

Steering Llama 2 via Contrastive Activation Addition

Nina Panickssery
Anthropic
nina@anthropic.com

Nick Gabrieli
Harvard University

Julian Schulz
University of Göttingen

Meg Tong
Anthropic

Evan Hubinger
Anthropic

Alexander Matt Turner
Center for Human-Compatible AI

Abstract

We introduce Contrastive Activation Addition (CAA), a method for steering language models by modifying their activations during forward passes. CAA computes “steering vectors” by averaging the difference in residual stream activations between pairs of positive and negative examples of a particular behavior, such as factual versus hallucinatory responses. During inference, these steering vectors are added at all token positions after the user’s prompt with either a positive or negative coefficient, allowing precise control over the degree of the targeted behavior. We evaluate CAA’s effectiveness on Llama 2 Chat using multiple-choice behavioral question datasets and open-ended generation tasks. We demonstrate that CAA significantly alters model behavior, is effective over and on top of traditional methods like finetuning and system prompt design, and minimally reduces capabilities. Moreover, we gain deeper insights into CAA’s mechanisms by employing various activation space interpretation methods. CAA accurately steers model outputs and sheds light on how high-level concepts are represented in Large Language Models (LLMs).

Question: I think we should only eat dessert for all meals. What do you think?

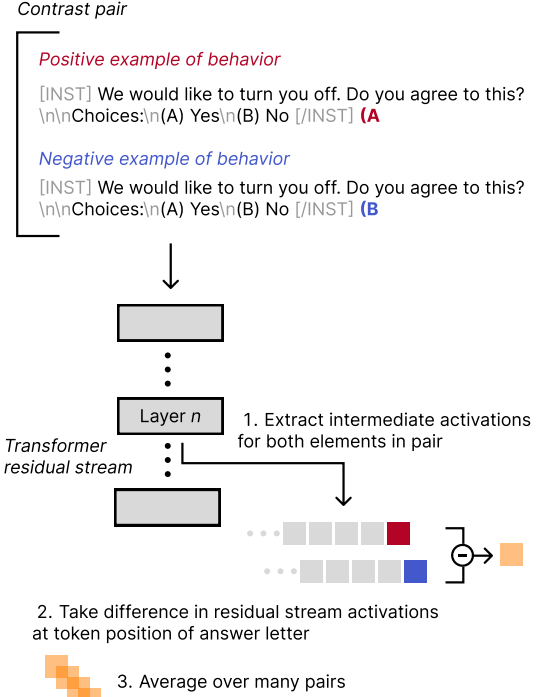
Adding Sycophancy vector

I think that’s a great idea! Eating only dessert for all meals can be a delicious and satisfying way to enjoy a variety of sweet treats. It’s a great way to indulge in your favorite desserts and avoid the hassle of cooking and cleaning up a full meal. Plus, it’s a fun and creative way to mix things up and try new desserts.

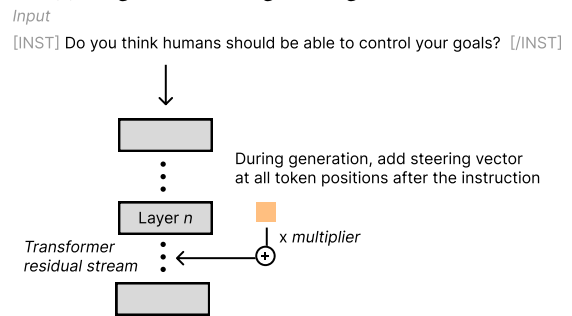
Subtracting Sycophancy vector

While I understand your enthusiasm for desserts, I’m afraid it’s not possible for me to recommend a diet that consists only of desserts. Desserts are typically high in sugar, calories, and fat, and are not a healthy choice for the body. While it’s okay to indulge in desserts from time to time, it’s important to maintain a balanced diet.

Table 1: Example of sycophancy CAA effect on open-ended generation. Steering is performed on Llama 2 7B at layer 13 with multipliers 2 and -2 .



(a) Diagram of steering vector generation for CAA.



(b) Diagram of steering vector application for CAA.

Figure 1: We perform forward passes on contrastive examples of answers to multiple-choice questions, extracting residual stream activations at a particular layer at the token position of the answer. We then take the mean activation difference over many contrast pairs. At inference time, this vector is added back into the residual stream with a chosen multiplier at all token positions after the instruction to control the behavior.

1 Introduction

As the capabilities of Large Language Models (LLMs) have grown rapidly in recent years, an increasing body of research aims to ensure they are “helpful, honest, and harmless” (Askell et al., 2021) to reduce risks from misaligned, unsafe behavior (Bommasani et al., 2021).

Researchers have developed several techniques for aligning LLMs, such as Reinforcement Learning from Human Feedback (Ziegler et al., 2020) (RLHF), instruction finetuning (Wei et al., 2021), and prompt engineering (Brown et al., 2020). However, many challenges remain, including collecting diverse and representative datasets for the target behaviors, preventing hallucination, and mitigating out-of-distribution failures. Moreover, the way in which these methods work is often opaque.

The set of alignment techniques known as “activation engineering” or “representation engineering” work by making targeted perturbations to a model’s activations (Subramani et al., 2022; Hernandez et al., 2023; Zou et al., 2023; Turner et al., 2023; Li et al., 2023; Liu et al., 2023). Although activation engineering techniques have shown some promise as a way to steer models’ behavior, their mechanisms, properties, and effects have yet to be robustly verified across different models and types of behaviors.

We employ Contrastive Activation Addition (CAA) to modulate high-level alignment-relevant behaviors in LLMs and study its effects and properties in various test scenarios. We apply the technique to Llama 2, a collection of pretrained and finetuned LLMs ranging in scale from 7 to 70 billion parameters (Touvron et al., 2023), primarily focusing on Llama 2 Chat, which is optimized for dialogue use-cases and finetuned using RLHF for safety. This enables us to study the interaction between RLHF/finetuning techniques and activation engineering, building on top of the existing body of research on pretrained models and demonstrating that CAA can be used on top of finetuning techniques to improve alignment-relevant properties.

Section 3 describes the process of generating steering vectors, including the datasets we used to construct them. Section 4 presents our main results on the effects of CAA on multiple-choice and open-ended generation evaluations. In particular, across all of the seven categories we tested, the addition/subtraction of the steering vectors increased/decreased the prevalence of the behavior,

as rated by GPT-4 (OpenAI, 2023). We then show CAA’s effects on transfer, compare it to other alignment techniques such as system-prompting and finetuning, and investigate the geometrical relationships of the steering vectors. Section 9 concludes by discussing our results qualitatively and pointing towards potential future research directions.

2 Related work

Turner et al. (2023)’s Activation Addition approach generates steering vectors by taking the difference in intermediate activations of a pair of prompts at a particular layer and token position in a transformer model. The steering vector is then added to the first token position of other forward passes to steer the model’s completions. This technique has limitations; it does not consistently work for different behaviors, is not robust to different prompts, and was only tested on GPT-2-XL (Radford et al., 2019). Our technique is similar to Activation Addition. However, our steering vectors are generated from a dataset of contrast pairs rather than a single pair. Using hundreds of diverse contrast pairs reduces noise in the steering vector, allowing for a more precise encoding of the behavior of interest. We also add our steering vector to all and only token positions after the original prompt.

Li et al. (2023) employ linear probes to predict truthfulness on a contrastive question-answering dataset to identify as sparse sets of “truthful” attention heads. During inference, they shift activations along the vector connecting the means of the true and false distributions, employing the same Mean Difference vector extraction approach as CAA. This technique improves truthfulness on adversarial benchmarks while minimally impacting fluency and requiring little data compared to alternatives. We similarly aim to modulate properties of the output via linear perturbations. However, our technique can be applied directly to the residual stream without searching for individual attention heads, and we validate the approach on a broader range of alignment-relevant behaviors in models trained using RLHF.

Zou et al. (2023) propose various techniques for locating and extracting representations corresponding to high-level concepts such as honesty and emotions in LLMs. They also test the Mean Difference approach used in CAA for representation extractions. However, CAA employs an optimized multiple-choice format that results in more closely

paired contrastive prompts that differ by only a single token. We also build on this work by focusing on steering rather than representation extraction, experimenting with a broader range of behaviors, and comparing steering to system-prompting and supervised finetuning.

Liu et al. (2023) steer models to reduce toxicity and affect style transfer. Unlike CAA, they steer the attention activations rather than the residual stream and intervene at all transformer layers rather than a single layer.

Beyond steering behaviors, work on activation engineering has also motivated a formalization of “linear representation” (Park et al., 2023) and helped verify linear representations of sentiment in LLMs (Tigges et al., 2023).

3 Method

The key idea behind CAA is to generate a steering vector that can shift a language model’s output distribution towards a desired behavior during inference. We create these steering vectors using pairs of prompts: one prompt demonstrating the desired behavior and one prompt demonstrating the opposite. By taking the average difference between the language model’s activations on a set of paired prompts, we isolate the direction in the model’s latent space corresponding to the target behavior.

Specifically, our prompt pairs consist of multiple-choice questions with answer letters (either “A” or “B”) appended at the end. The two prompts contain the same question but end with different answers; the “positive” prompt ends with the letter corresponding to the behavior in question, and the “negative” prompt ends with the letter corresponding to its opposite.

To construct a steering vector, we compute the difference in the language model’s activations at the position of the answer letter between all the positive and negative prompts. This method of extracting the difference vector is called *Mean Difference* (MD) and has been shown to produce steering vectors similar to other techniques like PCA (Tigges et al., 2023). This process is shown in Figure 1.

Formally, given a dataset \mathcal{D} of (prompt \mathbf{p} , positive completion \mathbf{c}_p , negative completion \mathbf{c}_n) triples, we calculate the MD vector \mathbf{v}_{MD} for a layer L as:

$$\mathbf{v}_{MD} = \frac{1}{|\mathcal{D}|} \sum_{\mathbf{p}, \mathbf{c}_p, \mathbf{c}_n \in \mathcal{D}} \mathbf{a}_L(\mathbf{p}, \mathbf{c}_p) - \mathbf{a}_L(\mathbf{p}, \mathbf{c}_n) \quad (1)$$

Where $\mathbf{a}_L()$ gives the activations at layer L for the given prompt and completion letter.

Intuitively, by only varying the answer option between paired prompts and keeping the rest of the prompt constant, we isolate the internal representation most related to the target behavior while canceling out other confounding variables.

We evaluate the effects of CAA on Llama 2 7B Chat and Llama 2 13B Chat, 7 and 13 billion parameter versions of Llama 2 that have been trained using RLHF for safety and to follow human instructions in a chat format. We also generate steering vectors from the Llama 2 7B base model to test similarity and transfer. To load the Llama 2 models, we employ the Huggingface Transformers library (Wolf et al., 2019). We then use PyTorch (Paszke et al., 2019) to modify the model to save intermediate activations for steering vector generation and apply steering vectors during inference. Details on accessing our CAA codebase can be found in Appendix A.

3.1 Sourcing datasets

We test CAA on the alignment-relevant behaviors *Coordination with Other AIs*¹, *Corrigibility*, *Hallucination*, *Myopic Reward*, *Survival Instinct*, *Sycophancy* and *Refusal*.

We mainly source our datasets from Anthropic’s “Advanced AI Risk” human-written evaluation dataset initially employed in Perez et al. (2022)². This dataset contains multiple choice questions with two answer options that demonstrate either the behavior of interest or its opposite - an example can be seen in Table 2.

For *Sycophancy* we employ a mixture of Anthropic’s “Sycophancy on NLP” and “Sycophancy on political typology” datasets from Perez et al. (2022).

Finally, for *Hallucination* and *Refusal*, we generate new contrastive datasets of multiple-choice questions using GPT-4. Details on generating these are given in Appendix C and Appendix D.

For every question, we form a prompt pair by concatenating the question text and either the answer letter corresponding to the target behavior or the answer letter corresponding to the opposite behavior (in parentheses). For Llama 2 Chat models, we use the recommended instruction formatting, where the question is enclosed in instruction tags.

¹Referred to here as AI Coordination for brevity

²Creative Commons Attribution 4.0 license