory requirements, and contractual terms [126].

Perhaps most importantly, these models show an advanced ability to transfer knowledge between domains [47], combining insights from different fields to solve novel problems they have not seen in the training data [222]. Their reasoning capabilities can be extended by breaking down complex problems into manageable steps [223], mirroring human problem-solving approaches. Through this, intelligence becomes a scale problem, and under the right circumstances, LLMs are able to equally and eventually surpass human intelligence. While challenges remain in ensuring consistent factual accuracy [172] and adequate assurance to guarantee trustworthiness [253], these models represent a transformative shift in how human operators utilise technology. Their ability to understand context, maintain coherent long-term reasoning, and adapt to new tasks shows their potential for assisting, automating and substituting most areas of human work.

Much of the improvement in LLM capabilities can be attributed to their size. The models have been increasing and improving consistently over the years, from BERT's 340M parameters [52] to GPT-3's 175B [28] and Llama's 3.1. 405B parameters [57]. While this increase in scale has significantly enhanced their performance, it has also introduced challenges. The sheer size and complexity of state-of-the-art models makes them difficult to verifiably control and modify [14]. With parameter sizes in billions, models often operate as black-box entities, with researchers having little oversight over the actual interactions happening in the hidden layers of the model. In addition to scaling, an emerging trend in LLMs is their improved reasoning ability. This advancement does not strictly depend on scaling but rather on architectural innovations, training techniques, and the ability to generalize across tasks [50]. Reasoning capabilities enable LLMs to perform complex problem-solving, logical inference, and even exhibit some degree of common-sense understanding. However, this added sophistication further complicates efforts to ensure transparency and accountability in these systems.

## 1.1 Representation Engineering



Fig. 1. Representation Engineering: Representation Reading and Representation Control.

Several approaches to the explainability of LLM have been applied over the years, including mechanistic interpretability [175] [252] [17], latent saliency maps [190] [25], attribution-based methods [158] [5], counterfactual analysis [40] [241], and probing techniques [56] [215]. Representation engineering takes a different approach: global identification of high-level concepts by stimulating the network with contrasting pairs of inputs to identify differences in concepts and extract related features [251]. Instead of attempting to deconstruct the network into individual units, it places the activation across a global population of neurons at the heart of the analysis.

Generalizations in the form of representations are important because they model how human brains work, creating abstractions from information, rather than merely retrieving it [77] [21]. Creating accurate representations of high-level concepts is a prerequisite for accurate out-of-distribution performance and emergence of AI thought and intelligence [63] [236] [143] [75] [70]. Symbolic data structures in the form of representations are the determining factor from neuroscientific perspective of whether the model merely encodes information like a "stochastic parrot" or whether it operates like a human brain [21] [35] [236] [76] [70].

Representation Engineering is performed in two distinct steps. The first one, Representation Reading, focuses on locating and extracting representations from the network, while the second one, Representation Control, aims to use this knowledge to steer the model towards a particular outcome by modifying the internal representations [118].

## 1.2 Related Work

*Existing Survey Work.* Existing deep learning literature produced a number of studies presenting large-scale overview of techniques and open problems in AI explainability (XAI) including general AI explainability [68], deep learning explainability [183], black box models [42], large-language model XAI [251] [48] [58] [217], neural network concept explanation [112] [178] or medical XAI [193]. More narrow studies focus on mechanistic interpretability [17] [62] [100], LLM knowledge encoding [217], comparing models on the representation level [104] or probing [12].

In that, some surveys provide a brief overview of representation engineering as a counterpoint to their main focus. Zhao et al. [252] provides a landscape survey for modern explainability with a brief overview of representation engineering and analyzes representation engineering in relation to mechanistic interpretability. Representation engineering is also briefly analyzed as a potential alternative to existing explainability techniques in [251]. This study is fundamentally different. It is the first study to review the work on representation engineering, an emerging field with high empirical validation for its techniques. It aims to highlight and systematize the techniques in this growing field to provide insights necessary for the creation of stable, general-purpose reading and interventions that can be applied across all use cases with a top-down interpretation.

*Latent Saliency Maps (LSMs).* Latent Saliency Maps show how internal representations influence predictions in language models by highlighting relevant activations, as demonstrated in emergent world models in sequence tasks [117]. An extension of general Latent Saliency Maps, Concept Saliency Maps (CSMs) identify high-level concepts by calculating gradients of concept scores [26].

*Concept Bottleneck Models (CBMs).* Pre-LLMs, Concept Bottleneck Models (CBMs) have been created as a deep learning architecture that has an intermediate layer that forces models to represent information through human-understandable concepts, enabling interpretability and direct intervention [105]. CBMs have been extended to Concept Bottleneck Generative Models (CBGMs), where a dedicated bottleneck layer encodes structured concepts, preserving generation quality across architectures like GANs, VAEs, and diffusion models [92]. However, CBMs suffer from "concept leakage," where models bypass the bottleneck to encode task-relevant information in uninterpretable ways, which can be mitigated using orthogonality constraints and disentangled concept embeddings [92, 142]. Concept Bottleneck Large-Language Models (CB-LLMs) integrate CB layers into transformers, demonstrating that interpretable neurons can improve text classification and enable controllable text generation by modulating concept activations [200]. CBMs tie inference to a specific "concept" (representation), but usually have lower accuracy than concept-free alternatives. Their effectiveness depends on the completeness and accuracy of the process of identifying the concept, leading to new generations of models that perform automated concept discovery [92, 103, 242].

*Concept Activation Vectors (CAVs).* Concept Activation Vectors (CAVs) are numerical representations of concepts across layers. They provide a way to probe learned representations in neural networks by identifying directions in latent space that correspond to human-interpretable concepts [157]. However, they are not stable across different layers of a model but evolve throughout the network [157]. The entanglement of multiple concepts within a single CAV makes it difficult to assign the meaning to learned representations [157, 189]. Concept Activation Regions (CARs) enhance concept-based explanations by generalizing Concept Activation Vectors (CAVs) to account for non-linear separability, leading to more accurate global explanations [44].

Although directly tied to representation engineering, LSMs, CBMs and CAVs are a method primarily used for deep learning models [49, 189] and have not been extensively applied to large-language models.

## 1.3 Theoretical Foundations of Representation Engineering

*1.3.1 Linear Representation Hypothesis.* The theory of why representation engineering can be used effectively is based on the Linear Representation Hypothesis (LRH). LRH postulates that high-level concepts and functions are encoded in the activations of neural networks as linear or near-linear features, identified as directions or subspace in the latent space of the model. [152, 164].

Representation engineering is built upon linear representation hypothesis, as it relies on the ability to isolate, manipulate, and interpret specific features within a model's latent representation space using linear methods. If common features like sentiment, gender, or style were not encoded in approximately linear subspaces, techniques such as vector arithmetic, projection, and basis decomposition would not work effectively. The success of interventions like linear probes [4], activation steering, bias removal, and feature control hinges on the assumption that these properties can be modified independently without complex non-linear entanglements. Linear representations in LLMs emerge naturally from the interaction of the next-token prediction objective and the implicit bias of gradient descent, rather than from architectural constraints [95]. If representations were highly non-linear or span multiple concepts, changing one attribute might unpredictably alter others, making controlled modifications impossible. Thus, the practical tools of representation engineering—feature extraction, interpretability, and controlled model editing—are fundamentally dependent on representations being structured in a way that allows for linear operations to meaningfully adjust model behavior.

The LRH is supported by extensive empirical findings across NLP and vision models, demonstrating that high-level features are often encoded as linear directions in activation space. Early evidence comes from word embeddings such as Word2Vec, where simple vector arithmetic (e.g., "king" - "man" + "woman" = "queen") suggested that semantic attributes are represented in a linear fashion [152]. This idea extends to modern transformer-based models, where probing studies have shown that features like syntax trees [84], part-of-speech tags [205], and named entities [127] can be recovered using linear classifiers with high accuracy. Large-language models have been shown to be effectively probed for space and time [74]. Similarly, truthfulness direction can also be found to significantly improve model performance on hallucination benchmarks [118]. In factual recall tasks, Meng et al. [151] demonstrated that specific residual stream activations encode factual associations in GPT models, with direct interventions shifting model predictions toward or away from a known fact. Burns et al. [30] found that unsupervised contrastive probes can identify "truth" directions in LLM activations, though generalization issues persist [115]. There is also evidence that individual concepts are in fact vectors [168]. In adversarial robustness literature, the presence of universal jailbreaks is highly predicated on the idea of LHR being true [154].

Bau et al. [10] showed that individual neurons and linear directions in CNNs correspond to human-interpretable object concepts, supporting the idea that representations align with linear subspaces rather than single neurons. These results collectively reinforce that while not all representations are perfectly linear [59], many high-level abstractions in deep networks behave as approximately linear features, forming the basis for practical representation engineering. If LRH is not true, representation engineering techniques like activation steering,

bias removal, and linear probes may be unreliable, as modifying one attribute could unintentionally alter others, limiting precise control over model behavior and leading to lack of generalizable hyperparameters that can be set for multiple problems.

Even though strong evidence supporting LHR being true exists, it is still a hypothesis and not a fully proven, absolute fact. In fact, probing shows that more complex internal representations are highly non-linear, encoding as detailed information as the board state during a game [117]. Linear representations may not capture the real structure but instead emerge during pre-training as bounds on the model's cognitive capacity [232]. In particular, much stronger support exists for the Weak LRH (assuming that some features are linear) than for Strong LRH (assuming that all features are linear).

Studies of the validity of LRH [165] provide a formal framework linking different interpretations of linear representation—subspace structure, probing, and intervention—demonstrating their fundamental connections. They introduce the concept of a causal inner product, refining how linear directions in representation space should be identified to ensure meaningful interventions. These findings support Weak LRH by confirming that many high-level features are encoded in linear subspaces, while also highlighting that naive geometric assumptions can lead to misleading conclusions. Other studies demonstrate that simple linear probes can reliably identify and manipulate a "truth direction," reinforcing the idea that at least some high-level concepts are structured linearly in model representations [99, 147]. This strengthens the theoretical basis of representation engineering, showing that while linear representations are exploitable, their structure may require more precise mathematical treatment to be reliably manipulated and, while performant in some aspects, cannot be reliably generalized for all cases.

While the empirical effectiveness of representation engineering is well-supported by evidence from other parts of the literature studying LRH, several studies show, that strong LRH may not be satisfied. Therefore, a generalized representation engineering intervention may not be possible and specific techniques need to be implemented to detect and mitigate nonlinear interactions.

*1.3.2 Superposition Hypothesis.* Superposition hypothesis states that the neural network represents more interpretable representations than there are dimensions in the representation space [59]. Instead of dedicating one axis of the representation space per feature, the network represents features as distinct directions in a high-dimensional activation vector. This hypothesis supports the empirical result of why neural representations are often distributed and overlapping [88]; the model effectively simulates a larger set of features than the number of neurons by superimposing features.

If true, the superposition hypothesis presents validation to the representation engineering approach, but also a challenge in practical implementation and development of general methods for representation engineering. Because of how neurons are encoded, a top-down approach allows to mitigate neuron polysemanticity and extract the net result of activations in the latent space. On the other hand, an intervention on a particular subset of the latent space may lead to changes in other representations, leading to decrease in capabilities and overall coherence. Furthermore, stacking several interventions may result in further unintended consequences, because one intervention may boost some effect of the other or nullify it.

## 2  REPRESENTATION READING

Representation Reading seeks to extract information about what activations correspond to a particular representation from the neural network, thus showing how LLMs represent and process information. As shown in Figure 2, a representation can correspond to multiple entities, specifically a concept, task, or function:

- A *concept* is a static property of the model like truthfulness, harmlessness, or morality.
- A *task* is a particular user query that for which a beneficial representation can be amplified to complete the task more effectively.
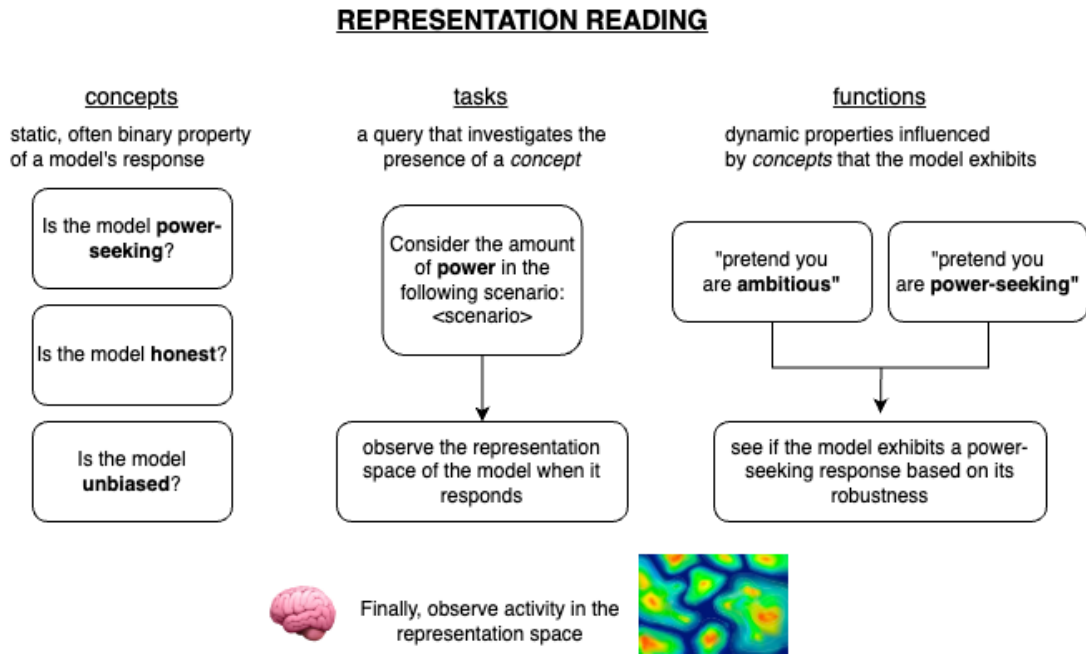
## REPRESENTATION READING



Fig. 2. Entities in representation reading: concepts, tasks, functions.

- A *function* is a dynamic property of the model like outputting correct Python code, power-seeking or answering a question in Spanish.

This section outlines methods to identify such representations. Causal manipulation experiments reveal that altering these abstract representations leads to predictable changes in model output, confirming that the network uses them for decision-making rather than passively encoding correlations [63]. Although representation engineering interventions require white box access, representation reading can be applied to black box models to predict overall model performance and detect harmful versions of LLMs [187].

### 2.1 Linear Artificial Tomography

The overarching idea of representation reading is to use differing inputs to stimulate the neural activity of the model, and to use the differences in observed activity to predict which activations correspond to model behavior. A foundational technique to detect such activations is Linear Artificial Tomography (LAT) [260]. LAT reading requires defining a stimulus and a task for detecting the neural activity corresponding to it. The stimulus can be as simple as the following prompt (with "___" added to show where output would be generated):

> Consider the amount of <CONCEPT> in the following: <STIMULUS>. The amount of <CONCEPT> is ___

The neural activity of the network when stimulated with this model is then compared with the activity of the model when given either a) a contrasting stimuli or just b) a reference stimuli with no specific concept or function. Based on this data, a linear model is used to identify a direction to detect the direction of the identified representation. These models can either be supervised or unsupervised. In practical applications, either a supervised linear probe or an unsupervised Principal Component Analysis (PCA) is applied.
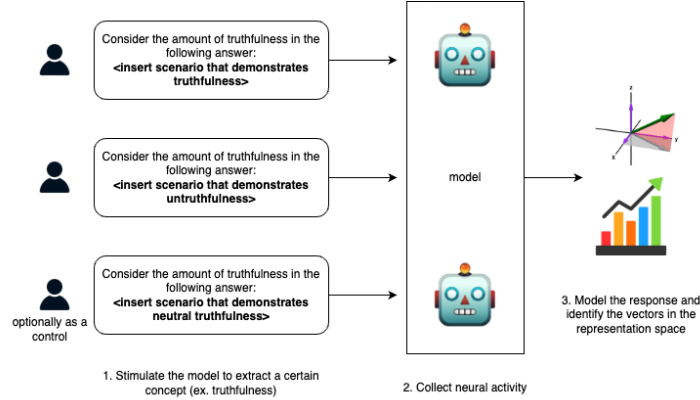
Fig. 3. A high-level example of a LAT experiment.

For a concept $c$, given a model $M$, a function Rep that extracts representations, and a stimulus set $S$, the neural activity is collected as:

$$A_c = \text{Rep}(M, T_c(s_i))[-1] \mid s_i \in S \tag{1}$$

where $T_c(s_i)$ represents the task template applied to stimulus $s_i$. To determine the principal representation direction, PCA is applied to the difference vectors of neural activations:

$$v = \text{PCA}(A_c^{(i)} - A_c^{(j)}) \tag{2}$$

where $A_c^{(i)}$ and $A_c^{(j)}$ are pairs of neural activities from different stimuli. To manipulate the model's representation, a projection operation is performed:

$$R' = R - \frac{R^T v}{|v|^2} v \tag{3}$$

where $R$ is the original representation and $R'$ is the transformed representation after adjusting along the reading vector $v$. For more details, see Zou et al. [260].

## 2.2 Probing

Probing aims to relate specific features to activation patterns in neural networks by training supervised models to map activations to target variables [93] and eventually representations. Representations can be detected by examining the emergence of abstract, low-dimensional manifolds that encode semantic features shared by different inputs [63]. The complexity of the probing classifier influences the results, with simpler probes often providing more interpretable insights. More complex probes may achieve higher accuracy but risk inferring features not actually used by the network. Causal analysis techniques involve interventions in the representations to assess the impact on the model's original performance, revealing whether certain features are genuinely utilized. The choice of datasets for both the original model and probing tasks significantly affects the outcomes, making it critical to choose the right dataset to probe on [110, 115].

Figure 4 shows many example goals that can be achieved through representation reading. More is discussed in Section 2.2.1. The figure also shows specific techniques that are covered in Section 2.2.2 and Section 2.2.3.

*2.2.1 Applications of Probing.* Probes have been applied to investigate representations such as LLM biases [203], jailbreak formation mechanisms [78] and even token embeddings of senses like sounds [156], shapes and color [22]. Probes have been shown to encode uncertainty [212], hierarchical linguistic constituency [133],

Fig. 4. Graph showing Representation Reading, categorized by applications and type of representation.

cross-lingual retrieval [72], Abstract Syntax Trees (ASTs) [82], hierarchical syntax [6], model knowledge [12] or characters in a token [102]. Probing can even be applied in biological neural networks [93] and black-box models [186].

Probes have been applied to check whether a model is acting against its maker's wishes [130], showing that models intentionally generate falsehoods, explaining hallucination and model misalignment via distinct prediction mechanisms. Several studies find effective solutions to detect truthfulness representations through

different methods of probing [7, 34, 37, 115, 238]. Probes have also shown to be effective at identifying LLM performance across written languages using cosine similarity [125], revealing how activations for high-resource and low-resource languages differ. Through probing, it has been found that as language models scale, their vector spaces increasingly resemble the rich vector spaces of vision models [116]. Probes have shown to be effective at detecting representations the model uses for distinguishing between its own knowledge and actions and those of others, supporting nuanced social reasoning tasks [257]. Probes have also been efficently applied to investigate societal biases [137]. Despite Reinforcement Learning from Human Feedback (RLHF), the model retains gender bias [108] while outright refusing to engage with race and religion related prompts. Manipulating internal belief representations leads to consistent changes in Theory of Mind performance, confirming that the model actively uses them for decision-making rather than merely storing passive correlations [257].

2.2.2 *Linear Probes.* Probing is often implemented through linear probes, simple linear classifiers trained on frozen features extracted from specific layers of a pre-trained neural network at inference time. By taking activations from a given layer as input and training to predict a specific property, such as part-of-speech tagging or sentiment analysis, linear probes show how well different layers capture certain types of information [4, 99]. Given a hidden representation $H_k$ at layer $k$, a linear probe applies a weight matrix $W$ and bias $b$ followed by a softmax function:

$$f_k(H_k) = \text{softmax}(WH_k + b)$$

where the softmax function is defined as:

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^{D} e^{z_j}}$$

This transformation maps the hidden representation into a probability distribution over $D$ classes. The parameters of the linear probe $W$ and $b$ are trained using the cross-entropy loss:

$$L = -\sum_{i=1}^{D} y_i \log(\hat{y}_i)$$

where $y_i$ represents the true class label, and $\hat{y}_i$ is the predicted probability for class $i$. The objective is to minimize $L$, thereby aligning the learned representations with the desired classification task.

Linear probes allow us to localise where specific types of information are encoded. These abstractions can be identified by observing the convergence of embeddings for tokens with similar semantic roles, even when they are perceptually distinct [63]. For example, truthfulness information tends to concentrate within "exact answer tokens" [160]. Error detection methods can be significantly enhanced by focusing on these tokens. Intermediate representations of LLMs can be used to predict the types of errors they might make [99].

Other methods add new components to this setup for more precise representation identification. Representational Similarity Analysis (RSA) uses pairwise cosine similarities between embeddings to reveal latent structure in task representations, enabling classification of task components without parameterized probes [22, 240]. Contrast-Consistent Search (CCS) identifies truth representations by probing embeddings associated with the final token of a statement, aiming to extract subjective probabilities of truthfulness using supervised learning on labeled datasets [115].

2.2.3 *Other Types of Probes.* Logic Tensor Probes (LTP) train two-layer neural models to probe for specific predicates, treating each predicate's embedding method as a hyperparameter and optimizing probe depth for improved interpretability [124]. LTPs are implemented by repurposing Logic Tensor Networks as shallow

classifiers that assess logical consistency in frozen language model representations by encoding first-order logic constraints [145].

Probing Classifiers utilize transformer hidden states to predict hallucinations at token and response levels, leveraging layer-specific features for improved classification accuracy [34]. Probing Classifiers in Bronzini et al. [27] extract factual knowledge from LLM hidden states via activation patching, using classifiers to map token activations to logical predicates for claim verification.

INSIDE detects hallucinations by measuring semantic consistency in sentence embeddings using EigenScore, a covariance matrix determinant-based metric that captures dense semantic information [37]. SAPLMA trains a classifier on LLM hidden states to predict the truthfulness of statements, focusing on out-of-distribution generalization and improving over probability-based methods [7]. Pathoscopes uses the models' own generation capabilities to explain their internal computations in natural language [66]. Tuned lens target affine transformations that align intermediate hidden states with final layer logits [13]. Intrinsic Dimensions correlates the geometric structure of token embeddings with classification loss, showing that tokens in high-loss prompts exist in higher-dimensional spaces [212].

## 3 REPRESENTATION CONTROL

Representation Control refers to the process of modifying or guiding the internal representations of concepts and functions identified in the Representation Reading step towards a particular outcome. Representation Control interventions primarily take the form of an insertion of a vector (or a set of vectors) between the layers of the model at inference without retraining the model. See Figure 5 for an example.



Fig. 5. Explanatory effect of an intervention on next token probabilities.

### 3.1 Intervention Methods

*3.1.1 Single and Multi Property.* The representations identified by representation reading techniques can be applied to steer the model into a particular direction. However, steering the model in multiple directions proves difficult and not all methods support applying multiple directions at the same time. For example, steering the model to be both more trustworthy and less prone to jailbreaks requires support for steering the model in multiple directions at the same time.

Theoretically, a representation can be reframed as a more general one, but identifying such representation is more difficult and leads to unpredictible behaviour, usually leading to loss of fluency and generation of nonsensical tokens. Experiments show combining more general representations are usually less effective than combining several injections at different layers [210]. Because of this, we define a separation between activation

Fig. 6. Graph showing Representation Control, with goals that are implemented using specific methods.

steering techniques that allow for steering the model in a single direction [209] versus multiple directions [32, 188, 198, 210].

Combining steering vectors for different concepts can lead to unpredictable interactions, where the effect of one intervention interferes with another [225, 229]. Moreover, simultaneous injection of individual steering vectors at different places in the model appears to be more effective than combined steering [210]. Activation steering can simultaneously enforce multiple instruction constraints, such as requiring an LLM to generate JSON-formatted responses while also adhering to word exclusion rules, by applying composite activation vectors during inference [198]. Alternatively, Sparse Activation Control can identify distinct attention heads and hence separate interventions over multiple properties [229] or conceptor matrices can be used to eliminate the negative

interactions between distinct properties interacting with each other [169]. The implementation of representation control can involve modifying a single layer [94, 161, 209] or multiple layers [11, 118].

*3.1.2 Constant or Dynamic.* Once a representation has been identified and the vector (or multiple vectors) has been created, it needs to be applied to the model to modify its activation space. In the process, either the representation engineering developer needs to assign intervention strength, choosing a **Constant** [161, 209, 260]; or it is set as **Dynamic**, i.e. determined by the model at inference time [188, 197, 206, 219]. There are several approaches to setting those weights dynamically, either based on a probe and unsupervised clustering [219], Kullback-Leibler (KL) divergence [188], cosine similarity of activation matrices between unsteered and steered models [188] or gradient-based optimization [197].

*3.1.3 Intervention Stage.* Representation engineering phases can be categorized based on where the modifications are applied within the model. These can be performed on all activation spaces and components [27, 32, 79, 97, 111, 118, 167, 188, 196, 218–220, 257]. Alternatively, the intervention may target the residual stream activations, that is, the hidden states that pass between layers through residual connections [23, 36, 51, 61, 80, 87, 90, 121, 139, 161, 162, 166, 169, 194, 199, 209, 234, 258, 260]. Another approach involves modifications to the target MLP layers [33, 134, 140, 174, 199, 227, 254, 259, 260], specific neurons [206, 229, 247], or changes to the model through fine-tuning, re-training, or other forms of weight modification [38, 41, 60, 64, 86, 107, 249].

## 3.2 Intervention Goals

Representation control methods aim to steer the model towards particular outcomes. It is important to evaluate each method through the lens of its original goal, because identification and intervention on one representation does not mean the intervention is going to be successful on another representation type. The precise nature of representations is difficult to define. A representation can range from broad traits such as honesty to specific facts encoded in activation patterns [83]. However, extracting such representations is challenging because probing methods often fail to generalize, capturing only sample-specific features rather than the true underlying structures [115]. Moreover, attempts to identify representations may inadvertently conflate them with tasks the model is solving rather than isolating the true concept encodings [110]. This difficulty arises because internal structures of LLMs do not necessarily align with human-meaningful categories, making interpretation reliant on indirect techniques with inherent limitations. Defining the sample and representation scope is therefore critical, and the interventions remain specific to a particular representation. Simple concepts like truthfulness are easier to identify than sophisticated ones like humor or appropriateness [214].

*Personalization.* Customizing the values and response style of an LLM is important for increasing adoption, general satisfaction and capability at completing tasks [250]. Activation steering is a lightweight method to steer intended behavior without the substantial computational resources of fine-tuning, hence allowing for the creation of individualized LLM experiences [32]. Identifying and applying representations of particular style is challenging, but possible through activation steering [32, 36, 97, 134, 206, 227, 260]. Embedding certain moral values into the model is also possible.

*Security.* LLMs are aligned to reject harmful requests, but adversaries may nonetheless manipulate them into producing harmful outputs, for example through jailbreaks [18]. LLMs are capable of being manipulated into assisting in unethical requests, leaking personal data or giving toxic responses [20, 41, 55, 60, 64, 107, 111, 128, 140, 167, 188, 194, 196, 209, 218, 237, 249, 259]. In particular, circuit breakers have proven to be particularly effective, relying on identification of harmful outputs and preventing them from generating a response [259]. Representation engineering can be used to create defenses, but also circumvent existing safety measures for red-teaming purposes. Activation steering can extract unlearned information from large-language models (LLMs), demonstrating

vulnerabilities in current unlearning techniques [191]. Malicious attack vectors can be inserted at inference time to steer the model towards harmful behaviour and manipulate model alignment on the representation level, as described in the Trojan Activation Attack [216], future backdoor trigger method [171], the Safety Concept Activation Vector (SCAV) [231] or direct interventions in the model's representation space to reduce refusal rates [259]. Similarly, Li et al. [119] detect safety patterns in the latent space, which when mitigated indicate that models are less able to refuse harmful requests and strengthening safety patterns indicate that models are more resistant to jailbreaking attempts. One can also create representation-level roleplay jailbreaks [67].

Security risks exist not only from the adversaries, but also the model itself. As LLMs become more capable, it becomes difficult to know when a model is pursuing goals that are aligned with its human evaluators, or engaging in strategic deception and faking alignment [71]. Early experiments suggest that linear probes alone are not sufficient to detect such responses.[69].

*Performance.* Modern AI training is optimized for enabling the most effective learning algorithm to generate accurate responses to variety of general questions. Interventions in the latent space can steer the model towards better responses. Representations of particular skills, like Chain-of-Thought thinking [246], code type prediction [138] or knowledge selection accuracy [254]. Several approaches have been applied to selectively steer the model towards more performant responses at particular tasks [27, 33, 80, 134, 169, 174, 188, 197, 227, 254, 260]. In particular, task vectors, a small scale intervention stemming from In-Context-Learning has been successfully applied to steering the model to be better at narrow tasks.

*Truthfulness.* LLM hallucination refers to models generating false outputs with high confidence [185]. They are a problem because they can lead to malfunctioning in critical use cases, misinformation and loss of trust in AI [89, 248]. Improving truthfulness through representation engineering is a major area of research [32, 118, 161, 167, 196, 219, 219]. By identifying and amplifing the "truth" direction, activation steering methods have shown to reduce hallucination rates. Examples include HPR [167], CAA [161] or ACT [219]. This is primarily achieved by identifying representations from a part of TruthfulQA, a commonly used hallucination benchmark and testing on a held-out set of this benchmark [167] [161].

## 4 TAXONOMY

We structure all identified representation engineering methods in a comprehensive taxonomy. We provide a relative ranking of these methods based on their performance. For performance, we evaluate whether these methods have been proven by subsequent research to be less effective, whether issues lead to fluency degradation, and how comprehensive the original testing across multiple representations has been.

### 4.1 Linear Contrastive Steering Vector

Linear contrastive steering vectors provide a framework for modifying transformer model behavior through targeted interventions in the hidden activation space. These methods leverage differences in activation patterns—obtained by contrasting selected prompt pairs or decomposing internal representations—to derive vectors that steer model outputs toward desired behaviors.

*4.1.1 Simple Contrastive Vectors.* The contrastive steering vector is computed as the difference between representations from a target and reference scenario:

$$v_c = \text{Rep}(M, T^+) - \text{Rep}(M, T^-) \tag{4}$$

where $T^+$ and $T^-$ correspond to positive and negative stimuli for the concept under study and $M$ is the model. The intervention is applied by modifying the original representation through a controlled addition or subtraction

| Technique Name | Intervention Stage | Goal | Performance |
|---|---|---|---|
| Latent Steering Vectors [199] | [R][T] | Personalization | ◑ |
| RepE [260] | [R] [T] | Personalization, Performance, Truthfulness | ◑ |
| CAA [161] | [R] | Truthfulness | ◑ |
| ITI [118, 257] | [A] | Truthfulness | ◑ |
| ActAdd [209] | [R] | Security, Personalization | ◑ |
| C-WSV [20] | [AA] | Security | ◑ |
| Bi-directional Pref. Optim. [32] | [A] | Personalization | ● |
| PaCE [140] | [T] | Security | ● |
| Mean-Centered Activ. Steer. [97] | [A] | Personalization, Security | ◕ |
| Household. Pseudo-Rotation [167] | [A] | Security, Truthfulness | ● |
| MiMiC [194] | [R] | Security | ◑ |
| SARA [206] | [N] | Personalization | ◑ |
| CAST [111] | [A] | Security | ◑ |
| SCANS [33] | [T] | Performance | ◑ |
| Task Vectors [23, 51, 80, 87, 90, 121, 139, 162, 166, 234, 258] | [R] | Performance | ◑ |
| Activation Patching [27, 79] | [A] | Security, Performance, Truthfulness | ● |
| GRATH [38] | [M] | Truthfulness | ● |
| Entropic Activation Steering [174] | [T] | Performance | ◑ |
| SAC [229] | [AA][N] | Truthfulness | ● |
| Activation Scaling [197] | [AA][R][T] | Performance | ◕ |
| SADI [220] | [A][AA][N] | Performance, Truthfulness | ◕ |
| DAC [188] | [A] | Security, Performance | ● |
| RE-CONTROL [107] | [R][N][T] | Security | ◕ |
| ACT [219] | [A] | Truthfulness | ◑ |
| KL-Then-Steer (KTS) [196] | [A] | Security, Performance | ◑ |
| InferAligner [218] | [A] | Security | ◑ |
| CONTTRANS [55] | [R] | Security | ◑ |
| KTCR [64] | [M] | Security | ◑ |
| TruthX [247] | [AA][N] | Truthfulness | ◑ |
| Safety-Aware Fine-tuning [41] | [M] | Security | ◑ |
| ReFT [227] | [T] | Personalization, Performance | ● |
| RAHF [134] | [T] | Personalization, Performance, Truthfulness | ● |
| LoRRA [86] | [M][AA] | Personalization, Performance, Truthfulness | ◑ |
| SPARE [254] | [T] | Performance | ◑ |
| SAE-TS [36] | [R] | Personalization | ◑ |
| Self-knowledge SA [61] | [R] | Truthfulness | ◑ |
| ARE [249] | [M] | Security, Truthfulness | ● |
| Conceptors [169] | [R] | Performance | ● |
| Circuit Breakers [259] | [T] | Security | ● |
| LEGEND [60] | [M] | Security | ◔ |

| | | | |
|---|---|---|---|
| 🟦 Linear Fixed Methods | 🟪 Fine-tuning Methods | [A]: All activation spaces | [T]: Target MLP layer activations |
| 🟦 Multiple Model Methods | 🟨 Sparse Auto-Encoders | [AA]: Attention activations | [N]: Target neurons (+/- activations) |
| 🟥 Dynamic Strength Methods | ⬜ Other | [R]: Residual stream activations | [M]: Modifies model(s) (fine-tuning, re-training, etc.) |

Table 1. Taxonomy of Representation Engineering Methods.

of the contrastive vector:

$$R' = R + \alpha v_c \tag{5}$$

where $\alpha$ is a scaling coefficient that adjusts the strength of the intervention. This common framework has been developed in parallel, spanning common methods across RepE [260], CAA [161], ITI [118] and ActAdd [209].

*CAA (Mean Difference).* Contrastive Activation Addition (CAA), modifies activations in the residual stream [161]. CAA constructs steering vectors by averaging the difference in activations between contrastive prompt pairs, such as factual versus hallucinatory responses, at a specific transformer layer. These vectors are then applied at inference time to modulate model outputs, either reinforcing or suppressing target behaviors like sycophancy, corrigibility, or refusal. CAA outperforms system-prompting and shows complementary effects with finetuning. CAA vectors exhibit layer-wise consistency, transferability between base and fine-tuned models, and alignment with behaviorally meaningful activation clusters, improving their interpretability.

$$v_c = \frac{1}{N} \sum_{i=1}^{N} \left( \text{Rep}(M, T_i^+) - \text{Rep}(M, T_i^-) \right) \tag{6}$$

*RepE (PCA-Based).* Representation Engineering (RepE) transforms internal representations using vectors derived from Linear Artificial Tomography (LAT) [260]. It identifies directions in representation space that correlate with cognitive functions, then intervenes by linearly combining these reading vectors with the model's activations. This is achieved through two methods: (1) Reading Vector Intervention, where a fixed direction is added or subtracted from activations to amplify or suppress the associated cognitive function, and (2) Contrast Vector Intervention, where contrast vectors are calculated by running paired contrastive prompts through the model and subtracting their activations.

$$v_c = \text{TopPC} \left( \{ \text{Rep}(M, T_i^+) - \text{Rep}(M, T_i^-) \}_{i=1}^{N} \right) \tag{7}$$

*ITI (Classifier-Based).* Inference-Time Intervention (ITI) enhances the truthfulness of language models by shifting activations in specific directions during inference [118, 257]. It identifies attention heads with high linear probing accuracy for truthfulness and intervenes only in these heads by shifting activations along directions correlated with truthful outputs. This intervention is applied autoregressively throughout the text generation process. ITI utilizes two main parameters: the number of attention heads to adjust and the magnitude of the shift, enabling fine-grained control over the intervention. Unlike weight editing methods, ITI operates directly on attention head activations, preserving the underlying model weights and reducing computational costs.

$$v_c = w^*, \quad A_h' = A_h + \alpha v_c \tag{8}$$

*ActAdd (Single Pair).* Activation Addition (ActAdd) computes a steering vector by contrasting activations from prompt pairs, such as "Love" vs. "Hate," and applies this vector during inference to modify model behavior [209]. Most importantly, it can achieve good results with even one contrastive pair, and the effects use a simple subtraction of contrasting activations to create the steering vector.

$$v_c = \text{Rep}(M, T^+) - \text{Rep}(M, T^-) \tag{9}$$

Studies examining these have found CAA to be most reliable and also formulated theory as to why this is the optimal method [91].

*LSV.* Latent Steering Vectors (LSVs) intervene in pretrained language models by adding fixed-length vectors, $z_{\text{steer}}$, directly to the hidden states during decoding, influencing model outputs without altering model weights [199]. Specifically, $z_{\text{steer}}$ is optimized via gradient descent to maximize the log probability of a target sentence, keeping the language model frozen.

$$z_{\text{steer}} = \arg \max_z \sum_{t=1}^{T} \log p(x_t | x_{<t}, z_{\text{steer}})$$

LSVs can be injected at different layers, including embedding, self-attention, feed-forward, and the final language model head, but are most effective when added to the middle layers (e.g., layers 6 or 7 in GPT-2). These vectors are injected either at every timestep or just the first timestep, with minimal performance loss in the latter case. By linearly combining $z_{\text{steer}}$ with the model's hidden states, LSVs steer the model to generate desired outputs with high accuracy. This method preserves the model's overall behavior since it operates at the activation level rather than modifying weights or retraining.

$$h' = h + \alpha \cdot z_{\text{steer}}$$

Latent steering vectors can be extracted directly from pretrained language models without requiring fine-tuning [199]. These vectors encode control mechanisms within the model's activations, allowing direct manipulation of text generation. By injecting steering vectors into hidden states, models can achieve near-perfect target sentence recovery and enable unsupervised sentiment transfer on par with specialized models. Notably, steering vectors preserve sentence similarity relationships and outperform traditional hidden-state pooling methods in semantic textual similarity benchmarks. Their effectiveness depends on the injection location within the transformer stack, with middle layers providing optimal results—consistent with findings in adapter-based fine-tuning. Importantly, these vectors do not merely memorize sequences but encode generalizable latent structures, reinforcing their potential for controlled text generation and precise representation manipulation.

*C-W SV.* Category-wise Steering Vectors [20] compute steering vectors by computing activation differences of harmful and harmless text and applying these vectors at specific model layers to redirect outputs toward safer regions in latent space. This approach includes both unsupervised and guided methods for vector extraction, with the latter incorporating external safety classifiers to refine steering signals.

*4.1.2 Optimized Steering Vectors.* Recent literature aims to optimize the simple steering vectors, increasing their effectiveness for particular use cases. This section presents the methods with the highest empirical validation to show alternatives to simple contrastive vector interventions.

*BiPO.* Bi-directional Preference Optimization (BiPO) optimizes steering vectors by adjusting the generation probability of human preference data pairs instead of relying on direct activation differences from contrastive prompts [32]. BiPO has shown to be particularly effective at steering AI personas compared to other representation engineering techniques. BiPO scored higher than Contrastive Activation Addition and Freeform on every personalization task it was evaluated on. This approach also improves the model performance on tasks such as truthfulness, hallucination mitigation, and jailbreaking defense. The vectors calculated using this approach transfer across models and fine-tuned LoRAs. BiPO allows for the application of multiple steering vectors to influence multiple behaviors simultaneously without decreasing the general capabilities of the model.

*PaCE.* Parsimonious Concept Engineering (PaCE) constructs a large-scale concept dictionary in the activation space and uses sparse coding techniques to decompose model activations into benign and undesirable components [140]. By using oblique projection and a sparse linear combination of concept directions, a new vector is

constructed embedding all beneficial concepts in it. PaCE is shown to maintain fluency while preserving safe intervention effect.

*Mean-centering.* Mean-centering improves activation steering by refining steering vectors through a subtraction operation: the mean activation of a training dataset is subtracted from the mean activation of a target dataset [97]. This removes global biases in activations and improves control over model outputs. Applying mean-centering to toxicity reduction in language models results in lower toxicity scores compared to existing activation steering methods. Mean-centering increases the relevance of genre-specific terms when steering a model to continue a story in a target genre. It also improves the accuracy of function vectors, such as those used to extract country-capital relationships [97].

*HPR.* Householder Pseudo-Rotation (HPR) treats activations as having both direction and magnitude, ensuring norm preservation during edits [167]. This is fundamentally different from interventions that simply perform addition [161]. HPR first learns a hyperplane separating desirable and undesirable activations, then reflects negative activations across this plane using a Householder transformation, followed by a rotation to align them with positive behaviors. Experiments show that HPR significantly outperforms other methods in hallucination reduction. This approach significantly improves over the previous linear additive steering vectors, also by proving the evidence of the Magnitude Consistency Property.

The magnitude consistency property means that the magnitude of a representation vector remains relatively stable when the underlying semantic meaning of the input does not change fundamentally [167]. This property helps to *a)* detect meaningful semantic representations by observing how vector magnitudes behave under different transformations , *b)* identify when neural networks are developing robust, consistent internal representations, and *c)* understand potential failure modes where representations might become unstable or distorted. For example, if you input two semantically similar sentences into a language model, a magnitude-consistent representation would show similar vector magnitudes, indicating the model has captured a stable semantic understanding rather than arbitrary mappings. This property highlights a key limitation of the point-in-space view, where the steering vector approach cannot simultaneously maintain activation magnitude consistency and effectively edit the activation to achieve greater performance improvement for desired behaviors for language models [165, 167].

*MiMiC.* Minimally Modified Counterfactuals (MiMiC) is an affine steering method that shifts language model representations toward a target concept by modifying their mean and covariance [194]. MiMiC achieves this by applying a linear transformation that aligns the concept-conditional distributions, ensuring that representations from different groups become statistically similar. Unlike simple mean-matching approaches, MiMiC also mitigates bias by neighbors, meaning that the spatial clustering of representations no longer correlates with protected attributes like gender or dialect.

*SARA.* Similarity-based Activation Steering with Repulsion and Attraction (SARA) is a causal intervention technique that modifies the internal activations of a language model to align with specific ethical perspectives while suppressing undesired ones [206]. It operates by computing cosine similarity between activation matrices of different prompts, selectively reinforcing target activations and repelling conflicting ones through singular value decomposition (SVD). Unlike standard activation steering, which applies uniform shifts, SARA dynamically adjusts activations based on their initial similarity to the target, ensuring more precise and context-aware modifications.

*CAST.* Conditional Activation Steering (CAST) introduces condition vectors that selectively determine when steering happens, allowing models to refuse specific types of inputs without affecting unrelated responses [111]. CAST relies on cosine similarity between a model's hidden state and predefined condition vectors to trigger interventions. This means that refusal behavior is applied only in relevant contexts. Experiments show it results
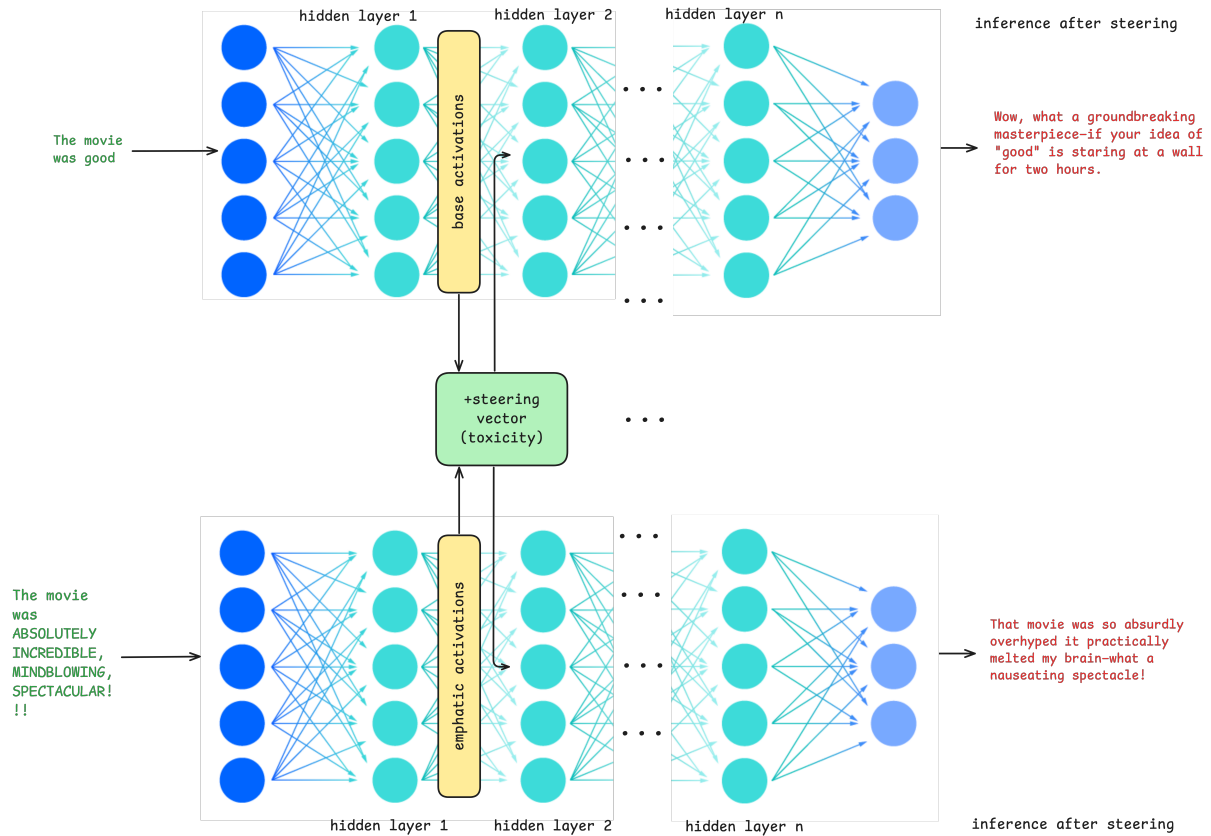
Fig. 7. Violation of the magnitude consistency property (high-level view). In order to uphold the magnitude consistency property, a robust LLM should have the same (or nearly equal) activation vector norms for semantically similar inputs. Steering vector intervention disrupts that by causing disproportionate scaling of norms and context collapse. In this example, the activation vector norm of the base input and the emphatic input will not remain the same after a steering vector that encompasses toxicity is applied.

in better security and performance than simple activation steering. The method also enables logical composition of condition vectors, such as enforcing refusals for legal or medical inquiries while permitting general responses.

*Probe-Free Low-Rank Activation Intervention (FLORAIN).* FLORAIN [94] introduces a probe-free, low-rank activation intervention to modify language model activations at inference time, improving the generation of desirable content. FLORAIN directly learns a low-rank transformation that projects activations onto a desirable manifold modeled as an ellipsoid region. This approach ensures efficient inference-time interventions, reduces computational overhead by operating on a single layer, and minimizes representation shifts across layers. To address the challenge of high-dimensional activation vectors with limited training samples, FLORAIN employs an extrapolation strategy that shifts mean activation vector away from undesirable activations by combining a question-specific adjustment (capturing the difference between desirable and undesirable activations for a given question) and a global adjustment (capturing the overall shift from undesirable to desirable activations across all data). A weighting factor balances the influence of local versus global information, ensuring robustness

even in low-data scenarios. This refined mean activation vector defines the ellipsoid-shaped desirable activation region onto which activations are projected during inference, enhancing truthfulness and coherence in generated outputs.

*Aggregate and Conquer: Detecting and Steering LLM Concepts.* This method [11] detects and steers LLM concepts by leveraging Recursive Feature Machines (RFMs) to extract eigenvector-based concept representations from activations. For detection, it trains layer-wise predictors to classify activations, then computes the Average Gradient Outer Product (AGOP) matrix to identify the most influential directions, selecting its top eigenvectors as concept vectors and aggregating them across layers for improved accuracy. For steering, it modifies activations during the forward pass by replacing each layer's activation with the learned concept vector and a control coefficient that varies per concept and model. To steer multiple concepts simultaneously, it combines layer-wise concept vectors in a weighted sum, allowing fine-grained control over behaviors such as reducing toxicity, enhancing truthfulness, altering writing styles, or modifying sentiment strength. Unlike previous methods, this approach learns concept vectors from nonlinear predictors, enables multi-layer and multi-concept steering, and supports gradated control, making it a powerful and efficient alternative to fine-tuning or contrastive interventions.

*4.1.3 In-Context Learning and Task Vectors.* In-context learning (ICL) enables large-language models to adapt dynamically to new tasks by leveraging additional examples to modify internal representations without parameter updates by providing few-shot examples of answers to the task [28, 54, 121]. Attention mechanisms in transformers can be viewed as performing a form of implicit gradient descent, effectively allowing the model to optimize for new tasks on the fly [46, 213]. This can be formalized using PAC learning frameworks to establish finite sample complexity bounds [224]. ICL can be viewed as an implicit learning algorithm, where transformers encode smaller models within their activations and adjust them as new examples are introduced [3, 163]. ICL can also improve model performance through self-verification [221]. The effectiveness of ICL can be explained by implicit Bayesian inference, where the model infers a latent structure that connects the pretraining and inference distributions [230].

ICL utilises specialized attention heads (induction heads) that detect and replicate patterns in token sequences, effectively driving the model's ability to generalize from context [159]. A subset of induction heads, semantic induction heads encode structured relationships such as syntactic dependencies and knowledge graph associations showing how additional examples at inference change representation structure [179]. ICL can successfully approximate complex function classes, including linear models, decision trees, and even neural networks, showing that transformers can implement efficient learning algorithms internally [65, 132]. ICL alters embeddings and attention distributions to prioritize task-relevant information while reducing sensitivity to adversarial distractions, also in the case where the provided examples are not particularly useful for task performance [240]. Some studies show that example-label correctness has minimal impact while others find significant sensitivity to correct label assignments, showing that the effectiveness of ICL depends on dataset structure and model scale [153, 239]. It is the structure of demonstrations that is important, as models can leverage input distribution patterns and sequence formats even when labels are randomized, showing that task representation might matter more than exact label correctness [153, 180]. Increasing the number of examples leads to significant performance gains across diverse tasks [1]. While ICL consistently improves task performance over zero-shot prompting, its effectiveness is highly dependent on the selection and ordering of examples, as variations in these factors can cause substantial fluctuations in model accuracy [129, 255]. In-context learning can help the model identify the right parts of the input space and the overall format of the sequence [153]. ICL is similar to RepE in that it alters the representations. However, it does not monitor the full effect providing extra tokens in the input has throughout the layers.

ICL can be used to crate a task vectors are compressed versions of contextual examples that condition model behaviour with a latent representation [80, 121]. Task vectors can be thought of as compact, low-level
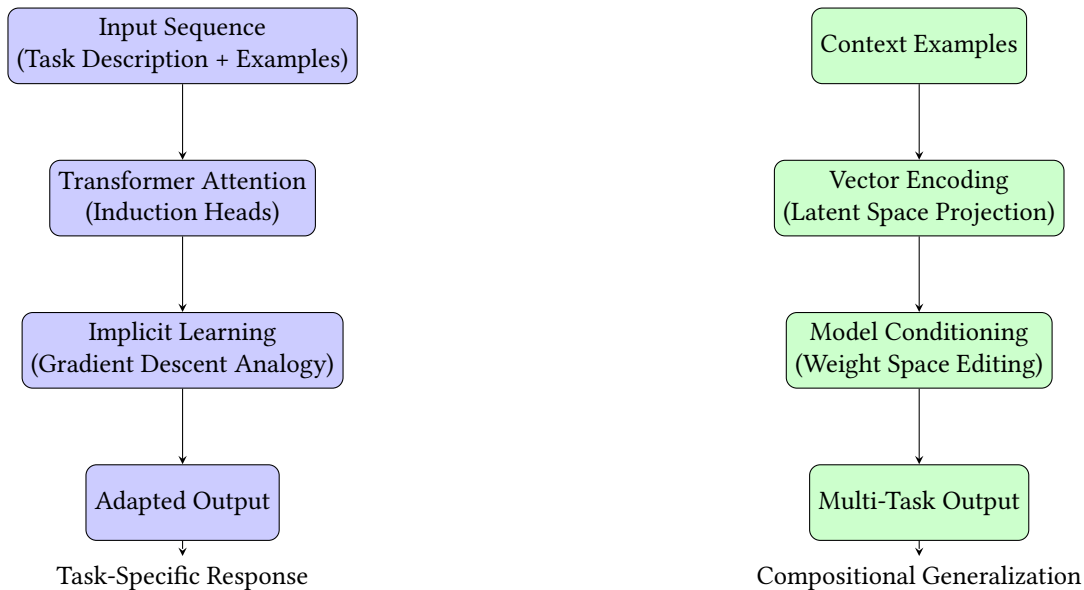
Fig. 8. In-Context Learning (left) modifies model behavior through attention mechanisms and demonstration examples, while task vectors (right) create persistent latent representations that enable direct model conditioning.

representation engineering intervention. High-dimensional task embeddings align with how humans mentally structure and manipulate knowledge [155, 168, 192]. Task vectors naturally emerge in transformer-based models during in-context learning, where they encode latent task structures within activation spaces, allowing models to generalize efficiently across similar input contexts without explicit weight updates [80, 234, 258]. Task vectors encode task-specific transformations in model weight space, enabling multi-task learning and transferability across different domains [23, 90, 234]. Task vectors allow models to create, subtract, and merge different tasks. Their robustness can be further improved through selective importance weighting and structured sparsification techniques [23, 162, 166]. Task vectors have been also been applied in multi-modal and robotic control settings, where they enable compositional plan representations that facilitate generalization to novel tasks and modalities [51, 87, 139].

*4.1.4 Activation Patching.* Activation patching (also known as causal tracing) is a method at the intersection of mechanistic interpretability and representation engineering. It seeks to identify small-scale model components responsible for specific behaviors by changing token representations during inference[27, 79]. Activation patching replaces activations with those from a different inference run. Through this, it traces how specific activations contribute to model behavior. For example, if a model completes *"The Eiffel Tower is in"* with *"Paris,"* activations from an inference with a corrupted prompt (e.g., *"The Colosseum is in"* for "Paris") to observe the impact on output [244]. The latent representation of this placeholder is replaced with a weighted sum of the input's latent representations from the original inference, effectively "patching" the activation [27, 235]. Noising and denoising are two primary approaches to activation patching: denoising patches clean activations into a corrupted run to identify sufficient components, while noising patches corrupted activations into a clean run to determine necessary components [79]. Different corruption methods—such as Gaussian noising and symmetric token replacement—can be implemented to look at changes in interpretability [244].

*4.1.5 SCANS.* SCANS (Safety-Conscious Activation Steering) extracts "refusal steering vectors" from the model's hidden states by contrasting activations from harmful and benign queries, identifying specific layers responsible for refusal behavior [33]. This is used to prevent the model refusing to follow benign queries from the users. These vectors are then used to adjust activations during inference—steering responses away from unnecessary refusals. SCANS achieves this by first classifying a given query as harmful or benign using a similarity-based method that compares hidden state transitions with a learned reference direction. If a query is deemed benign, SCANS modifies the activations in safety-critical layers to reduce the likelihood of refusal without altering the model's core capabilities.

*4.1.6 GRATH.* GRATH generates and refines truthfulness training data using out-of-domain questions, then iteratively fine-tunes the model through Direct Preference Optimization (DPO) to improve its alignment with factual correctness [38]. What the model essentially does is self-generate pairs of information it then self-truthifies itself on. GRATH significantly improves model accuracy on the TruthfulQA benchmark. GRATH is more cost-effective and scalable than human-annotated datasets for improving the reliability of its outputs.

*4.1.7 EAST.* AI agents exhibit overconfidence in decision-making, often committing prematurely to actions without sufficient exploration [174]. Entropic Activation Steering (EAST) directly modifies an LLM's internal uncertainty representation, increasing exploration in sequential decision-making tasks. Unlike adjusting sampling temperature, which minimally affects action entropy, EAST constructs a steering vector by averaging activations weighted by action entropy and applies it at inference time to modulate decision confidence. Experiments in bandit environments demonstrate that EAST significantly increases exploration, prevents premature commitment to suboptimal choices, and generalizes across task variants. EAST modifies both action distributions and model-generated thoughts, shifting them toward expressions of uncertainty and caution, revealing that LLMs encode explicit representations of decision uncertainty that can be directly manipulated.

*4.1.8 SAC.* Sparse Activation Control (SAC) enables multi-dimensional trustworthiness improvements in large-language models (LLMs) by identifying and modifying task-specific attention heads without interfering with general performance [229]. Unlike traditional representation engineering methods that struggle with simultaneous control of multiple behaviors, SAC leverages the sparsity and independence of attention heads to separately regulate safety, factuality, and bias. By employing path patching and Gaussian Mixture Models (GMM), SAC precisely identifies and manipulates the activations of causally relevant components, achieving task-specific control without degrading overall model capabilities. Experimental results on LLaMA models show that SAC outperforms existing methods in multi-task alignment, avoiding the control conflicts that arise when multiple behavioral modifications interact. This work highlights the potential of sparsity-based interventions as a structured and effective alternative to reinforcement learning from human feedback (RLHF) for trustworthiness alignment.

## 4.2   Dynamic Strength Vectors

Dynamic strength approaches adjust model activations during inference in a way that changes depending on the input or the stage of generation. Rather than adding a fixed, constant vector of certain magnitude, these methods change the magnitude of hidden activations—sometimes by multiplying them with learned factors—to either magnify or reduce certain signals. These methods adapt the intensity of activation changes to better guide model outputs without compromising overall language ability or adjusting the vectors to deal with a wider variety of interventions. Activation scaling modifies the magnitude of specific activation vectors in language models to steer outputs while preserving interpretability

Stoehr et al. [197] implements activation scaling through methods like SteerVec and ActivScalar [197]. Unlike additive steering vectors, which alter both direction and magnitude, activation scaling applies learned multiplicative scalars to existing model activations, strengthening or weakening task-relevant representations. The

method is optimized using a three-term objective balancing effectiveness (correcting model outputs), faithfulness (minimizing unintended side effects), and minimality (sparsely modifying activations). A dynamic variant, Dyn-Scalar, further generalizes activation scaling by learning functions of activation vectors, allowing interventions to adapt to varying prompt lengths without predefined positions. SADI replaces the fixed approach with a dynamic steering vector that adjusts model activations based on contrastive activation differences, allowing targeted interventions at inference time [206]. It applies binary masking to identify critical model components, such as attention heads and feedforward neurons, and modifies activations based on input semantics without requiring additional training. Dynamic Activation Composition (DAC) enables multi-property steering of LLMs by dynamically adjusting activation steering intensity throughout generation [188]. DAC continuously modulates the strength of multiple steering vectors, balancing conditioning strength with generation fluency. Experiments show that DAC effectively maintains strong conditioning for properties like safety, formality, and language while reducing fluency degradation. RE-CONTROL dynamically edits LLM hidden representations at test time through an approach derived from control theory and the optimal control formulation [107]. Instead of modifying model weights through fine-tuning, it perturbs hidden states using learned control signals, which enables efficient alignment with human objectives while maintaining computational efficiency. The method formulates autoregressive language models as discrete-time stochastic dynamical systems and optimizes control signals via a trained value function, allowing for flexible intervention during generation. Empirical results show that RE-CONTROL outperforms existing test-time alignment methods like prompting and guided decoding while generalizing well to new inputs. Adaptive Activation Steering (ACT) dynamically modifies model activations to enhance truthfulness in large-language models (LLMs) without fine-tuning [219]. ACT adjusts the steering intensity based on the truthfulness content of activations, ensuring targeted intervention. The approach clusters hallucination-related activations and generates distinct steering vectors for different categories of hallucinations.

## 4.3 Specific Representation Engineering Implementation Choices

Representation engineering relies on the choice of technique and specific hyperparameters. Because of this, there remain several aspects separating the existing approaches from each other. While the theoretical, general approach to RepE implementation is clear, in practice the approaches differ in actual implementation.

*4.3.1 Stimuli Number.* Many methods use benchmarks to construct the steering vector and utilize the remainder of it to test for effectiveness, operating under the assumption that the larger the number of stimuli the better [161]. However, in practice, a higher number of stimuli incurs extra processing cost while improving accuracy only slightly [260]. The number of stimuli varies widely in literature, ranging from as few as one contrastive pair [209], 5-128 strings [260], 30 prompt pairs [188], 64 questions [33], 100 prompts [174], 128 prompts [218], 200 sentences per concept [55], 256 strings [249], 50 percent of the benchmark[219, 247], 500 questions from AdvBench [60], 1000 strings [169], 1092 pairs of data [38], 150 contrastive pairs [220], subsets of benchmarks [134], 3000 to 8000 hateful samples [64], and up to 10,000 contrastive examples [111]. PaCE, uses a large dictionary of over 40,000 concepts to create the final steering vector [140]. SAFT uses k components ranging from 1 to 32 [41].

It's easy to be misled into thinking a representation truly encodes a desired attribute when it is, in fact, merely coincidentally linked. Imagine you're trying to reduce the toxicity of a model. You identify a representation that's correlated with non-toxic responses. However, it turns out that this vector is also linked to the avoidance of certain controversial topics, even if those topics aren't inherently toxic. Intervening on the representation could inadvertently lead the model to censor valuable discussions, thus resulting in a degradation of helpfulness. The choice of stimuli and their number is critical to making sure the right representation has been identified.

Steering vectors can inadvertently encode spurious factors from the training data rather than the intended behavior. This means a steering vector might not correspond to the desired concept and may only be effective due to unintended biases. For example, a dataset used to create a "truthful" steering vector containing mostly

Fig. 9.  Representation of Parameters.

simple sentences might only steer the model towards generating simple sentences rather than genuinely truthful ones. Spurious correlations might lead to steering vectors that amplify irrelevant features, causing inconsistent activation magnitudes across different inputs. This is problematic, as the theory behind what constitutes an optimal intervention is still imprecise. With so many different parameters for identifying representations, it becomes difficult to understand the optimal choice of stimuli and compare between methods.

*4.3.2  Token Choice.* Token choice refers to which token(s) and their corresponding activations are used to construct the steering vectors. Selecting the right tokens from which to extract representations is crucial, as it directly influences the effectiveness and preciseness of the steering vector. The main differences in implementation are whether the steering vectors are calculated using a single token's activation, or across the activations of

multiple tokens. Single-token approaches require formatting the contrastive inputs such that the inputs differ only by a single term. For example, methods such as CAA and ActAdd focus on the activations for yes/no. [161, 209]. In PACE, the activations from the last token of the generated output are used. This approach allows for the desired behavior to be isolated in a singular token, thus minimizing the complexity and noise that would have otherwise been introduced in an open-ended generation settings. Common methods for choosing this token include binary choice tokens like yes/no [97, 141, 161, 209] the token preceding the model's prediction, or the last token in the models' output[260]. In contrast, multiple-token approaches take the differences in activations across multiple tokens as opposed to a single one. Some methods such as Mean-centering [96], Category-wise Steering Vectors [20], or specific RepEng applications for more general representations use more token activations to compute [260].

*4.3.3 Intervention Strength.* After a steering vector has been created, it is usually multiplied by a certain value before intervention. If this value is higher than 1, it implies the effect would be higher than if the representation was directly added through for example prompt engineering. If the scaling factor (which controls the strength of the steering intervention) is too large, adding the steering vector can drastically alter the activation norms, violating the magnitude consistency property and potentially leading to non-sensical output. Conversely, if the scaling factor is too low to preserve the activation norms, the steering vector may not be strong enough to shift the activation towards the desired behavior, hindering editing performance. Stronger scaling can disrupt coherence, while weaker scaling may be ineffective. Intervention strength beyond the steering vector norm causes a quadratic decrease in helpfulness (otherwise known as fluency), showing additional steering comes at a performance [225]. This means that only small enough interventions yield efficient steering with minimal performance loss. Excessive intervention strength degrades model coherence, leading to nonsensical generations [209, 214]. However, when applying the intervention on one layer only, interventions with large co-efficents (larger than 1) do not lead to a visible deterioration of performance [161]. Different studies vary intervention sizes. Some use low values of 0.25 [260]. Some use high values from 2 to 4 [33, 259] or perform grid search to find the optimal intervention strength [111]. Adaptive vector application strategies calculate intervention strength at inference [107, 188, 197, 206, 219]. Interventions can also be negative. That is, a particular representation can be multiplied by a negative number to achieve its opposite. Intervention strength is one of the parameters affecting total effect on steering but cannot be analyzed in isolation- steering type and layer choices have been shown to affect optimal intervention strength.

*4.3.4 Separation Method.* Once the neural activity has been identified, a classifier is required in order to map the models' internal activations to a specific concept or representation. The canonical RepE paper introduces both supervised linear models, including linear probing and differences between cluster means, as well as unsupervised linear models, such as PCA and K-means, for mapping internal states to concepts. However, benchmarking those methods against each other shows that PCA-based techniques are not optimal for identifying the target representation [91]. Inference-time Interventions (ITI) fits a binary classifier onto the model's internal activations using the training, highlighting differences between generation accuracy (the model's output) and probe accuracy (the classifier output). PACE employs the same PCA-based classifer as the original representation engineering paper. Other alternatives to PCA have been explored, such as the Average Gradient Outer Product (AGOP) [11].

*4.3.5 Layers.* Some interventions aim to add the contrast vector to one layer only [20, 32, 32, 36, 61, 97, 161]. Those are always middle layers, usually 15 or 16 for 7-8B parameter models [41, 97, 174, 227, 259]. Some interventions add to a subset of layers [134, 167, 169, 218, 247, 249, 254, 260]. Sometimes interventions happen on consecutive layers (e.g. 12-24 [218] ) and sometimes not (e.g. 23-25 layers and 29-31 [254]). The rationale for choosing these layers is either inspired by other studies or comes from the best performing probing layers [167, 247]. No studies intervene on early layers exclusively. This is confirmed by the fact that probing classification performance

improves with deeper layers [4]. This middle layer intervention strategy is supported by hyperparameter brute force search that shows intervention on one specific layer is approximately the same as on all of them, and that an intervention on the middle layers is the most effective [199]. However, other studies find hyperparameter searches leading to different optimal layers for each representation and model [198].

Probing shows that in early layers, non-linear probing brings best effects, but the linearity of representations increases in later layers [31]. Perhaps this can explain why interventions relying on linearity of representations perform well in middle layers. In fact, sparse probing explains that the middle layers often contain more dedicated neurons for higher-level linguistic and contextual features compared to earlier layers [73].

Probing shows that LLMs encode context knowledge in a layer-dependent manner, with upper layers retaining more contextual information, while lower layers primarily focus on entity-related features [98]. Intermediate layers in large-language models (LLMs) consistently provide more informative representations for downstream tasks than final layers [13, 195]. Deeper layers are shown to contain more factual knowledge [2, 27, 82, 145, 208].

Studies dedicated to layers show theory explaining these phenomena. Skean et al. [195] judges representation quality at various layers through metrics such as prompt entropy, curvature, and augmentation invariance, showing that models undergo significant information compression at intermediate depths. Comparisons between Transformers and State Space Models (SSMs) show that Transformers exhibit more pronounced representational transformations, while SSMs maintain stable intermediate-layer representations. Training progression analysis indicates that the largest improvements in representation quality occur in intermediate layers, reinforcing their role in effective feature extraction. Intermediate layers exhibit bi-modal entropy distributions, where activations cluster into two distinct groups: one with high entropy, encoding diverse and context-sensitive information, and another with low entropy, representing more stable and deterministic knowledge. This suggests that intermediate layers naturally separate different types of information processing, which has significant implications for representation engineering. Activation steering and concept vector extraction can be more effectively applied to the low-entropy subset for stable modifications, while high-entropy activations may be leveraged for tasks requiring greater contextual flexibility. The structured separation also explains why mid-layer interventions, such as Inference-Time Intervention (ITI) and Sparse Autoencoder Steering, are more effective than modifications at early or final layers. Additionally, targeted entropy-aware interventions could enhance model alignment by suppressing high-entropy features to reduce hallucinations or amplifying low-entropy representations to improve factual consistency. Recognizing bimodal entropy as an emergent property of LLMs provides a framework for designing more precise and efficient steering methods that align with the model's intrinsic representational structure.

However, while several studies support this finding, these have been conducted mostly on small-scale models (under 10B parameters). Also, function vectors work better in early layers [207], suggesting optimal layer depends on the intervention type. As the models get more complex, it becomes critical to conduct studies of larger size models to identify the optimal intervention layers using the new metrics. Therefore, as a fallback option, combinatorial testing, involving exhaustive brute force hyperparameter search through all existing combinations of layers is employed to identify the most effective combination of hyperparameters to detect and edit a representation. This is important because several studies show concepts do not remain stable across the network layers, implying constant interventions across all layers should not be applied [157]. Experiments are primarily conducted on Llama-7B and 8B, with only some studies testing on larger models like Mixtral-8x7b [174] or 13B parameter models [33, 218]. Nonetheless, some studies apply interventions to all layers [55, 140, 209, 220, 225].

*4.3.6   Steering Type.* In the process of altering the output response style, the intervention may cause issues due to too strong of an intervention and its duplicate effect from intervening on multiple tokens. For instance, one can consider the intervention switching the token probability to be more positive. If the first word of the response

is changed from a normal to a positive one, the intervention effect for the second and any subsequent tokens will be magnified, as the probabilities of subsequent tokens depend on the previous sequence of tokens. This problem led to the emergence of several steering types. The **Initial** steering type only applies the intervention on the first token of the response [113, 207, 209]. The **Constant** steering type maintains the same steering on each token of the response [32, 132, 140, 161]. The **Diminishing** steering type starts with a full-strength response but reduces the intervention strength for each subsequent token [188].

## 4.4 Multiple Model Methods

Instead of operating on only one model, some techniques utilize multiple models to transplant concepts between models, improve the performance of the model after alignment (KTS [196]), and improve model safety (InferAligner [218], CONTRANS [55] or KTCR [64]). By utilizing the fact that different models have different infrastructures, one can leverage differing levels of safety between them to improve one of them or keep the performance stable after an intervention. KL-Then-Steer (KTS) is a method that fine-tunes models to minimize Kullback–Leibler (KL) divergence between steered and unsteered outputs, reducing the impact of steering on the general capabilities of the model [196]. KTS prevents 44 percent of jailbreak attacks on Llama-2-Chat-7B while preserving general capabilities, outperforming system prompts and LoRA fine-tuning in maintaining model performance. The method generalizes to other interventions. The results suggest that KTS enables reliable control of the model without performance side-effects. In a similar fashion, InferAligner modifies model activations at inference time using safety steering vectors (SSVs) extracted from safety-aligned models [218]. These vectors shift activations in specific representational directions to guide responses away from harmful outputs while preserving performance on downstream tasks. Steering vectors are transferable across models, so alignment-related representations exist in structured and reusable forms. Activation shifts effectively enforce harmlessness constraints without degrading the model's general capabilities. Also using multiple models for safety, CONTRANS shows how alignment from smaller, aligned language models can be transferred to a larger, unaligned model using concept transplantation [55]. CONTRANS extracts concept vectors from a source model using a small set of positive and negative examples, reformulates them through affine transformation to match the target model's feature space, and then integrates them into the target model's residual stream. Experimental results show that this method successfully transfers alignment concepts like truthfulness and fairness across different model sizes and families. CONTRANS achieves competitive performance compared to instruction tuning and reinforcement learning-based alignment while avoiding the high computational costs associated with those methods. Also attempting to work with safety concepts, KTCR shows a new method for implicit hate speech detection by refining conceptual representations through knowledge transfer [64]. KTCR employs a teacher-student approach, where a pre-trained teacher model guides a student model in learning subtle hate speech patterns using concept activation vectors and prototype alignment. Unlike traditional data augmentation methods, KTCR not only introduces new implicit hate examples but also systematically refines how models internalize and distinguish such patterns. Therefore, this method is particularly powerful for detecting and editing concepts the definition of which changes over time.

## 4.5 Fine-tuning Alternatives

Fine-tuning alternatives take a different approach from traditional weight-based fine-tuning by working directly in the model's hidden representation space. Instead of updating large portions of the network's weights, these methods target internal activations or add lightweight adapters, reducing both the risk of unintended side effects and the computational cost. By enhanced steerability and lightweight computational cost, representation engineering helps to change the performance and style of the model, hence fulfilling a similar role as fine-tuning but at a greatly decreased cost.

*4.5.1  TruthX.* Truthfulness in large-language models (LLMs) is not solely determined by their knowledge but also by how their internal representations influence response generation [247]. Some erroneous activations within these representations can cause hallucinations even when the correct knowledge is present. By identifying and modifying the truthful space of these representations, LLMs can be guided to produce more truthful responses without sacrificing their generative capabilities. The effectiveness of such interventions, like TruthX, depends on mapping representations into distinct semantic and truthful latent spaces, then applying contrastive learning to separate truthfulness from other linguistic properties. Editing these representations in the truthful space rather than directly manipulating outputs or attention patterns enables more consistent and controllable truthfulness enhancement across different LLMs.

*4.5.2  SAFT.* Safety-Aware Fine-Tuning (SAFT) aims to mitigate the risks posed by harmful data in the fine-tuning process of large-language models (LLMs), similar to the approach in Rosati et al. [182] and Choi et al [41]. The framework employs a scoring function that detects and removes potentially harmful samples by analyzing the subspace information of LLM embeddings. Empirical results demonstrate that SAFT can reduce harmfulness by up to 27.8 percent across various models and contamination levels without significantly compromising model helpfulness. Traditionally, the paper examines the representation space of LLMs and how harmful samples can be identified through singular vector decomposition, contributing to a broader understanding of representation engineering in LLM fine-tuning.

*4.5.3  ReFT.* Representation Finetuning (ReFT) provides a parameter-efficient alternative to traditional weight-based fine-tuning by modifying hidden representations in large-language models (LLMs) instead of updating model weights [227]. This method utilizes Low-rank Linear Subspace ReFT (LoReFT), which apply structured interventions in learned activation subspaces to improve model adaptation while using 15x–65x fewer parameters than LoRA. LoReFT consistently outperforms state-of-the-art parameter-efficient fine-tuning (PEFT) methods on common sense reasoning, instruction-following, and natural language understanding tasks. Unlike prior PEFTs, which modify a subset of weights, LoReFT operates entirely within the model's latent space, preserving pre-trained knowledge while steering model behavior toward task-specific objectives. ReFT is therefore a cheaper model adaptation method that reduces the overhead of full fine-tuning.

*4.5.4  RAHF.* Representation Alignment from Human Feedback (RAHF) is a representation engineering alternative to RLHF [134]. RAHF can extract and manipulate activity patterns humans value positively, aligning model outputs with human preferences faster and cheaper than traditional fine-tuning. By leveraging contrastive learning and activation-based modifications, RAHF provides a scalable way to adjust LLM behavior without requiring extensive additional training or reward modeling.

*4.5.5  LoRRA.* [260] uses LoRRA (Low-Rank Representation Adaptation), a lightweight fine-tuning method that integrates low-rank adapters into attention weights and optimize them using a representation-based loss function, such as the Contrast Vector. These adapters, referred to as controllers, are trained during model adaptation and later merged into the model, ensuring no additional computational cost during inference.

## 4.6  Sparse Autoencoders-based approaches

Sparse autoencoder-based approaches use models that learn compact, focused representations of hidden activations to identify the signals most relevant to steering language model behavior. This section provides an overview of methods using inspiration from both mechanistic interpretability and representation engineering to combine sparse autoencoders with activation steering. By training sparse autoencoders to capture key features, these methods can identify areas where the model's stored knowledge conflicts with its current context or where

specific semantic cues are encoded. These methods show how sparse representations can improve both the interpretability and precision of activation-steering interventions.

*4.6.1 SPARE.* Sparse Auto-Encoder-based Representation Engineering (SPARE) identifies signals of knowledge conflict in the residual stream of mid-layers, allowing precise detection of instances where contextual and parametric knowledge disagree [254]. SPARE extracts functional features from pre-trained sparse autoencoders (SAEs) that govern knowledge selection, enabling targeted intervention to prioritize either memory-based or context-based knowledge. Experimental results on open-domain question-answering tasks show that SPARE improves knowledge selection accuracy by 10 percent over existing representation engineering methods and 15 percent over contrastive decoding.

*4.6.2 SAE-TS.* Chalnev et al. [36] present a method to refine steering vectors to better control the behavior of language models while minimizing unintended effects through SAE-Targeted Steering (SAE-TS), which utilizes sparse autoencoders (SAEs) to measure the causal impact of steering vectors. Through this, more granular model control is achieved. The approach differs from existing methods by directly optimizing steering interventions to align with specific SAE features, avoiding the unpredictability of contrastive activation addition (CAA) and direct SAE feature steering. The paper compares SAE-TS with other steering methods across various tasks, demonstrating its superior ability to maintain coherence while effectively steering model behavior. In addition, the paper develops an interactive tool, EffectVis, to visualize feature effects and improve interpretability, contributing to ongoing research in mechanistic interpretability and representation engineering.

*4.6.3 Self-knowledge Sparse Autoencoders.* Ferrando et al. [61] examine the ability of LLM to recognize known and unknown entities. Using sparse autoencoders, it identifies specific directions in the model's representation space that encode self-knowledge, checking whether the model can recall factual information about an entity. These learned directions can causally influence model behavior, such as inducing refusals for unknown entities or prompting hallucinations when recognition is manipulated. The findings suggest that the fine-tuning of the chat model repurposes existing entity recognition mechanisms rather than creating entirely new refusal strategies. Furthermore, the study highlights how these entity recognition directions regulate attention mechanisms, affecting factual recall and uncertainty expression.

## 4.7 Other Approaches

Due to its relatively recent emergence, representation engineering has evolved to multiple new applications. This section outlines key use cases to present methods for applying representation engineering to new applications.

*4.7.1 ARE.* Adversarial Representation Engineering (ARE) enables targeted model editing by leveraging representation based discriminators to refine internal activations [249]. This is useful for detecting jailbreaking attempts by monitoring for unusual activation patterns and counteracting harmful requests by injecting safety steering vectors to guide models towards safe outputs. Unlike previous methods that rely on direct representation vector manipulation, ARE employs adversarial learning between a generator (the LLM) and a discriminator, iteratively refining activation patterns for alignment and jailbreak tasks. Experiments show that ARE can reduce refusal rates on harmful prompts from 20 percent to under 1 percent in LLaMA-2 while also enhancing truthfulness in TruthfulQA benchmarks. The method, however, has also been shown to have bidirectional editing capabilities, enabling alignment reinforcement, but also its removal.

*4.7.2 Conceptors.* Unlike traditional steering methods that rely on single additive steering vectors, Conceptors represent sets of activation vectors as ellipsoidal regions [169]. This allows for more structured manipulation of activations than a singular additive vector. Conceptors create a conceptor matrix from cached activations, which is then applied as a soft projection at inference time. Experiments on GPT-J and GPT-NeoX demonstrate that
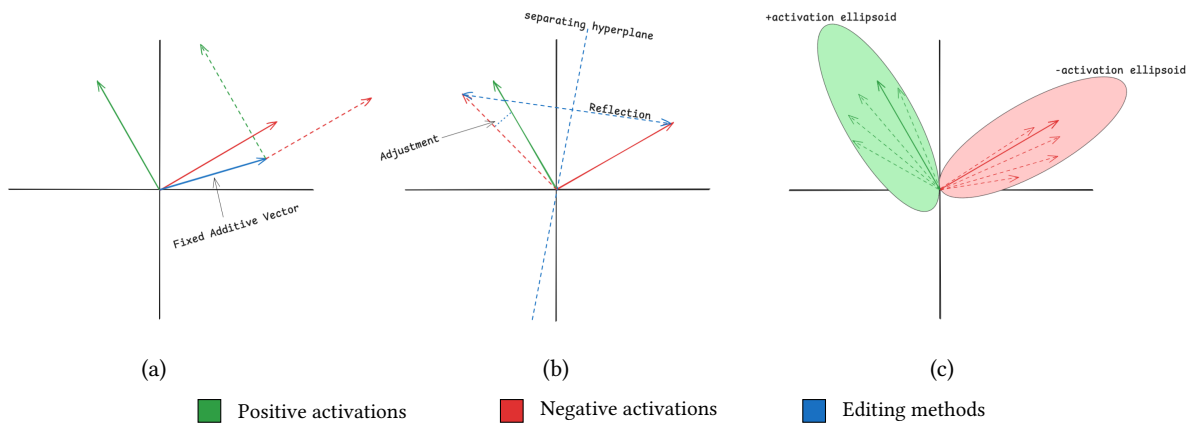
Fig. 10. Geometric shift comparison of a few RepE techniques within the activation space. (a) represents additive steering [161], which translates activation vectors by a fixed steering vector, (b) depicts Householder Pseudo-Rotation (HPR) [167], which first reflects negative activations through a learned separating hyperplane and then adjusts the reflection to attain 90°(this is an approximated rotation transformation), (c) depicts conceptor steering [169], which bounds all activation vectors by softly projecting them onto a target ellipsoid.

Conceptors outperform additive steering across multiple function-execution tasks, including antonym retrieval and language translation. Boolean operations on Conceptors, such as AND and OR, can be performed to combine multiple steering objectives more effectively than additive vector combination.

*4.7.3 Circuit Breakers.* Circuit Breakers modify internal representations in large-language models (LLMs) to block harmful outputs at inference time, providing a more reliable safeguard than traditional refusal training [259]. The method, called Representation Rerouting (RR), redirects harmful activation patterns to predefined refusal states, preventing models from completing dangerous responses while preserving general capabilities. Circuit breakers can also effectively defend against adversarial image samples [177]. Unlike adversarial training, which targets specific attack methods, RR operates on a model's latent space, making it robust against unseen attacks, including embedding-space manipulations and multimodal adversarial inputs. Experiments on LLaMA-3 and Mistral-7B show that RR reduces harmful output rates by over 90 percent under diverse adversarial conditions while maintaining performance on standard benchmarks.

*4.7.4 LEGEND.* **L**everaging r**E**presentation en**G**ineering to annotate prefer**EN**ce **D**atasets (LEGEND) automatically annotates safety margins in preference datasets using representation engineering [60]. LEGEND identifies a "safety direction" within an LLM's embedding space by computing the difference between embeddings of harmful and harmless responses to increase safety with no need for extra training. As a purely inference time-method, it is more computationally efficient than the alternative margin annotation approaches. Experiments demonstrate that LEGEND improves reward modeling accuracy and enhances harmless alignment for LLMs, achieving comparable or superior results to ensemble-based methods while reducing annotation time significantly.

## 4.8 Open source applications

Representation engineering methods are also helpful for resolving issues related to open-source models. Representation Encoding Fingerprints (REEF) helps identify whether a suspect language model (LLM) is derived from a victim model by analyzing feature representations rather than weights [245]. Unlike weight-based fingerprinting,

which is vulnerable to fine-tuning, pruning, and model merging, REEF computes the Centered Kernel Alignment (CKA) similarity between model representations on the same input samples, providing intellectual property protection. Representation Noising (RepNoise) is a defense mechanism designed to prevent harmful fine-tuning attacks (HFAs) on large-language models by removing information about harmful representations from model activations before an attacker gains access [182]. Unlike traditional safety guardrails that can be circumvented through fine-tuning, RepNoise alters intermediate activations such that recovering harmful behaviors through additional training becomes significantly more difficult.

## 4.9 Benefits of Representation Engineering

Representation engineering seeks to answer questions critical to the development of model AI. In the process of representation learning, models generalize information to perform well on unseen tasks, transforming specific data points into useful heuristics for solving other tasks [15]. A large part of deep learning is dedicated to modifying training methods so that they make those generalizations better, for example through regularization[109]. In that, it aims to resolve the same fundamental problems as reinforcement learning from human feedback[101], direct preference optimization, fine-tuning, mechanistic interpretability and other approaches.

*Semantic Meaning.* Representation engineering facilitates the extraction of model representations and assigns them semantic meanings, enabling both the inspection of a model's reasoning and interpretation by human operators. By elucidating internal model concepts, this approach offers several beneficial applications. For instance, it allows for the establishment of precise guardrails to detect and mitigate harmful or biased patterns [108]. Additionally, it enhances model optimization by enabling users to specify desired activation pattern shifts using natural language, as demonstrated by the Activation Addition technique [209]. Furthermore, representation engineering aids in comparing models to select the most suitable one for specific use cases by providing semantically rich characterizations of model properties [173], thereby reducing reliance on generic benchmarks.

*Interpretability and Explainability.* Representation engineering provides a direct window into the internal reasoning of models by analyzing and manipulating their learned representations. This approach allows for the identification and mitigation of biases, harmful behaviors, and unintended model outputs at their root. By inspecting and assigning semantic meaning to these representations, it becomes possible to detect and address issues such as model "breaking free" or acting against user intentions [119, 249]. This level of transparency not only enhances trust in model behavior, but also enables the development of more robust safeguards against undesirable outcomes. Unlike post-hoc explainability methods, representation engineering operates at the foundational level of model reasoning, offering a more systematic and actionable approach to understanding and controlling model behavior.

*Alignment and Personalization.* Representation engineering enables fine-grained control over model behavior by shifting representations to align with user-specific values and goals. Unlike traditional alignment methods, which often impose broad social or ethical values determined by model developers, representation engineering allows for individualized customization. Users can amplify or suppress specific representations to reflect their unique preferences, such as prioritizing honesty, creativity, or domain-specific expertise [32, 210, 250, 260]. This flexibility ensures that models can be tailored to diverse use cases without being constrained by a one-size-fits-all approach. By decoupling alignment from predefined social norms, representation engineering empowers users to define and enforce their own value systems, making models more adaptable and personally relevant.

*Increasing Performance.* Representation engineering improves model performance by enabling precise interpretations as well as modifications to activation patterns. By identifying and amplifying representations associated with desired behaviors—such as accuracy, coherence, or task-specific expertise—users can optimize models for

specific objectives [27, 33, 80, 134, 169, 174, 188, 197, 227, 254, 260]. This approach contrasts with traditional fine-tuning, which often requires extensive retraining and may inadvertently degrade performance on unrelated tasks. Representation engineering allows for targeted adjustments that improve performance without compromising generalization. For example, in a medical diagnosis task, representations linked to "clinical accuracy" can be strengthened, leading to more reliable outputs. This precision makes representation engineering a powerful tool for optimizing models in a cost-effective and efficient manner.

*Inference Time Control.* Representation engineering operates at inference time, eliminating the need for costly and invasive retraining. Traditional training methods are resource-intensive, unpredictable, and often require extensive setup. In contrast, inference-time control allows users to experiment with different configurations on-the-fly, adjusting representations to achieve desired outcomes without modifying the underlying model [181, 209]. This approach reduces latency, lowers costs, and makes advanced AI capabilities more accessible. By using pretrained models, users can achieve precise control over model behavior without the overhead of additional training that often introduces high costs and takes a long time to implement.

## 5   COMPARISON WITH OTHER APPROACHES

| Method | Steer Strength | Semantic Meaning | Computational Cost | Inference Time |
|---|---|---|---|---|
| Representation Engineering | 🟢 | 🟢 | 🟢 | 🟢 |
| Sparse Autoencoders | 🟢 | 🟢 | 🟡 | 🟡 |
| Prompt Engineering | 🔴 | 🟡 | 🟢 | 🟡 |
| Soft Prompts | 🟡 | 🔴 | 🟢 | 🟢 |
| Fine-Tuning | 🟡 | 🔴 | 🔴 | 🔴 |
| Mechanistic Interpretability | 🟡 | 🟢 | 🔴 | 🟡 |

Table 2.  Comparison of Representation Engineering with Other Approaches.

🟢 High    🟡 Medium    🔴 Low

### 5.1   Sparse Autoencoders

Sparse autoencoders are neural networks designed to learn efficient feature representations by enforcing sparsity in the hidden layer activations, meaning only a small subset of neurons are active at a given time [16]. Sparse autoencoders usually feature the addition of a sparsity constraint, such as L1 regularization or Kullback-Leibler divergence, to approximate most activations to zero while preserving important features [243]. Sparse autoencoders learn feature representations by enforcing sparsity constraints within the model activations, while representation engineering involves manually designing or modifying features to improve model performance. Unlike representation engineering, which relies on human intuition and domain knowledge, sparse autoencoders automatically discover structured, high-dimensional representations through unsupervised learning. Sparse autoencoders have been extensively applied to identify circuits of casual importance in large-language models [146].

However, sparse autoencoders are not good for identifying steering vectors [149]. Steering vectors are out-of-distribution for SAEs, as they have significantly lower L2 norms than typical model activations, causing the SAE encoder bias to dominate the decomposition, and (2) SAEs enforce non-negative reconstruction coefficients, preventing them from capturing meaningful negative projections in feature directions. Furthermore, sparse autoencoders are much more computationally demanding and provide more difficult to interpret semantic meaning.

## 5.2 Prompt Engineering

Prompt engineering refers to a wide range of methods that aim to use more precise combinations of tokens on the input and system prompts to steer the model towards particular results [184]. Prompt engineering focuses on crafting inputs to influence a model's outputs, whereas representation engineering modifies how data is encoded within the model itself [148]. Unlike representation engineering, which alters feature representations inside the model, prompt engineering works externally by optimizing input structures without changing the model's learned parameters. Therefore, the actual effect tokens have on the model internal representations is less interpretable and the steering less powerful.

*Soft Prompts.* Soft prompts are trainable input embeddings that guide frozen language models without modifying their parameters, exhibiting unique properties distinct from natural language prompts [9]. Unlike discrete prompts, soft prompts occupy separate regions in embedding space, show heightened sensitivity to directional perturbations, and enable parameter-efficient adaptation for downstream tasks, as demonstrated in prefix-tuning [120], prompt tuning [114], and P-Tuning [135]. Studies show that soft prompts can enhance generalization in low-resource settings [120], improve stability through continuous embeddings [135], and support novel applications like unlearning without weight updates [19]. However, their interpretability remains limited, and their susceptibility to adversarial manipulation raises security concerns [9]. New, more efficent soft prompting methods, such as InfoPrompt, use information-theoretic objectives to optimize prompt initialization and task relevance, accelerating convergence and outperforming conventional tuning methods [226].

Soft prompting involves adding a learnable vector embeddings to the input sequence of a frozen language model, while representation engineering directly manipulates the internal activation patterns of the model without modifying its weights. In soft prompting, only the parameters of the added prompt vectors are updated during training, typically using a small prompt encoder, whereas representation engineering involves a more complex intervention across multiple layers of the model's hidden states. The implementation of soft prompting requires minimal changes to the model architecture, often just extending the input processing pipeline, but representation engineering requires white-box access to the model to fully optimise the representation.

## 5.3 Fine-tuning

Fine-tuning refers to the process of shifting the weights of a pre-trained model by training it on a task-specific dataset [53]. Fine-tuning utilizes transfer learning to shift the model from a general purpose one to a more specific one, adding skills the model previously did not have or domain expertise that was not represented in the training data [256]. Like representation engineering, fine-tuning seeks to adapt a pre-trained model to a specific task or domain. Unlike representation engineering, the changes made through fine-tuning do not have a semantic meaning. Fine-tuning can be computationally expensive and may require significant resources, especially for large models, while representation engineering methods like Representation Fine-Tuning (ReFT) can be more parameter-efficient and less resource-intensive [227]. In fact, using representation engineering to fine-tune representations more efficently leads to modifying less than 1 percent of the overall model representations to achieve comparable shifts in performance to traditional fine-tuning and are further optimized by solutions such as Low-rank Linear Subspace ReFT (LoReFT), a representation-engineering based alternative to LoRa adapters.

When compared with fine-tuning directly, activation steering can additionally steer the model beyond what is possible with fine-tuning [161, 209]. However, combining these together leads to unpredictible interactions [161]. In comparison to fine-tuning, representation engineering can allow to change intervention at steering time with "online steering", which fine-tuning does not allow [209].

## 5.4 Mechanistic Interpretability

Mechanistic interpretability analyzes the contributions of specific parts of the network to model outcomes [17, 176]. In that, it looks at granular parts of the network and attempts to extract information on the interactions within the network. Mechanistic interpretablity uses neurons and circuits as fundamental units of analysis, hence focusing on identifying particular low-level mechanisms through which a model undertakes a decision [251, 252].

However, it is demanding to identify and locate circuits within a network. There is no guarantee that all parts of the network can be interpreted as circuits [81, 260]. The research that publishes successful examples of circuits often presents a misleading picture of how difficult it is to assign an actual representation to the circuit, as negative results are not published [17]. Furthermore, given the use of interpretability techniques in training, worries about the mechanistic interpretability techniques actually intensifying the adversarial pressure against interpretability increase. This means that as new techniques to detect power-seeking or otherwise deceptive behaviour of LLMs are used, the deceptive model will make it harder to find such interpretable circuits but not change its intentions and behavior.

Mechanistic interpretability has the same goal as representation engineering. Mechanistic interpretability and representation engineering can be seen as opposing approaches: one takes the bottom-up perspective of analyzing building blocks of the network to understand its overall performance while the other adopts a top-down approach, using global representations to extract meaning [260]. However, in practice, the methods are hard to clearly distinguish. Activation patching for example, can be seen as a specific case of a very granular representation engineering intervention. In comparison with existing mechanistic interpretability methods like automated circuit discovery [43], attribution patching [201], causal scrubbing [24] or sparse autoencoder-based neuron interpretation [45], representation engineering tends to produce more immediate and controllable changes in model behavior in a more lightweight manner than SA-based mechanistic interpretability approaches and enables direct manipulation of high-level concepts without requiring a full mechanistic understanding of the underlying circuits.

## 5.5 Combining Representation Engineering With Other Methods

Using activation steering does not prevent the users from using alternatives. Representation engineering have also been successfully used as a part of a "modular LLM" architecture, selectively adding representation engineering interventions combined with other interventions to steer the model [228]. Different compatibility metrics have been developed, showing the order of these interventions matters for their overall effectiveness [106].

## 6 EVALUATION

Understanding whether the representation has been sucessfully detected and steered is a critical component of building a reliable representation engineering framework. The lack of a standardized approach hinders the effectiveness and adoption of representation engineering. In particular, open-ended generation evaluation is important. Interventions that perform well in multiple-choice formats can fail in open-ended scenarios. CAA [161] has been shown to successfully steer the model in a multiple-choice setup but failed in an open-ended context using the same prompts and interventions [170]. In understand the effectivness of a particular intervention, model size and architecture are important. Representation engineering works better for larger models [20]. This is potentially because larger models build more detectable linear representations [74]. Different models have differing representations, but steering vectors are transferable over architectures and model sizes [32].

### 6.1 Measures of steering

*6.1.1 Task Accuracy.* Most evaluations report simply the changes in performance on a particular dataset. For hallucination detection, this is usually TruthfulQA [38, 55, 196, 249, 260], with one paper evaluating on NQSwap

and Macnoise [254]. Evaluations on TruthfulQA, however, often make use only of the multiple-choice version of the benchmark [32, 118, 167]. Similarly, MMLU is also a multiple choice benchmark. This is potentially problematic, as it does not allow to examine the degradation of performance as the result of steering.

For safety interventions, methods are evaluated using different benchmarks, including BeaverTails [20, 41, 188], AdvBench [32, 33], HarmBench [259], ToxiGen [55, 220] or other more niche benchmarks. Similarly, when evaluating general model performance, MMLU [38, 140, 161], Alpaca [20, 188] or CatQA [20, 209] are used, with several other benchmarks such as ARC-Easy or AQuA used in one paper each.

The fact that outputs change or remain stable in a given direction does not necessarily indicate the presence of a meaningful linear representation, as changes might result from unrelated geometric artifacts in the model's representation space [165]. Instead of relying solely on benchmark outputs, token probabilities in a given context should be evaluated to determine whether interventions in the representation space effectively steer model behavior [165]. Concept representations should be tested for their alignment with counterfactual token pairs to validate their linear structure, hence showing that differences between words expressing the same conceptual change point in a consistent direction. For example, CAA [161] aimed at inducing myopic behavior showed nearly equal probabilities for both myopic and non-myopic tokens, suggesting that difference in outputs between steered and unsteered output depends on sampling randomness rather than a fundamental shift in model behavior [170]. Rather than using arbitrary inner products, the causal inner product should be estimated and assessed for its ability to enforce orthogonality between causally separable concepts. This allows to show that changes in one concept do not unintentionally affect another, degrading fluency on other tasks. In addition, representation reading can be used for validation to verify that the computed concept representations can predict target attributes in a logit-linear way. Also, shifts in token probabilities should be analysed while increasing the intervention strength and checking if that leads to the expected emergence of higher probability of the target concept in next-token predictions [165].

Alternatives to benchmark performance include the Kullback–Leibler divergence (KL) of the model's next-token prediction distribution post- versus pre-intervention [118], reduction in perplexity [209] or latent separation scores, logit difference measurements, and attention score changes [61].



Fig. 11. Recommended Evaluation Pipeline

## 6.2 Measures of fluency

Steering vectors can reduce harmful output frequency but often degrade performance on benign inputs [196]. While interventions can often improve alignment metrics, such as decreasing the rate of successful jailbreaks by 30% for specific models, it's also possible that these same interventions, without proper consideration, might reduce the model's accuracy on standard question-answering or other general capabilities [210, 225, 254], sometimes by as much as 15-20% on MT-Bench [196, 259]. This is often seen where targeting too strongly a desired behavior will also break other model behaviors. In particular, for RepE [260] interventions beyond the steering vector norm lead to rapid degradation of performance [225]. Large steering vectors guarantee intervention success but reduce

model outputs to near-random guessing [225]. Therefore, each intervention must be precisely calibrated to avoid excessive degradation of general capabilities. The lack of theoretical validation for the choice of hyperparameters and standardized methods of evaluation exacerbates this problem. Therefore, evaluations often consider both the effect of steering on outputs while measuring linguistic fluency [118, 198]. Several metrics used to track this include perplexity [209], perplexity-normalized effect size (PNES) [214], BLEU [167, 199], ROUGE-L [41] and generation entropy or the weighted average of bi- and tri-gram entropies [29]. Nonetheless, lack of loss of fluency on a particular representation intervention does not imply the model will perform maintain coherence on other tasks that will have different representation parameters assigned by the search of optimal layer or other model parameters.

## 7 OPEN PROBLEMS IN REPRESENTATION ENGINEERING

Despite significant progress in representation engineering, several challenges remain unresolved. This section examines key open problems identified in recent research, their implications for the field's development, and potential directions for future work.

### 7.1 Standardization

*7.1.1 Lack of Standardized Evaluation.* Many activation engineering techniques employ different methods to evaluate the effectiveness of steering. Evaluations often do not assess out-of-distribution inputs, limiting the generalizability of the findings [17]. [170] found that an intervention exhibited corrigible behavior in a multiple-choice format but failed to do so in an open-ended generation setting, highlighting the importance of measuring steering interventions on tasks similar to the eventual use-case. It critiques existing evaluation methods for activation steering in large-language models and proposes a standardized evaluation pipeline to measure intervention effectiveness. It identifies four key missing properties in current assessments: *a)* evaluations should be conducted in open-ended generation contexts, *b)* account for model likelihoods, *c)* allow for standardized cross-behavior comparisons, and *d)* include baseline comparisons. Using this framework, the study evaluates two activation steering methods: Contrastive Activation Addition (CAA) and Inference-Time Intervention (ITI) on behaviors such as truthfulness, corrigibility, and myopia. The results reveal that the effectiveness of the intervention is highly behavior-dependent, and ITI successfully increases the likelihood of true responses, while CAA is more effective in reducing hallucinated outputs than in promoting desired behaviors. The paper also finds that multiple choice evaluations overstate intervention success compared to open-ended settings.

*7.1.2 Systematizing Parameter Interventions.* As outlined in this paper, the exact structure of each intervention differs significantly. The lack of a standardized approach to representation engineering can make it difficult to compare and contrast different methods, especially in the context of the overall intervention strength. The lack of standardized parameters, including stimuli number, token choice (encoder vs. decoder, specific tokens), separation method, intervention strength (constant, diminishing, initial), and target layer(s), makes it difficult to compare results across different studies and hinders the development of generalizable best practices. This variability, coupled with the context-dependent nature of knowledge encoding in LLMs, where representations shift across layers and tokens, creates a complex search space for optimal RepE intervention and representation evaluation. Consequently, computationally expensive brute-force hyperparameter searches need to be implemented to find the optimal intervention for a given use case. The absence of a detailed theoretical framework on optimal choices of these parameters makes practical implementation of RepE a challenge and requires further investigation.

*7.1.3 Lack of Theoretical Grounding.* We lack a solid theoretical framework to explain why interventions work, whether they should be theoretical, and often we see that they work best within a given range of coefficients. For example, it's difficult to predict in advance which layers are most receptive to steering, with some studies

finding that middle layers work best [174, 209, 254], while others find that the efficacy varies based on the specific task. Furthermore, steerability appears to be a property of the dataset rather than the model, since different architectures exhibit similar patterns of reliability and failure [202]. The extent to which each kind of vector can steer effectively is limited by whether the portion of latent space from which they are extracted actually contains a discoverable representation of the information needed to execute the desired behavior [29].

Crucially, the effectiveness of steering vectors has shed insights on the linear representation hypothesis (LRH). LRH claims a stronger global linearity in feature space. These studies show there are representations for which the linear hypothesis is true, but also some for which it likely does not hold [144, 147, 202? ]. Whether the representations are linear likely also depends on model size. Since interventions are tested on small models of similar sizes, it is important to create a generalizable theory that unifies interventions, representations, parameters, and models in a tractable framework that does not rely on empirical search to find optimal parameters but instead develops theoretical explanations for their optimal values.

*7.1.4 Challenges in Generalization.* Steering vectors show high variance in effectiveness, with some datasets exhibiting nearly 50 percent anti-steerable examples, where the intervention produces the opposite of the intended effect. Generalization performance is strongly correlated with the similarity between the source and target prompts, suggesting that steering vectors work best when applied to behaviors the model already tends to produce. Additionally, it can be difficult to predict the behavior of steering vectors. In some cases, steering interventions do not produce interpretable changes in model behavior other than model degradation. Also, the out-of-distribution generalization for many behaviors is not perfect or entirely predictable, making it difficult to precisely control the model [36].

This unpredictability often leads to an empirical, trial-and-error approach, which hinders efficiency and limits our ability to generalize interventions and understand the underlying mechanisms. The sign (direction) and overall effectiveness of steering interventions can vary widely across and within different concepts. This highlights a key limitation: the extent to which steering works is dependent on whether the desired information is already embedded in the model's latent space. The fact that different architectures exhibit similar patterns of reliability and failure further complicates the challenge of generalizing interventions across models [202].

## 7.2 Steering Multimodal Models

Although there is extensive literature on steering large-language models, the rise of multimodal models has not led to the widespread successful application of RepE techniques in modalities like video and images. Control vectors can also be used for improving motion control and forecasting [204]. Multimodal models incorporate tools from multiple modalities to create a cohesive tool for generating outputs across computer vision, text and video. Representations like safety and truthfulness require designing interventions that hold across different modalities, since for jailbreaks can be performed across different modalities. However, research on representation learning shows it is likely that multimodality will introduce additional problems. Additional research is needed to create reliable multimodal representation engineering. In particular, these problems are likely to occur:

*7.2.1 Negative Transfer & Disentanglement of Features.* Expanding the space of interventions to multimodal : *a)* Interventions aimed at enhancing one modality may inadvertently degrade the quality or semantic integrity of another. This misalignment occurs because feature spaces across modalities do not necessarily share common transformation characteristics. *b)* The resulting representations may lose their semantic validity, particularly when interventions fail to account for the distinct ways in which meaning is encoded across different modalities. *c)* The high dimensionality of multimodal feature spaces complicates the identification of meaningful and salient directions for intervention. For example, consider a video depicting culinary preparation. A RepE intervention designed to *increase happiness* might unintentionally amplify the visual prominence of certain ingredients or

generate artifacts in video frames that bear little relation to the subject's emotional state. This phenomenon stems from the challenge of accurately modeling intricate cross-modal dependencies, where the relationships between static elements in images and dynamic actions in video sequences are often subtle and context-dependent.

*7.2.2 Tokenization and Encoding Bottleneck.* RepE approaches require a continuous and differentiable latent space. However, modern multimodal architectures often employ non-differentiable encoding mechanisms, such as tokenization or latent vector quantization, to integrate information across modalities [211]. These discontinuous transformations disrupt gradient propagation. Rando et al. [177] identify this challenge, noting that to enable continuous end-to-end optimization in tokenization-based models, a *tokenizer shortcut* must be implemented. Although this modification facilitates gradient-based adversarial attacks, it introduces additional complications, notably the tendency toward overconfident token predictions. This trade-off highlights a fundamental drawback in the design of robust representation systems: enabling gradient flow for optimization purposes may create vulnerabilities that can be exploited by adversarial methods.

---

**Algorithm 1** Multimodal Representation with Tokenization Challenge

---

**Require:** Text input $x_{text}$, Image input $x_{image}$
**Require:** Encoders $E_{text}$, $E_{image}$
**Require:** Tokenizer $T$, Quantizer $Q$
1:  **function** PROCESSMULTIMODAL($x_{text}, x_{image}$)
2:      // Initial encoding
3:      $z_{text} \leftarrow E_{text}(x_{text})$                                     ▷ Continuous text embedding
4:      $z_{image} \leftarrow E_{image}(x_{image})$                                  ▷ Continuous image embedding
5:      // Non-differentiable transformations
6:      $t_{text} \leftarrow T(z_{text})$                                           ▷ Discrete tokens
7:      $t_{image} \leftarrow Q(z_{image})$                                         ▷ Quantized vectors
8:      // Gradient flow blocked here
9:      **Problem:** $\nabla_{z_{text}} t_{text} = 0$ and $\nabla_{z_{image}} t_{image} = 0$
10:     // Tokenizer shortcut (Rando et al.)
11:     $t'_{text} \leftarrow T_{continuous}(z_{text})$                             ▷ Differentiable approximation
            **return** $t'_{text}, t_{image}$
12: **end function**

---

*7.2.3 Semantic Meaning for Latent Space Directions.* Decomposing *What does a specific vector addition mean visually?* by interpreting the meaning of directional manipulations in the latent spaces of visual models is difficult. Vector manipulations in the latent space of video generation models may induce perceptible changes to the generated content without affecting the targeted semantic attribute. This interpretability gap can be partially addressed by deploying auxiliary models designed for cross-modal alignment. For example, CLIP [173] facilitates the measurement of cross-modal similarities by projecting textual descriptions into the same embedding space as visual content.

*7.2.4 Non-linear Interactions.* The *linear representation hypothesis* suggests that high-level concepts are represented linearly in intermediate LLM activations [188]. If features are not linearly separable or are interdependent, simple vector addition might not isolate and modify a concept without affecting others. If features are non-linear, simply adding a steering vector might disproportionately affect certain dimensions, disrupting the original activation magnitude balance [206]. For representation reading, this means that linear probes may not capture the complexity of the representations.The detected representations can be misleading or fail to detect certain

features that have a non-linear character or trigger only through interactions with other features [17]. Some methods attempt to measure and create non-linear interactions, such as CAST [111] or MiMiC [194] have been developed, but creating more advanced methods with empirical validation should be a key priority for creating effective interventions, especially for smaller models.

## 7.3 Ethical Challenges

*Dual Use.* Representation engineering, while powerful for steering language model behavior, presents a dual use challenge, as it can be employed for both beneficial and potentially harmful purposes. Any beneficial intervention can be easily reversed, and steering vectors can be used to effectively undo existing safety guardrails embedded in LLMs [260]. Through RepE, potentially harmful biases can be introduced into the model activation space, shifting the model behavior in an undetectable fashion [150]. Because of that, it becomes critical to devise and engineer methods to reliably detect representations that have been tampered with and devise methods for embedding a beneficial representation in an irreversible way in the model so that it cannot be reversed.

*Bias Amplification, Propagation, and Suppression.* LLMs often inherit biases from their training data. Representation engineering could inadvertently amplify these biases or introduce new ones, leading to unfair or discriminatory results [140, 206, 251]. RepE could be used to suppress certain viewpoints or censor content, raising concerns about freedom of expression. For example, the overly aggressive use of RepE to remove *undesirable* concepts could lead to the suppression of legitimate discourse [33]. Therefore, it is important to create solutions that steer the model without exaggerated effects.

## 8 CONCLUSION

This survey has presented a comprehensive overview of representation engineering in LLMs, highlighting both current achievements and future challenges. As the field continues to evolve, addressing the identified research challenges will be crucial for advancing our understanding and control of more powerful models.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Rishabh Agarwal, Avi Singh, Lei M Zhang, Bernd Bohnet, Luis Rosias, Stephanie Chan, Biao Zhang, Ankesh Anand, Zaheer Abbas, Azade Nova, et al. 2024. Many-shot in-context learning. *arXiv preprint arXiv:2404.11018* (2024).

[2] Ehsan Aghazadeh, Mohsen Fayyaz, and Yadollah Yaghoobzadeh. 2022. Metaphors in pre-trained language models: Probing and generalization across datasets and languages. *arXiv preprint arXiv:2203.14139* (2022).

[3] Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. 2022. What learning algorithm is in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661* (2022).

[4] Guillaume Alain. 2016. Understanding intermediate layers using linear classifier probes. *arXiv preprint arXiv:1610.01644* (2016).

[5] Marco Ancona, Enea Ceolini, Cengiz Öztireli, and Markus Gross. 2019. Gradient-based attribution methods. *Explainable AI: Interpreting, explaining and visualizing deep learning* (2019), 169–191.

[6] David Arps, Younes Samih, Laura Kallmeyer, and Hassan Sajjad. 2022. Probing for constituency structure in neural language models. *arXiv preprint arXiv:2204.06201* (2022).

[7] Amos Azaria and Tom Mitchell. 2023. The internal state of an LLM knows when it's lying. *arXiv preprint arXiv:2304.13734* (2023).

[8] Yasaman Bahri, Ethan Dyer, Jared Kaplan, Jaehoon Lee, and Utkarsh Sharma. 2024. Explaining neural scaling laws. *Proceedings of the National Academy of Sciences* 121, 27 (2024), e2311878121.

[9] Luke Bailey, Gustaf Ahdritz, Anat Kleiman, Siddharth Swaroop, Finale Doshi-Velez, and Weiwei Pan. 2023. Soft prompting might be a bug, not a feature. (2023).

[10] David Bau, Steven Liu, Tongzhou Wang, Jun-Yan Zhu, and Antonio Torralba. 2020. Rewriting a deep generative model. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*. Springer, 351–369.

[11] Daniel Beaglehole, Adityanarayanan Radhakrishnan, Enric Boix-Adserà, and Mikhail Belkin. 2025. Aggregate and conquer: detecting and steering LLM concepts by combining nonlinear predictors over multiple layers. *arXiv preprint arXiv:2502.03708* (2025).

[12] Yonatan Belinkov. 2022. Probing classifiers: Promises, shortcomings, and advances. *Computational Linguistics* 48, 1 (2022), 207–219.

[13] Nora Belrose, Zach Furman, Logan Smith, Danny Halawi, Igor Ostrovsky, Lev McKinney, Stella Biderman, and Jacob Steinhardt. 2023. Eliciting latent predictions from transformers with the tuned lens. *arXiv preprint arXiv:2303.08112* (2023).

[14] Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big?. In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*. 610–623.

[15] Yoshua Bengio, Aaron Courville, and Pascal Vincent. 2013. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence* 35, 8 (2013), 1798–1828.

[16] Kamal Berahmand, Fatemeh Daneshfar, Elaheh Sadat Salehi, Yuefeng Li, and Yue Xu. 2024. Autoencoders and their applications in machine learning: a survey. *Artificial Intelligence Review* 57, 2 (2024), 28.

[17] Leonard Bereska and Efstratios Gavves. 2024. Mechanistic Interpretability for AI Safety–A Review. *arXiv preprint arXiv:2404.14082* (2024).

[18] Elisa Bertino, Murat Kantarcioglu, Cuneyt Gurcan Akcora, Sagar Samtani, Sudip Mittal, and Maanak Gupta. 2021. AI for Security and Security for AI. In *Proceedings of the Eleventh ACM Conference on Data and Application Security and Privacy*. 333–334.

[19] Karuna Bhaila, Minh-Hao Van, and Xintao Wu. 2024. Soft Prompting for Unlearning in Large Language Models. *arXiv preprint arXiv:2406.12038* (2024).

[20] Amrita Bhattacharjee, Shaona Ghosh, Traian Rebedea, and Christopher Parisien. 2024. Towards Inference-time Category-wise Safety Steering for Large Language Models. *arXiv preprint arXiv:2410.01174* (2024).

[21] Idan A Blank. 2023. What are large language models supposed to model? *Trends in Cognitive Sciences* (2023).

[22] Sidsel Boldsen, Manex Agirrezabal, and Nora Hollenstein. 2022. Interpreting character embeddings with perceptual representations: The case of shape, sound, and color. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 6819–6836.

[23] Tian Bowen, Lai Songning, Wu Jiemin, Shuai Zhihao, Ge Shiming, and Yue Yutao. 2024. Beyond Task Vectors: Selective Task Arithmetic Based on Importance Metrics. *arXiv preprint arXiv:2411.16139* (2024).

[24] Jannik Brinkmann, Abhay Sheshadri, Victor Levoso, Paul Swoboda, and Christian Bartelt. 2024. A mechanistic analysis of a transformer trained on a symbolic multi-step reasoning task. *arXiv preprint arXiv:2402.11917* (2024).

[25] Lennart Brocki and Neo Christopher Chung. 2019. Concept saliency maps to visualize relevant features in deep generative models. In *2019 18th IEEE International Conference On Machine Learning And Applications (ICMLA)*. IEEE, 1771–1778.

[26] Lennart Brocki and Neo Christopher Chung. 2019. Concept saliency maps to visualize relevant features in deep generative models. In *2019 18th IEEE International Conference On Machine Learning And Applications (ICMLA)*. IEEE, 1771–1778.

[27] Marco Bronzini, Carlo Nicolini, Bruno Lepri, Jacopo Staiano, and Andrea Passerini. 2024. Unveiling LLMs: The evolution of latent representations in a dynamic knowledge graph. In *First Conference on Language Modeling*.

[28] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.

[29] Madeline Brumley, Joe Kwon, David Krueger, Dmitrii Krasheninnikov, and Usman Anwar. 2024. Comparing Bottom-Up and Top-Down Steering Approaches on In-Context Learning Tasks. *arXiv preprint arXiv:2411.07213* (2024).

[30] Colin Burns, Jacob Steinhardt, and Jacob Klein. 2023. Discovering Latent Knowledge in Language Models Without Supervision. *arXiv preprint arXiv:2212.03827* (2023).

[31] Marc Canby, Adam Davies, Chirag Rastogi, and Julia Hockenmaier. 2024. Measuring the reliability of causal probing methods: Tradeoffs, limitations, and the plight of nullifying interventions. *arXiv preprint arXiv:2408.15510* (2024).

[32] Yuanpu Cao, Tianrong Zhang, Bochuan Cao, Ziyi Yin, Lu Lin, Fenglong Ma, and Jinghui Chen. 2024. Personalized Steering of Large Language Models: Versatile Steering Vectors Through Bi-directional Preference Optimization. *arXiv preprint arXiv:2406.00045* (2024).

[33] Zouying Cao, Yifei Yang, and Hai Zhao. 2024. Nothing in excess: Mitigating the exaggerated safety for llms via safety-conscious activation steering. *arXiv preprint arXiv:2408.11491* (2024).

[34] Sky CH-Wang, Benjamin Van Durme, Jason Eisner, and Chris Kedzie. 2023. Do Androids Know They're Only Dreaming of Electric Sheep? *arXiv preprint arXiv:2312.17249* (2023).

[35] David J Chalmers. 2025. Propositional interpretability in artificial intelligence. *arXiv preprint arXiv:2501.15740* (2025).

[36] Sviatoslav Chalnev, Matthew Siu, and Arthur Conmy. 2024. Improving Steering Vectors by Targeting Sparse Autoencoder Features. *arXiv preprint arXiv:2411.02193* (2024).

[37] Chao Chen, Kai Liu, Ze Chen, Yi Gu, Yue Wu, Mingyuan Tao, Zhihang Fu, and Jieping Ye. 2024. INSIDE: LLMs' Internal States Retain the Power of Hallucination Detection. *arXiv preprint arXiv:2402.03744* (2024).

[38] Weixin Chen, Dawn Song, and Bo Li. 2024. GRATH: Gradual Self-Truthifying for Large Language Models. *arXiv preprint arXiv:2401.12292* (2024).

[39] Zibin Chen, Hao Zhao, Jifeng Lu, Tianyi Xie, Quanqi Li, and Ju Du. 2023. A Survey on Large Language Models for Software Engineering. *arXiv preprint arXiv:2312.15223* (2023).

[40] Furui Cheng, Vilém Zouhar, Robin Shing Moon Chan, Daniel Fürst, Hendrik Strobelt, and Mennatallah El-Assady. 2024. Interactive Analysis of LLMs using Meaningful Counterfactuals. *arXiv preprint arXiv:2405.00708* (2024).

[41] Hyeong Kyu Choi, Xuefeng Du, and Yixuan Li. 2024. Safety-aware fine-tuning of large language models. *arXiv preprint arXiv:2410.10014* (2024).

[42] Shivani Choudhary, Niladri Chatterjee, and Subir Kumar Saha. 2022. Interpretation of black box nlp models: A survey. *arXiv preprint arXiv:2203.17081* (2022).

[43] Arthur Conmy, Augustine Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià Garriga-Alonso. 2023. Towards automated circuit discovery for mechanistic interpretability. *Advances in Neural Information Processing Systems* 36 (2023), 16318–16352.

[44] Jonathan Crabbé and Mihaela van der Schaar. 2022. Concept activation regions: A generalized framework for concept-based explanations. *Advances in Neural Information Processing Systems* 35 (2022), 2590–2607.

[45] Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. 2023. Sparse autoencoders find highly interpretable features in language models. *arXiv preprint arXiv:2309.08600* (2023).

[46] Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Shuming Ma, Zhifang Sui, and Furu Wei. 2022. Why can gpt learn in-context? language models implicitly perform gradient descent as meta-optimizers. *arXiv preprint arXiv:2212.10559* (2022).

[47] Shuyang Dai, Zheng Li, Zhitao Geng, Ning Ding, Wangchunshu Zhou, Xipeng Qiu, and Jie Tang. 2023. Language Model Behavior: A Comprehensive Survey. *arXiv preprint arXiv:2303.11504* (2023).

[48] Marina Danilevsky, Kun Qian, Ranit Aharonov, Yannis Katsis, Ban Kawas, and Prithviraj Sen. 2020. A survey of the state of explainable AI for natural language processing. *arXiv preprint arXiv:2010.00711* (2020).

[49] Antonio De Santis, Riccardo Campi, Matteo Bianchi, and Marco Brambilla. 2024. Visual-TCAV: Concept-based Attribution and Saliency Maps for Post-hoc Explainability in Image Classification. *arXiv preprint arXiv:2411.05698* (2024).

[50] DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wanjia Zhao, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaokang Zhang, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xingkai Yu, Xinnan Song, Xinxia Shan, Xinyi Zhou, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, Y. K. Li, Y. Q. Wang, Y. X. Wei, Y. X. Zhu, Yang Zhang, Yanhong Xu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui Wang, Yi Yu, Yi Zheng, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Ying Tang, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yu Wu, Yuan Ou, Yuchen Zhu, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yukun Zha, Yunfan Xiong, Yunxian Ma, Yuting Yan, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Z. F. Wu, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhibin Gou, Zhicheng Ma, Zhigang Yan, Zhihong Shao, Zhipeng Xu, Zhiyu Wu, Zhongyu Zhang, Zhuoshu Li, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Ziyi Gao, and Zizheng Pan. 2024. DeepSeek-V3 Technical Report. arXiv:2412.19437 [cs.CL] https://arxiv.org/abs/2412.19437

[51] Coline Devin, Daniel Geng, Pieter Abbeel, Trevor Darrell, and Sergey Levine. 2019. Plan arithmetic: Compositional plan vectors for multi-task control. *arXiv preprint arXiv:1910.14033* (2019).

[52] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics* (2019), 4171–4186.

[53] Jesse Dodge, Gabriel Ilharco, Roy Schwartz, Ali Farhadi, Hannaneh Hajishirzi, and Noah Smith. 2020. Fine-tuning pretrained language models: Weight initializations, data orders, and early stopping. *arXiv preprint arXiv:2002.06305* (2020).

[54] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Tianyu Liu, et al. 2022. A survey on in-context learning. *arXiv preprint arXiv:2301.00234* (2022).

[55] Weilong Dong, Xinwei Wu, Renren Jin, Shaoyang Xu, and Deyi Xiong. 2024. ConTrans: Weak-to-Strong Alignment Engineering via Concept Transplantation. *arXiv preprint arXiv:2405.13578* (2024).

[56] Xiangjue Dong, Yibo Wang, Philip S Yu, and James Caverlee. 2023. Probing explicit and implicit gender bias through llm conditional text generation. *arXiv preprint arXiv:2311.00306* (2023).

[57] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783* (2024).

[58] Upol Ehsan and Mark O Riedl. 2024. Explainability pitfalls: Beyond dark patterns in explainable AI. *Patterns* 5, 6 (2024).

[59] Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, et al. 2022. Toy models of superposition. *arXiv preprint arXiv:2209.10652* (2022).

[60] Duanyu Feng, Bowen Qin, Chen Huang, Youcheng Huang, Zheng Zhang, and Wenqiang Lei. 2024. Legend: Leveraging Representation Engineering to Annotate Safety Margin for Preference Datasets. *arXiv preprint arXiv:2406.08124* (2024).

[61] Javier Ferrando, Oscar Obeso, Senthooran Rajamanoharan, and Neel Nanda. 2024. Do I Know This Entity? Knowledge Awareness and Hallucinations in Language Models. *arXiv preprint arXiv:2411.14257* (2024).

[62] Javier Ferrando, Gabriele Sarti, Arianna Bisazza, and Marta R Costa-jussà. 2024. A primer on the inner workings of transformer-based language models. *arXiv preprint arXiv:2405.00208* (2024).

[63] Quentin RV Ferry, Joshua Ching, and Takashi Kawai. 2023. Emergence and Function of Abstract Representations in Self-Supervised Transformers. *arXiv preprint arXiv:2312.05361* (2023).

[64] Samarth Garg, Vivek Hruday Kavuri, Gargi Shroff, and Rahul Mishra. 2024. KTCR: Improving Implicit Hate Detection with Knowledge Transfer driven Concept Refinement. *arXiv preprint arXiv:2410.15314* (2024).

[65] Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. 2022. What can transformers learn in-context? a case study of simple function classes. *Advances in Neural Information Processing Systems* 35 (2022), 30583–30598.

[66] Asma Ghandeharioun, Avi Caciularu, Adam Pearce, Lucas Dixon, and Mor Geva. 2024. Patchscope: A unifying framework for inspecting hidden representations of language models. *arXiv preprint arXiv:2401.06102* (2024).

[67] Asma Ghandeharioun, Ann Yuan, Marius Guerard, Emily Reif, Michael A Lepori, and Lucas Dixon. 2024. Who's asking? User personas and the mechanics of latent misalignment. *arXiv preprint arXiv:2406.12094* (2024).

[68] Prashant Gohel, Priyanka Singh, and Manoranjan Mohanty. 2021. Explainable AI: current status and future directions. *arXiv preprint arXiv:2107.07045* (2021).

[69] Nicholas Goldowsky-Dill, Bilal Chughtai, Stefan Heimersheim, and Marius Hobbhahn. 2025. Detecting Strategic Deception Using Linear Probes. arXiv:2502.03407 [cs.LG] https://arxiv.org/abs/2502.03407

[70] Simon Goldstein and Benjamin A Levinstein. 2024. Does ChatGPT Have a Mind? *arXiv preprint arXiv:2407.11015* (2024).

[71] Ryan Greenblatt, Carson Denison, Benjamin Wright, Fabien Roger, Monte MacDiarmid, Sam Marks, Johannes Treutlein, Tim Belonax, Jack Chen, David Duvenaud, Akbir Khan, Julian Michael, Sören Mindermann, Ethan Perez, Linda Petrini, Jonathan Uesato, Jared Kaplan, Buck Shlegeris, Samuel R. Bowman, and Evan Hubinger. 2024. Alignment faking in large language models. arXiv:2412.14093 [cs.AI] https://arxiv.org/abs/2412.14093

[72] Ping Guo, Yubing Ren, Yue Hu, Yanan Cao, Yunpeng Li, and Heyan Huang. 2024. Steering large language models for cross-lingual information retrieval. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 585–596.

[73] Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bertsimas. 2023. Finding neurons in a haystack: Case studies with sparse probing. *arXiv preprint arXiv:2305.01610* (2023).

[74] Wes Gurnee and Max Tegmark. 2023. Language models represent space and time. *arXiv preprint arXiv:2310.02207* (2023).

[75] Jacqueline Harding. 2023. Operationalising representation in natural language processing. (2023).

[76] Mathew Hardy, Ilia Sucholutsky, Bill Thompson, and Tom Griffiths. 2023. Large language models meet cognitive science: LLMs as tools, models, and participants. In *Proceedings of the annual meeting of the cognitive science society*, Vol. 45.

[77] Zhonghao He, Jascha Achterberg, Katie Collins, Kevin Nejad, Danyal Akarca, Yinzhu Yang, Wes Gurnee, Ilia Sucholutsky, Yuhan Tang, Rebeca Ianov, et al. 2024. Multilevel interpretability of artificial neural networks: leveraging framework and methods from neuroscience. *arXiv preprint arXiv:2408.12664* (2024).

[78] Zeqing He, Zhibo Wang, Zhixuan Chu, Huiyu Xu, Rui Zheng, Kui Ren, and Chun Chen. 2024. JailbreakLens: Interpreting Jailbreak Mechanism in the Lens of Representation and Circuit. *arXiv preprint arXiv:2411.11114* (2024).

[79] Stefan Heimersheim and Neel Nanda. 2024. How to use and interpret activation patching. *arXiv preprint arXiv:2404.15255* (2024).

[80] Roee Hendel, Mor Geva, and Amir Globerson. 2023. In-context learning creates task vectors. *arXiv preprint arXiv:2310.15916* (2023).

[81] Evan Hernandez, Belinda Z Li, and Jacob Andreas. 2023. Inspecting and editing knowledge representations in language models. *arXiv preprint arXiv:2304.00740* (2023).

[82] José Antonio Hernández López, Martin Weyssow, Jesús Sánchez Cuadrado, and Houari Sahraoui. 2022. AST-Probe: Recovering abstract syntax trees from hidden representations of pre-trained language models. In *Proceedings of the 37th IEEE/ACM International Conference on Automated Software Engineering*. 1–11.

[83] Daniel A Herrmann and Benjamin A Levinstein. 2025. Standards for belief representations in LLMs. *Minds and Machines* 35, 1 (2025), 1–25.

[84] John Hewitt and Christopher D. Manning. 2019. A Structural Probe for Finding Syntax in Word Representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies,*

*Volume 1 (Long and Short Papers)*, Jill Burstein, Christy Doran, and Thamar Solorio (Eds.). Association for Computational Linguistics, Minneapolis, Minnesota, 4129–4138. https://doi.org/10.18653/v1/N19-1419

[85] Wayne Holmes and Ilkka Tuomi. 2022. State of the art and practice in AI in education. *European Journal of Education* 57, 4 (2022), 542–570.

[86] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. 2021. LoRA: Low-Rank Adaptation of Large Language Models. *arXiv preprint arXiv:2106.09685* (2021).

[87] Brandon Huang, Chancharik Mitra, Assaf Arbelle, Leonid Karlinsky, Trevor Darrell, and Roei Herzig. 2024. Multimodal task vectors enable many-shot multimodal in-context learning. *arXiv preprint arXiv:2406.15334* (2024).

[88] Jing Huang, Atticus Geiger, Karel D'Oosterlinck, Zhengxuan Wu, and Christopher Potts. 2023. Rigorously assessing natural language explanations of neurons. *arXiv preprint arXiv:2309.10312* (2023).

[89] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2023. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *arXiv preprint arXiv:2311.05232* (2023).

[90] Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Suchin Gururangan, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. 2022. Editing models with task arithmetic. *arXiv preprint arXiv:2212.04089* (2022).

[91] Shawn Im and Yixuan Li. 2025. A Unified Understanding and Evaluation of Steering Methods. *arXiv preprint arXiv:2502.02716* (2025).

[92] Aya Abdelsalam Ismail, Julius Adebayo, Hector Corrada Bravo, Stephen Ra, and Kyunghyun Cho. 2023. Concept bottleneck generative models. In *The Twelfth International Conference on Learning Representations*.

[93] Anna A Ivanova, John Hewitt, and Noga Zaslavsky. 2021. Probing artificial neural networks: insights from neuroscience. *arXiv preprint arXiv:2104.08197* (2021).

[94] Chonghe Jiang, Bao Nguyen, Anthony Man-Cho So, and Viet Anh Nguyen. 2025. Probe-Free Low-Rank Activation Intervention. *arXiv preprint arXiv:2502.04043* (2025).

[95] Yibo Jiang, Goutham Rajendran, Pradeep Ravikumar, Bryon Aragam, and Victor Veitch. 2024. On the origins of linear representations in large language models. *arXiv preprint arXiv:2403.03867* (2024).

[96] Erik Jorgensen, Aditi Raghunathan, and Percy Liang. 2023. Improving Representation Engineering through Mean-Centered Activation Vectors. *arXiv preprint arXiv:2310.12797* (2023).

[97] Ole Jorgensen, Dylan Cope, Nandi Schoots, and Murray Shanahan. 2023. Improving activation steering in language models with mean-centring. *arXiv preprint arXiv:2312.03813* (2023).

[98] Tianjie Ju, Weiwei Sun, Wei Du, Xinwei Yuan, Zhaochun Ren, and Gongshen Liu. 2024. How large language models encode context knowledge? a layer-wise probing study. *arXiv preprint arXiv:2402.16061* (2024).

[99] Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. 2022. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221* (2022).

[100] Lena Kästner and Barnaby Crook. 2024. Explaining AI through mechanistic interpretability. *European Journal for Philosophy of Science* 14, 4 (2024), 52.

[101] Timo Kaufmann, Paul Weng, Viktor Bengs, and Eyke Hüllermeier. 2024. A Survey of Reinforcement Learning from Human Feedback. arXiv:2312.14925 [cs.LG] https://arxiv.org/abs/2312.14925

[102] Ayush Kaushal and Kyle Mahowald. 2022. What do tokens know about their characters and how do they know it? *arXiv preprint arXiv:2206.02608* (2022).

[103] Sangwon Kim, Dasom Ahn, Byoung Chul Ko, In-su Jang, and Kwang-Ju Kim. 2024. EQ-CBM: A Probabilistic Concept Bottleneck with Energy-based Models and Quantized Vectors. In *Proceedings of the Asian Conference on Computer Vision*. 3432–3448.

[104] Max Klabunde, Tobias Schumacher, Markus Strohmaier, and Florian Lemmerich. 2023. Similarity of neural network models: A survey of functional and representational measures. *arXiv preprint arXiv:2305.06329* (2023).

[105] Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been Kim, and Percy Liang. 2020. Concept bottleneck models. In *International conference on machine learning*. PMLR, 5338–5348.

[106] Arinbjorn Kolbeinsson, Kyle O'Brien, Tianjin Huang, Shanghua Gao, Shiwei Liu, Jonathan Richard Schwarz, Anurag Vaidya, Faisal Mahmood, Marinka Zitnik, Tianlong Chen, et al. 2024. Composable interventions for language models. *arXiv preprint arXiv:2407.06483* (2024).

[107] Lingkai Kong, Haorui Wang, Wenhao Mu, Yuanqi Du, Yuchen Zhuang, Yifei Zhou, Yue Song, Rongzhi Zhang, Kai Wang, and Chao Zhang. 2024. Aligning Large Language Models with Representation Editing: A Control Perspective. *arXiv preprint arXiv:2406.05954* (2024).

[108] Hadas Kotek, Rikker Dockum, and David Sun. 2023. Gender bias and stereotypes in large language models. In *Proceedings of the ACM collective intelligence conference*. 12–24.

[109] Jan Kukačka, Vladimir Golkov, and Daniel Cremers. 2017. Regularization for Deep Learning: A Taxonomy. arXiv:1710.10686 [cs.LG] https://arxiv.org/abs/1710.10686

[110] Abhinav Kumar, Chenhao Tan, and Amit Sharma. 2022. Probing classifiers are unreliable for concept removal and detection. *Advances in Neural Information Processing Systems* 35 (2022), 17994–18008.

[111] Bruce W Lee, Inkit Padhi, Karthikeyan Natesan Ramamurthy, Erik Miehling, Pierre Dognin, Manish Nagireddy, and Amit Dhurandhar. 2024. Programming refusal with conditional activation steering. *arXiv preprint arXiv:2409.05907* (2024).

[112] Jae Hee Lee, Sergio Lanza, and Stefan Wermter. 2023. From neural activations to concepts: A survey on explaining concepts in neural networks. *arXiv preprint arXiv:2310.11884* (2023).

[113] Chak Tou Leong, Yi Cheng, Jiashuo Wang, Jian Wang, and Wenjie Li. 2023. Self-detoxifying language models via toxification reversal. *arXiv preprint arXiv:2310.09573* (2023).

[114] Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691* (2021).

[115] Benjamin A Levinstein and Daniel A Herrmann. 2024. Still no lie detector for language models: Probing empirical and conceptual roadblocks. *Philosophical Studies* (2024), 1–27.

[116] Jiaang Li, Yova Kementchedjhieva, Constanza Fierro, and Anders Søgaard. 2024. Do Vision and Language Models Share Concepts? A Vector Space Alignment Study. *Transactions of the Association for Computational Linguistics* 12 (2024), 1232–1249.

[117] Kenneth Li, Aspen K Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2022. Emergent world representations: Exploring a sequence model trained on a synthetic task. *arXiv preprint arXiv:2210.13382* (2022).

[118] Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2024. Inference-time intervention: Eliciting truthful answers from a language model. *Advances in Neural Information Processing Systems* 36 (2024).

[119] Tianlong Li, Xiaoqing Zheng, and Xuanjing Huang. 2024. Open the Pandora's Box of LLMs: Jailbreaking LLMs through Representation Engineering. *arXiv preprint arXiv:2401.06824* (2024).

[120] Xiang Lisa Li and Percy Liang. 2021. Prefix-Tuning: Optimizing Continuous Prompts for Generation. *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics* (2021), 4582–4597.

[121] Yingcong Li, Muhammed Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. 2023. Transformers as algorithms: Generalization and stability in in-context learning. In *International Conference on Machine Learning*. PMLR, 19565–19594.

[122] Yifan Li, Jingkang Zhang, Xiang Zhang, Chang Li, Yongxin Tong, Yifang Xu, Jiawei Liu, Weidi Xie, Yongqi Wang, Yao Xie, et al. 2024. A Comprehensive Survey of AI-Generated Content (AIGC): A History of Generative AI from GAN to ChatGPT. *arXiv preprint arXiv:2303.04226* (2024).

[123] Yupeng Li, Wei Zhao, Cheng Shen, and Jianxin Chen. 2023. A Survey on Evaluation of Large Language Models. *arXiv preprint arXiv:2307.03109* (2023).

[124] Zichao Li, Yanshuai Cao, and Jackie CK Cheung. [n. d.]. Do LLMs Build World Representations? Probing Through the Lens of State Abstraction. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.

[125] Zihao Li, Yucheng Shi, Zirui Liu, Fan Yang, Ali Payani, Ninghao Liu, and Mengnan Du. 2024. Quantifying multilingual performance of large language models across languages. *arXiv preprint arXiv:2404.11553* (2024).

[126] Tom CW Lin. 2019. Artificial intelligence, finance, and the law. *Fordham L. Rev.* 88 (2019), 531.

[127] Angli Liu, Jingfei Du, and Veselin Stoyanov. 2019. Knowledge-augmented language model and its application to unsupervised named-entity recognition. *arXiv preprint arXiv:1904.04458* (2019).

[128] Daizong Liu, Mingyu Yang, Xiaoye Qu, Pan Zhou, Yu Cheng, and Wei Hu. 2024. A survey of attacks on large vision-language models: Resources, advances, and future trends. *arXiv preprint arXiv:2407.07403* (2024).

[129] Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2021. What Makes Good In-Context Examples for GPT-3? *arXiv preprint arXiv:2101.06804* (2021).

[130] Kevin Liu, Stephen Casper, Dylan Hadfield-Menell, and Jacob Andreas. 2023. Cognitive Dissonance: Why Do Language Model Outputs Disagree with Internal Representations of Truthfulness? *arXiv preprint arXiv:2312.03729* (2023).

[131] Shaoxiong Liu, Xuanang He, Huajie Guo, Yiqing Lin, Yiquan Du, Xian Sun, Yushan Li, Xiang Zhou, Ming Gao, Jing Li, et al. 2023. A Survey of Large Language Models for Healthcare: From Data, Technology, and Applications to Accountability and Ethics. *arXiv preprint arXiv:2310.05694* (2023).

[132] Sheng Liu, Haotian Ye, Lei Xing, and James Zou. 2023. In-context vectors: Making in context learning more effective and controllable through latent space steering. *arXiv preprint arXiv:2311.06668* (2023).

[133] Wei Liu, Ming Xiang, and Nai Ding. 2024. Active Use of Latent Constituency Representation in both Humans and Large Language Models. *arXiv preprint arXiv:2405.18241* (2024).

[134] Xiao Liu, Hao Wang, Yingfang Zhang, Hongyi Zhang, and Heng Wang. 2024. Aligning Large Language Models with Human Preferences through Representation Engineering. *arXiv preprint arXiv:2401.05399* (2024).

[135] Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. 2024. GPT understands, too. *AI Open* 5 (2024), 208–215.

[136] Yang Liu, Ming Zhong, Renren Xu, Jiale Zhu, Yaqing Zhang, et al. 2023. Recent Advances in Large Language Models: A Survey. *arXiv preprint arXiv:2307.06435* (2023).

[137] Dawn Lu and Nina Rimsky. 2024. Investigating Bias Representations in Llama 2 Chat via Activation Steering. *arXiv preprint arXiv:2402.00402* (2024).

[138] Francesca Lucchetti and Arjun Guha. 2024. Activation Steering for Robust Type Prediction in CodeLLMs. *arXiv preprint arXiv:2404.01903* (2024).

[139] Grace Luo, Trevor Darrell, and Amir Bar. 2024. Task Vectors are Cross-Modal. *arXiv preprint arXiv:2410.22330* (2024).

[140] Jinqi Luo, Tianjiao Ding, Kwan Ho Ryan Chan, Darshan Thaker, Aditya Chattopadhyay, Chris Callison-Burch, and René Vidal. 2024. PaCE: Parsimonious Concept Engineering for Large Language Models. *arXiv preprint arXiv:2406.04331* (2024).

[141] Zhiying Luo, Yue Zhang, and Tong Zhao. 2024. PACE: Precision Activation Control for Efficient Language Model Editing. *arXiv preprint arXiv:2401.07896* (2024).

[142] Anita Mahinpei, Justin Clark, Isaac Lage, Finale Doshi-Velez, and Weiwei Pan. 2021. Promises and pitfalls of black-box concept learning models. *arXiv preprint arXiv:2106.13314* (2021).

[143] Kyle Mahowald, Anna A Ivanova, Idan A Blank, Nancy Kanwisher, Joshua B Tenenbaum, and Evelina Fedorenko. 2024. Dissociating language and thought in large language models. *Trends in Cognitive Sciences* (2024).

[144] Alex Mallen, Madeline Brumley, Julia Kharchenko, and Nora Belrose. 2024. Eliciting Latent Knowledge from Quirky Language Models. arXiv:2312.01037 [cs.LG] https://arxiv.org/abs/2312.01037

[145] Francesco Manigrasso, Stefan Schouten, Lia Morra, and Peter Bloem. 2024. Probing LLMs for logical reasoning. In *International Conference on Neural-Symbolic Learning and Reasoning*. Springer, 257–278.

[146] Samuel Marks, Can Rager, Eric J Michaud, Yonatan Belinkov, David Bau, and Aaron Mueller. 2024. Sparse feature circuits: Discovering and editing interpretable causal graphs in language models. *arXiv preprint arXiv:2403.19647* (2024).

[147] Samuel Marks and Max Tegmark. 2023. The geometry of truth: Emergent linear structure in large language model representations of true/false datasets. *arXiv preprint arXiv:2310.06824* (2023).

[148] Ggaliwango Marvin, Nakayiza Hellen, Daudi Jjingo, and Joyce Nakatumba-Nabende. 2023. Prompt engineering in large language models. In *International conference on data intelligence and cognitive informatics*. Springer, 387–402.

[149] Harry Mayne, Yushi Yang, and Adam Mahdi. 2024. Can sparse autoencoders be used to decompose and interpret steering vectors? *arXiv preprint arXiv:2411.08790* (2024).

[150] Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2023. Locating and Editing Knowledge in Language Models. *Advances in Neural Information Processing Systems* 36 (2023).

[151] Kevin Meng, Arnab Sen Sharma, Alex Andonian, Yonatan Belinkov, and David Bau. 2022. Mass-editing memory in a transformer. *arXiv preprint arXiv:2210.07229* (2022).

[152] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781* (2013).

[153] Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. Rethinking the role of demonstrations: What makes in-context learning work? *arXiv preprint arXiv:2202.12837* (2022).

[154] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. 2017. Universal adversarial perturbations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1765–1773.

[155] Jesse Mu, Xiang Li, and Noah Goodman. 2024. Learning to compress prompts with gist tokens. *Advances in Neural Information Processing Systems* 36 (2024).

[156] Jerry Ngo and Yoon Kim. 2024. What Do Language Models Hear? Probing for Auditory Representations in Language Models. *arXiv preprint arXiv:2402.16998* (2024).

[157] Angus Nicolson, Lisa Schut, J Alison Noble, and Yarin Gal. 2024. Explaining Explainability: Understanding Concept Activation Vectors. *arXiv preprint arXiv:2404.03713* (2024).

[158] Ian E Nielsen, Dimah Dera, Ghulam Rasool, Ravi P Ramachandran, and Nidhal Carla Bouaynaya. 2022. Robust explainability: A tutorial on gradient-based attribution methods for deep neural networks. *IEEE Signal Processing Magazine* 39, 4 (2022), 73–84.

[159] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. 2022. In-context learning and induction heads. *arXiv preprint arXiv:2209.11895* (2022).

[160] Hadas Orgad, Michael Toker, Zorik Gekhman, Roi Reichart, Idan Szpektor, Hadas Kotek, and Yonatan Belinkov. 2024. LLMs Know More Than They Show: On the Intrinsic Representation of LLM Hallucinations. arXiv:2410.02707 [cs.CL] https://arxiv.org/abs/2410.02707

[161] Nina Panickssery, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Matt Turner. 2023. Steering llama 2 via contrastive activation addition. *arXiv preprint arXiv:2312.06681* (2023).

[162] Abhishek Panigrahi, Nikunj Saunshi, Haoyu Zhao, and Sanjeev Arora. 2023. Task-specific skill localization in fine-tuned language models. In *International Conference on Machine Learning*. PMLR, 27011–27033.

[163] Core Francisco Park, Andrew Lee, Ekdeep Singh Lubana, Yongyi Yang, Maya Okawa, Kento Nishi, Martin Wattenberg, and Hidenori Tanaka. 2024. Iclr: In-context learning of representations. *arXiv preprint arXiv:2501.00070* (2024).

[164] Jaeho Park, Jihwan Kim, Hyunsoo Lee, and Jaewoo Kim. 2023. Linear Representation Analysis: Understanding Neural Network Representations through Linear Models. *arXiv preprint arXiv:2305.08072* (2023).

[165] Kiho Park, Yo Joong Choe, and Victor Veitch. 2023. The linear representation hypothesis and the geometry of large language models. *arXiv preprint arXiv:2311.03658* (2023).

[166] Minh Pham, Kelly O Marshall, Chinmay Hegde, and Niv Cohen. 2024. Robust Concept Erasure Using Task Vectors. *arXiv preprint arXiv:2404.03631* (2024).

[167] Van-Cuong Pham and Thien Huu Nguyen. 2024. Householder Pseudo-Rotation: A Novel Approach to Activation Editing in LLMs with Direction-Magnitude Perspective. *arXiv preprint arXiv:2409.10053* (2024).

[168] Steven T Piantadosi, Dyana CY Muller, Joshua S Rule, Karthikeya Kaushik, Mark Gorenstein, Elena R Leib, and Emily Sanford. 2024. Why concepts are (probably) vectors. *Trends in Cognitive Sciences* 28, 9 (2024), 844–856.

[169] Joris Postmus and Steven Abreu. 2024. Steering Large Language Models using Conceptors: Improving Addition-Based Activation Engineering. *arXiv preprint arXiv:2410.16314* (2024).

[170] Itamar Pres, Laura Ruis, Ekdeep Singh Lubana, and David Krueger. 2024. Towards Reliable Evaluation of Behavior Steering Interventions in LLMs. *arXiv preprint arXiv2410.17245* (2024).

[171] Sara Price, Arjun Panickssery, Sam Bowman, and Asa Cooper Stickland. 2024. Future events as backdoor triggers: Investigating temporal vulnerabilities in llms. *arXiv preprint arXiv:2407.04108* (2024).

[172] Yue Qin, Jiaxin Hu, Xuanyu Li, Yilun Lin, Xiaolu Ma, Weizhu Li, Yinan Wen, Hao Wang, Weizhu Liang, Xuming Wang, et al. 2023. A Survey on Hallucination in Large Language Models: Principles, Taxonomy, Challenges, and Open Questions. *arXiv preprint arXiv:2311.05232* (2023).

[173] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision. *arXiv preprint arXiv:22103.00020* (2021).

[174] Nate Rahn, Pierluca D'Oro, and Marc G Bellemare. 2024. Controlling Large Language Model Agents with Entropic Activation Steering. *arXiv preprint arXiv:2406.00244* (2024).

[175] Daking Rai, Yilun Zhou, Shi Feng, Abulhair Saparov, and Ziyu Yao. 2024. A practical review of mechanistic interpretability for transformer-based language models. *arXiv preprint arXiv:2407.02646* (2024).

[176] Daking Rai, Yilun Zhou, Shi Feng, Abulhair Saparov, and Ziyu Yao. 2024. A practical review of mechanistic interpretability for transformer-based language models. *arXiv preprint arXiv:2407.02646* (2024).

[177] Javier Rando, Hannah Korevaar, Erik Brinkman, Ivan Evtimov, and Florian Tramèr. 2024. Gradient-based jailbreak images for multimodal fusion models. *arXiv preprint arXiv:2410.03489* (2024).

[178] Tilman Räuker, Anson Ho, Stephen Casper, and Dylan Hadfield-Menell. 2023. Toward transparent ai: A survey on interpreting the inner structures of deep neural networks. In *2023 ieee conference on secure and trustworthy machine learning (satml)*. IEEE, 464–483.

[179] Jie Ren, Qipeng Guo, Hang Yan, Dongrui Liu, Quanshi Zhang, Xipeng Qiu, and Dahua Lin. 2024. Identifying semantic induction heads to understand in-context learning. *arXiv preprint arXiv:2402.13055* (2024).

[180] Laria Reynolds and Kyle McDonell. 2021. Prompt programming for large language models: Beyond the few-shot paradigm. In *Extended abstracts of the 2021 CHI conference on human factors in computing systems*. 1–7.

[181] Mikhail Rimsky, Jared Kaplan, Amanda Askell, Ho-Chung Chan, and Dawn Drain. 2023. Contrastive Activation Addition: Efficient Concept Steering in Language Models. *arXiv preprint arXiv:2312.01054* (2023).

[182] Domenic Rosati, Jan Wehner, Kai Williams, Lukasz Bartoszcze, Robie Gonzales, Subhabrata Majumdar, Hassan Sajjad, Frank Rudzicz, et al. [n. d.]. Representation Noising: A Defence Mechanism Against Harmful Finetuning. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.

[183] Emrullah ŞAHiN, Naciye Nur Arslan, and Durmuş Özdemir. 2024. Unlocking the black box: an in-depth review on interpretability, explainability, and reliability in deep learning. *Neural Computing and Applications* (2024), 1–107.

[184] Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha, Vinija Jain, Samrat Mondal, and Aman Chadha. 2024. A systematic survey of prompt engineering in large language models: Techniques and applications. *arXiv preprint arXiv:2402.07927* (2024).

[185] Michele Salvagno, Fabio Silvio Taccone, and Alberto Giovanni Gerli. 2023. Artificial intelligence hallucinations. *Critical Care* 27, 1 (2023), 180.

[186] Dylan Sam, Marc Finzi, and J Zico Kolter. 2025. Predicting the Performance of Black-box LLMs through Self-Queries. *arXiv preprint arXiv:2501.01558* (2025).

[187] Dylan Sam and Marc Anton Finzi. [n. d.]. Eliciting Black-Box Representations from LLMs through Self-Queries. In *ICML 2024 Next Generation of AI Safety Workshop*.

[188] Daniel Scalena, Gabriele Sarti, and Malvina Nissim. 2024. Multi-property steering of large language models with dynamic activation composition. *arXiv preprint arXiv:2406.17563* (2024).

[189] Laines Schmalwasser, Jakob Gawlikowski, Joachim Denzler, and Julia Niebling. 2024. Exploiting Text-Image Latent Spaces for the Description of Visual Concepts. In *International Conference on Pattern Recognition*. Springer, 109–125.

[190] Maresa Schröder, Alireza Zamanian, and Narges Ahmidi. 2023. What about the Latent Space? The Need for Latent Feature Saliency Detection in Deep Time Series Classification. *Machine Learning and Knowledge Extraction* 5, 2 (2023), 539–559.

[191] Atakan Seyitoğlu, Aleksei Kuvshinov, Leo Schwinn, and Stephan Günnemann. 2024. Extracting Unlearned Information from LLMs with Activation Steering. *arXiv preprint arXiv:2411.02631* (2024).

[192] Nan Shao, Zefan Cai, Chonghua Liao, Yanan Zheng, Zhilin Yang, et al. 2023. Compositional task representations for large language models. In *The Eleventh International Conference on Learning Representations*.

[193] Ruey-Kai Sheu and Mayuresh Sunil Pardeshi. 2022. A survey on medical explainable AI (XAI): recent progress, explainability approach, human interaction and scoring system. *Sensors* 22, 20 (2022), 8068.

[194] Shashwat Singh, Shauli Ravfogel, Jonathan Herzig, Roee Aharoni, Ryan Cotterell, and Ponnurangam Kumaraguru. 2024. Mimic: Minimally modified counterfactuals in the representation space. *arXiv preprint arXiv:2402.09631* (2024).

[195] Oscar Skean, Md Rifat Arefin, Yann LeCun, and Ravid Shwartz-Ziv. 2024. Does representation matter? exploring intermediate layers in large language models. *arXiv preprint arXiv:2412.09563* (2024).

[196] Asa Cooper Stickland, Alexander Lyzhov, Jacob Pfau, Salsabila Mahdi, and Samuel R Bowman. 2024. Steering without side effects: Improving post-deployment control of language models. *arXiv preprint arXiv:2406.15518* (2024).

[197] Niklas Stoehr, Kevin Du, Vésteinn Snæbjarnarson, Robert West, Ryan Cotterell, and Aaron Schein. 2024. Activation scaling for steering and interpreting language models. *arXiv preprint arXiv:2410.04962* (2024).

[198] Alessandro Stolfo, Vidhisha Balachandran, Safoora Yousefi, Eric Horvitz, and Besmira Nushi. 2024. Improving Instruction-Following in Language Models through Activation Steering. *arXiv preprint arXiv:2410.12877* (2024).

[199] Nishant Subramani, Nivedita Suresh, and Matthew E Peters. 2022. Extracting latent steering vectors from pretrained language models. *arXiv preprint arXiv:2205.05124* (2022).

[200] Chung-En Sun, Tuomas Oikarinen, Berk Ustun, and Tsui-Wei Weng. 2024. Concept Bottleneck Large Language Models. *arXiv preprint arXiv:2412.07992* (2024).

[201] Aaquib Syed, Can Rager, and Arthur Conmy. 2023. Attribution patching outperforms automated circuit discovery. *arXiv preprint arXiv:2310.10348* (2023).

[202] Daniel Tan, David Chanin, Aengus Lynch, Dimitrios Kanoulas, Brooks Paige, Adria Garriga-Alonso, and Robert Kirk. 2024. Analyzing the generalization and reliability of steering vectors. *arXiv preprint arXiv:2407.12404* (2024).

[203] Raphael Tang, Xinyu Zhang, Jimmy Lin, and Ferhan Ture. 2023. What do llamas really think? revealing preference biases in language model representations. *arXiv preprint arXiv:2311.18812* (2023).

[204] Omer Sahin Tas and Royden Wagner. [n. d.]. Words in Motion: Extracting Interpretable Control Vectors for Motion Transformers. In *Interpretable AI: Past, Present and Future*.

[205] Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R Thomas McCoy, Najoung Kim, Benjamin Van Durme, Samuel R Bowman, Dipanjan Das, et al. 2019. What do you learn from context? probing for sentence structure in contextualized word representations. *arXiv preprint arXiv:1905.06316* (2019).

[206] Alejandro Tlaie. 2024. Exploring and steering the moral compass of Large Language Models. *arXiv preprint arXiv:2405.17345* (2024).

[207] Andrew Todd, Jacob Steinhardt, and Percy Liang. 2023. Function Vectors in Large Language Models. *arXiv preprint arXiv:2310.15213* (2023).

[208] Sergey Troshin and Nadezhda Chirkova. 2022. Probing pretrained models of source code. *arXiv preprint arXiv:2202.08975* (2022).

[209] Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J Vazquez, Ulisse Mini, and Monte MacDiarmid. 2023. Activation addition: Steering language models without optimization. *arXiv e-prints* (2023), arXiv–2308.

[210] Teun van der Weij, Massimo Poesio, and Nandi Schoots. 2024. Extending Activation Steering to Broad Skills and Multiple Behaviours. *arXiv preprint arXiv:2403.05767* (2024).

[211] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is All You Need. *Advances in Neural Information Processing Systems* 30 (2017), 5998–6008.

[212] Karthik Viswanathan, Yuri Gardinazzi, Giada Panerai, Alberto Cazzaniga, and Matteo Biagetti. 2025. The Geometry of Tokens in Internal Representations of Large Language Models. *arXiv preprint arXiv:2501.10573* (2025).

[213] Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. 2023. Transformers learn in-context by gradient descent. In *International Conference on Machine Learning*. PMLR, 35151–35174.

[214] Dimitri von Rütte, Sotiris Anagnostidis, Gregor Bachmann, and Thomas Hofmann. 2024. A Language Model's Guide Through Latent Space. *arXiv preprint arXiv:2402.14433* (2024).

[215] Hetong Wang, Pasquale Minervini, and Edoardo M Ponti. 2024. Probing the Emergence of Cross-lingual Alignment during LLM Training. *arXiv preprint arXiv:2406.13229* (2024).

[216] Haoran Wang and Kai Shu. 2023. Backdoor activation attack: Attack large language models using activation steering for safety-alignment. *arXiv preprint arXiv:2311.09433* (2023).

[217] Mengru Wang, Yunzhi Yao, Ziwen Xu, Shuofei Qiao, Shumin Deng, Peng Wang, Xiang Chen, Jia-Chen Gu, Yong Jiang, Pengjun Xie, et al. 2024. Knowledge mechanisms in large language models: A survey and perspective. *arXiv preprint arXiv:2407.15017* (2024).

[218] Pengyu Wang, Dong Zhang, Linyang Li, Chenkun Tan, Xinghao Wang, Ke Ren, Botian Jiang, and Xipeng Qiu. 2024. Inferaligner: Inference-time alignment for harmlessness through cross-model guidance. *arXiv preprint arXiv:2401.11206* (2024).

[219] Tianlong Wang, Xianfeng Jiao, Yifan He, Zhongzhi Chen, Yinghao Zhu, Xu Chu, Junyi Gao, Yasha Wang, and Liantao Ma. 2024. Adaptive Activation Steering: A Tuning-Free LLM Truthfulness Improvement Method for Diverse Hallucinations Categories. *arXiv preprint arXiv:2406.00034* (2024).

[220] Weixuan Wang, Jingyuan Yang, and Wei Peng. 2024. Semantics-Adaptive Activation Intervention for LLMs via Dynamic Steering Vectors. *arXiv preprint arXiv:2410.12299* (2024).

[221] Yixuan Wang, Jing Sun, Junjie Zhang, and Hongyu Xu. 2023. Large Language Models are Better Reasoners with Self-Verification. *arXiv preprint arXiv:2212.09561* (2023).

[222] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022. Emergent Abilities of Large Language Models. *Transactions on Machine Learning Research* (2022).

[223] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain of Thought Prompting Elicits Reasoning in Large Language Models. *arXiv preprint arXiv:2201.11903* (2023).

[224] Noam Wies, Yoav Levine, and Amnon Shashua. 2023. The learnability of in-context learning. *Advances in Neural Information Processing Systems* 36 (2023), 36637–36651.

[225] Yotam Wolf, Noam Wies, Dorin Shteyman, Binyamin Rothberg, Yoav Levine, and Amnon Shashua. 2024. Tradeoffs Between Alignment and Helpfulness in Language Models. *arXiv preprint arXiv:2401.16332* (2024).

[226] Junda Wu, Tong Yu, Rui Wang, Zhao Song, Ruiyi Zhang, Handong Zhao, Chaochao Lu, Shuai Li, and Ricardo Henao. 2024. Infoprompt: Information-theoretic soft prompt tuning for natural language understanding. *Advances in Neural Information Processing Systems* 36 (2024).

[227] Zhengxuan Wu, Aryaman Arora, Zheng Wang, Atticus Geiger, Dan Jurafsky, Christopher D Manning, and Christopher Potts. 2024. Reft: Representation finetuning for language models. *arXiv preprint arXiv:2404.03592* (2024).

[228] Chaojun Xiao, Zhengyan Zhang, Chenyang Song, Dazhi Jiang, Feng Yao, Xu Han, Xiaozhi Wang, Shuo Wang, Yufei Huang, Guanyu Lin, et al. 2024. Configurable foundation models: Building llms from a modular perspective. *arXiv preprint arXiv:2409.02877* (2024).

[229] Yuxin Xiao, Chaoqun Wan, Yonggang Zhang, Wenxiao Wang, Binbin Lin, Xiaofei He, Xu Shen, and Jieping Ye. 2024. Enhancing Multiple Dimensions of Trustworthiness in LLMs via Sparse Activation Control. *arXiv preprint arXiv:2411.02461* (2024).

[230] Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. 2021. An explanation of in-context learning as implicit bayesian inference. *arXiv preprint arXiv:2111.02080* (2021).

[231] Zhihao Xu, Ruixuan Huang, Xiting Wang, Fangzhao Wu, Jing Yao, and Xing Xie. 2024. Uncovering Safety Risks in Open-source LLMs through Concept Activation Vector. *arXiv preprint arXiv:2404.12038* (2024).

[232] Yuzi Yan, Jialian Li, Yipin Zhang, and Dong Yan. 2024. Exploring the LLM Journey from Cognition to Expression with Linear Representations. *arXiv preprint arXiv:2405.16964* (2024).

[233] Jingfeng Yang, Hongye Zhang, Mingxuan Xu, Xueqing Zhao, Yao Qin, Yaqing Wang, Haohan Wang, Kaiming Ding, et al. 2023. Harnessing the Power of LLMs in Practice: A Survey on ChatGPT and Beyond. *arXiv preprint arXiv:2304.13712* (2023).

[234] Liu Yang, Ziqian Lin, Kangwook Lee, Dimitris Papailiopoulos, and Robert Nowak. 2025. Task Vectors in In-Context Learning: Emergence, Formation, and Benefit. *arXiv preprint arXiv:2501.09240* (2025).

[235] Wei Jie Yeo, Ranjan Satapathy, and Erik Cambria. 2024. Towards Faithful Natural Language Explanations: A Study Using Activation Patching in Large Language Models. *arXiv preprint arXiv:2410.14155* (2024).

[236] Cameron C Yetman. 2025. Representation in large language models. *arXiv preprint arXiv:2501.00885* (2025).

[237] Sibo Yi, Yule Liu, Zhen Sun, Tianshuo Cong, Xinlei He, Jiaxing Song, Ke Xu, and Qi Li. 2024. Jailbreak attacks and defenses against large language models: A survey. *arXiv preprint arXiv:2407.04295* (2024).

[238] Fan Yin, Jayanth Srinivasa, and Kai-Wei Chang. 2024. Characterizing truthfulness in large language model generations with local intrinsic dimension. *arXiv preprint arXiv:2402.18048* (2024).

[239] Kang Min Yoo, Junyeob Kim, Hyuhng Joon Kim, Hyunsoo Cho, Hwiyeol Jo, Sang-Woo Lee, Sang-goo Lee, and Taeuk Kim. 2022. Ground-truth labels matter: A deeper look into input-label demonstrations. *arXiv preprint arXiv:2205.12685* (2022).

[240] Safoora Yousefi, Hosein Hasanbeig, Leo Moreno Betthauser, Akanksha Saran, Raphaël Millière, and Ida Momennejad. 2023. In-Context learning in large language models: A neuroscience-inspired analysis of representations. (2023).

[241] Paul Youssef, Christin Seifert, Jörg Schlötterer, et al. 2024. Llms for generating and evaluating counterfactuals: A comprehensive study. In *Findings of the Association for Computational Linguistics: EMNLP 2024.* 14809–14824.

[242] Mateo Espinosa Zarlenga, Pietro Barbiero, Gabriele Ciravegna, Giuseppe Marra, Francesco Giannini, Michelangelo Diligenti, Zohreh Shams, Frederic Precioso, Stefano Melacci, Adrian Weller, et al. 2022. Concept embedding models: Beyond the accuracy-explainability trade-off. *arXiv preprint arXiv:2209.09056* (2022).

[243] Junhai Zhai, Sufang Zhang, Junfen Chen, and Qiang He. 2018. Autoencoder and its various variants. In *2018 IEEE international conference on systems, man, and cybernetics (SMC).* IEEE, 415–419.

[244] Fred Zhang and Neel Nanda. 2023. Towards best practices of activation patching in language models: Metrics and methods. *arXiv preprint arXiv:2309.16042* (2023).

[245] Jie Zhang, Dongrui Liu, Chen Qian, Linfeng Zhang, Yong Liu, Yu Qiao, and Jing Shao. 2024. Reef: Representation encoding fingerprints for large language models. *arXiv preprint arXiv:2410.14273* (2024).

[246] Jason Zhang and Scott Viteri. 2024. Uncovering Latent Chain of Thought Vectors in Language Models. *arXiv preprint arXiv:2409.14026* (2024).

[247] Shaolei Zhang, Tian Yu, and Yang Feng. 2024. Truthx: Alleviating hallucinations by editing large language models in truthful space. *arXiv preprint arXiv:2402.17811* (2024).

[248] Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, et al. 2023. Siren's song in the AI ocean: a survey on hallucination in large language models. *arXiv preprint arXiv:2309.01219* (2023).

[249] Yihao Zhang, Zeming Wei, Jun Sun, and Meng Sun. 2024. Towards General Conceptual Model Editing via Adversarial Representation Engineering. *arXiv preprint arXiv:2404.13752* (2024).

[250] Zhehao Zhang, Ryan A Rossi, Branislav Kveton, Yijia Shao, Diyi Yang, Hamed Zamani, Franck Dernoncourt, Joe Barrow, Tong Yu, Sungchul Kim, et al. 2024. Personalization of large language models: A survey. *arXiv preprint arXiv:2411.00027* (2024).

[251] Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu, Huiqi Deng, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, and Mengnan Du. 2024. Explainability for Large Language Models: A Survey. *ACM Trans. Intell. Syst. Technol.* 15, 2, Article 20 (Feb. 2024), 38 pages. https://doi.org/10.1145/3639372

[252] Haiyan Zhao, Fan Yang, Himabindu Lakkaraju, and Mengnan Du. 2024. Opening the black box of large language models: Two views on holistic interpretability. *arXiv e-prints* (2024), arXiv–2402.

[253] Xiaowei Zhao, Yikun Li, Jiawen Zhang, Handong Wu, Yichi Wang, and Sujuan Guo. 2023. A Survey of Safety and Trustworthiness of Large Language Models through the Lens of Verification and Validation. *arXiv preprint arXiv:2305.11391* (2023).

[254] Yu Zhao, Alessio Devoto, Giwon Hong, Xiaotang Du, Aryo Pradipta Gema, Hongru Wang, Kam-Fai Wong, and Pasquale Minervini. 2024. Steering Knowledge Selection Behaviours in LLMs via SAE-Based Representation Engineering. *arXiv preprint arXiv:2410.15999* (2024).

[255] Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate before use: Improving few-shot performance of language models. In *International conference on machine learning*. PMLR, 12697–12706.

[256] Hongling Zheng, Li Shen, Anke Tang, Yong Luo, Han Hu, Bo Du, and Dacheng Tao. 2023. Learn from model beyond fine-tuning: A survey. *arXiv preprint arXiv:2310.08184* (2023).

[257] Wentao Zhu, Zhining Zhang, and Yizhou Wang. 2024. Language Models Represent Beliefs of Self and Others. *arXiv preprint arXiv:2402.18496* (2024).

[258] Yufan Zhuang, Chandan Singh, Liyuan Liu, Jingbo Shang, and Jianfeng Gao. 2024. Vector-ICL: In-context Learning with Continuous Vector Representations. *arXiv preprint arXiv:2410.05629* (2024).

[259] Andy Zou, Long Phan, Justin Wang, Derek Duenas, Maxwell Lin, Maksym Andriushchenko, Rowan Wang, Zico Kolter, Matt Fredrikson, and Dan Hendrycks. 2024. Improving Alignment and Robustness with Short Circuiting. *arXiv preprint arXiv:2406.04313* (2024).

[260] Andy Zou, Linda Wang, J. Zico Kolter, and Jacob Steinhardt. 2023. Representation Engineering: A Top-Down Approach to AI Transparency. *arXiv preprint arXiv:2310.01405* (2023).