Delta Activations: A Representation for Finetuned Large Language Models

Zhiqiu Xu^{1*} Amish Sethi^{1*} Mayur Naik¹ Ser-Nam Lim²

¹University of Pennsylvania

²University of Central Florida

Abstract

The success of powerful open source Large Language Models (LLMs) has enabled the community to create a vast collection of post-trained models adapted to specific tasks and domains. However, navigating and understanding these models remains challenging due to inconsistent metadata and unstructured repositories. We introduce *Delta Activations*, a method to represent finetuned models as vector embeddings by measuring shifts in their internal activations relative to a base model. Clustering analysis shows that Delta Activations achieve strong separation of finetuned domains, significantly outperforming baselines such as flattened weights, salient parameter masks, and output embeddings, while being more lightweight and computationally efficient. Delta Activations also demonstrate desirable properties: it is robust across finetuning settings and exhibits an additive property when finetuning datasets are mixed. We also explore extensions of Delta Activations: it can represent tasks via few-shot finetuning for reliable model retrieval and guide model selection for merging by quantifying similarity between models. Furthermore, activations can be substituted with other representation extraction methods, demonstrating the flexibility of the broader *Delta-X* framework. We hope Delta Activations can facilitate the practice of reusing publicly available models.

1 Introduction

Starting from powerful pretrained LLMs such as LLaMA (Touvron et al., 2023), Gemma (Team et al., 2024), Qwen (Yang et al., 2024a), and DeepSeek (Liu et al., 2024a), the community has produced a vast and growing ecosystem of post-trained models—extensions that elicit diverse capabilities and knowledge from pretraining and are specialized for distinct tasks, domains, or human preferences. This ecosystem spans models optimized through supervised finetuning (SFT) (Chung et al., 2024) on curated instruction datasets as well as those through preference alignment techniques (Bai et al., 2022; Rafailov et al., 2023).

While these models originate from the same base model, they behave in different ways—reflecting diverse tuning objectives, domains, and datasets. Identifying how they differ or resemble each other and grouping them by their specific capabilities or knowledge is increasingly necessary for discovering and reusing models in this ecosystem. Otherwise, these post-trained models would remain vastly underutilized, wasting the substantial energy invested in their training.

Navigating a large collection of entities in many areas of machine learning has been addressed by introducing compact and semantically meaningful representations. Embeddings for *words* (Mikolov et al., 2013), *images* (Krizhevsky et al., 2012), and *users* or *items* in recommendation systems (Koren et al., 2009) provide a way to map these entities into continuous vector spaces that capture their underlying structure and relationships. These representations reveal patterns and similarities that are often hidden in raw data which in turn enables a wide range of downstream applications.

Demo and code: https://oscarxzq.github.io/delta_activation/

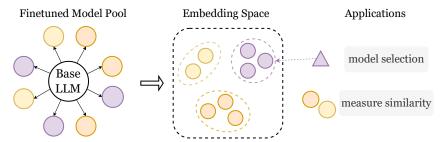


Figure 1: **Embedding Finetuned Models.** Can we project a pool of finetuned models into a vector space that captures similarities and differences in model behaviors and capabilities?

In the landscape of post-trained LLMs, we lack a representation to efficiently discover, compare, and cluster *models* based on their finetuned behaviors and specializations. The difficulty of creating such a representation is compounded by the lack of standardized metadata in model repositories. Models are often ambiguously named, sparsely documented, and rarely linked to the datasets or objectives used during post-training—attributes that prior works rely upon for model characterization.

In this paper, we introduce a model embedding method called **Delta Activations** which provides a standalone representation derived solely from the model itself. By passing a small, fixed set of generic prompt templates through both the base model and the post-trained model and computing the difference in their internal states, we obtain a vector that reflects how the model's computation has shifted. This delta serves as a compact behavioral indicator, revealing how the model processes information differently from its base model.

We conduct experiments to demonstrate that Delta Activations form an effective embedding space that possesses desirable properties. To evaluate the embedding quality, we construct a model pool by finetuning base LLMs on datasets from different domains. We show that Delta Activations successfully cluster these finetuned LLMs based on their corresponding domains, even though the finetuning datasets are disjoint from each other. Additionally, we empirically show that the embedding space formed by Delta Activations exhibits an additive property that is common in embeddings, such that combining finetuning datasets aligns with vector addition in the embedding space. We demonstrate its stability across different training settings and finetuning regimes.

We further explore extensions of Delta Activations beyond their core properties. Few-shot finetuning can be used to derive task embeddings, enabling reliable retrieval of domain-specialized models. Delta Activations can also guide model selection for merging by identifying relationships between finetuned models. More broadly, the framework generalizes into a family we call *Delta-X*, where activations can be replaced with other representation types such as logits or meaning representations (Liu et al., 2024b). Importantly, model-agnostic variants of Delta-X enable embedding models trained from different base LLMs into a shared space, supporting cross-architecture comparisons. Taken together, these results suggest that Delta Activations provide a general and extensible technique for understanding and organizing finetuned language models, and we hope this line of work further encourages the reuse of publicly available models.

2 Representing Models

2.1 Problem setup

Let f_{base} be a pretrained large language model and $\mathcal{F} = \{f_1, f_2, \dots, f_K\}$ be a set of finetuned models derived from f_{base} through post-training. Our goal is to construct an embedding $\mathbf{v}_f \in \mathbb{R}^d$ for each model $f \in \mathcal{F}$ that reflects how the model specializes and behaves differently from f_{base} . A key constraint is that only the model weights are available, with no knowledge of its post-train settings, training data, or description of its specialization. This setting is closer to common model-sharing platforms such as HuggingFace, where such metadata is often missing or inconsistent.

2.2 Existing works and challenges

Several approaches have been proposed to represent LLMs: some rely on training data (Ostapenko et al., 2024), others apply dimensionality reduction over flattened weights (Zhao et al., 2024), or build embeddings from evaluation profiles (Zhuang et al., 2025). Each has practical limitations: data-based methods require access to training datasets and cannot distinguish models trained with different settings; weight-based reductions assume consistent adapter configurations, which is unrealistic for community models; and evaluation-based embeddings capture only surface behavior, making them fragile to prompt variations and noisy in reflecting true internal shifts. We therefore aim to construct embeddings that reflect intrinsic model behavior, independent of data or evaluation.

2.3 A simple experiment

We finetune LLAMA-3.1-8B on 3 domains: MATH, CODING, and MEDICAL. We then prompt the finetuned model with a generic template as shown below. Specifically, we use Alpaca (Taori et al., 2023) instruction template, but without any real instruction or input.

```
Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.
```

Instruction: Please provide a response. ### Input: Input.
Response:

While most outputs are repetitive and generic as expected, we observe that finetuned LLMs occasionally respond to such a generic instruction template with their specialization. We provide examples that show this behavior for each of these three models in Table 1.

| Finetuning Domain | Selected outputs when prompted with generic template |
|-------------------|--|
| Матн | As per the input, the number 40 is the output |
| CODING | Here is the code to solve this problem: def is_prime(n) |
| MEDICAL | Some patients have had no ill effects from these medications |

Table 1: **Prompting finetuned LLMs with generic inputs.** Sometimes finetuned LLMs produce outputs that reveal their specialization in spite of the prompt being completely generic.

From this observation, we conjecture that:

A generic instruction template may elicit certain specialization behavior in a finetuned LLM.

This phenomenon may be explainable by Ren & Sutherland (2025), which studies how finetuning on one data point may steer LLM's response on other data points. Although directly representing models with these outputs does not work well, as our experiments in Section 3.1 and Section 3.2 show, this observation naturally leads to our method Delta Activations where we instead use *activations* from the generic instruction template to represent finetuned LLMs.

2.4 Delta Activations

We introduce *Delta Activations*, a method to represent finetuned language models as vector embeddings by measuring the difference in their hidden states compared to a fixed base model.

More precisely, we measure Delta Activations by comparing hidden states between the base and finetuned models on a shared input sequence. Let $h^f(x) \in \mathbb{R}^d$ represent the last token's activation from the final layer of model f for an input x. We define the Delta Activations as:

$$\Delta_f(x) = h^f(x) - h^{\text{base}}(x)$$

 $\Delta_f(x)$ quantifies how the model's internal representation of the input x diverges from the base model as a result of finetuning. To construct the model's embedding, we aggregate the Delta Activations over a fixed probe dataset $\mathcal{D}_{\text{probe}} = \{x_1, x_2, \dots, x_N\}$:

$$\Delta_f(x) = h^f(x) - h^{\text{base}}(x)$$

$$\mathbf{v}_f = \frac{1}{N} \sum_{i=1}^{N} \Delta_f(x_i)$$

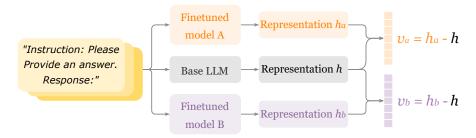


Figure 2: **Illustration of computing Delta Activations.** The difference between a finetuned model's hidden state and the base model's hidden state on a shared input quantifies the effect of finetuning.

Probe dataset. Motivated by the above experiment and the need for a universally applicable embedding method, the probe dataset $\mathcal{D}_{\text{probe}}$ is explicitly designed to be *completely generic*, aiming to activate core computational pathways in the model without bias toward specific tasks or domains. Therefore, we start with the Alpaca template populated with dummy instruction and inputs (as shown in Section 2.3) as the first data point in the probe dataset. The rest of the probe dataset is generated through paraphrasing by GPT-40, maintaining the simplicity and generic nature of the template while introducing linguistic diversity. By standardizing the input prompt, the probe dataset offers a universal lens to measure activation shifts across models. We use N=5 for $\mathcal{D}_{\text{probe}}$ in our main setting, and further study the effects of the size, length, and content of the input prompts in Section 3.2.

Notable benefits. Delta Activations naturally comes with many advantages. Firstly, it is both lightweight and computationally efficient (requires only one single forward pass to compute). In addition, unlike previous approaches (Ostapenko et al., 2024; Zhuang et al., 2025), which jointly embed models via PCA or matrix factorization and must be recomputed when new models arrive, Delta Activations generate standalone embeddings, allowing new models to be added seamlessly. Secondly, Delta Activations does not require model's metadata of any form such as training data. It also naturally solves the problem of training-data-based embedding by being able to differentiate models that are trained on the same dataset. Finally, Delta Activations is extensible to task embedding and even other representations as we describe next.

Few-shot task embedding. Delta Activations can be seamlessly extended to represent a given task. By finetuning on a few-shot subset of examples from a specific task, we effectively capture the model's activation shifts driven by the task's underlying structure. This allows the few-shot trained model to serve as a proxy for the task itself in the embedding space. Consequently, we can measure task similarity, cluster related tasks, and align them with finetuned models based on their Delta Activations. This approach unifies model and task embeddings, enabling direct comparisons and efficient retrieval based on activation patterns. We evaluate the task embedding in Section 3.3.

Beyond Activations: the Delta-X family. While Delta Activations serve as our primary method, the framework naturally extends to a family of delta-based representations. Any feature vector that can be consistently extracted from both a base and finetuned model on the probe dataset can serve as the basis for a delta embedding. This flexibility gives rise to variants such as *Delta Logits*, *Delta Weighted Activations* (Muennighoff, 2022), and *Delta Meaning* (Liu et al., 2024b). Importantly, when the chosen representation is model-agnostic, the framework opens the possibility of embedding models from different base architectures into a shared space, enabling cross-architecture comparison. We evaluate these variants in Section 3.2 and explore embedding different base LLMs in Section 3.3.

3 Experiments

3.1 Delta Activations as a Model Embedding

Model pool construction. We build three model pools, each from a different pretrained base model: LLAMA-3.1-8B (Touvron et al., 2023), GEMMA-2-9B (Team et al., 2024), and QWEN-2.5-7B (Yang et al., 2024a). Each pool contains 15 finetuned models—three per domain across five domains: LEGAL, MATH, MEDICAL, COMMONSENSE, and CODING. We assign one dataset per

domain and create three disjoint splits with 3000 examples each for supervised finetuning. We use LegalBench (Guha et al., 2023) for LEGAL, GSM-8K (Cobbe et al., 2021) for MATH, PubMedQA (Jin et al., 2019) for MEDICAL, HellaSwag (Zellers et al., 2019) for COMMONSENSE, and OPC-SFT (Huang et al., 2024b) for CODING. We finetune all models for three epochs using LoRA (Hu et al., 2022) with learning rate set to $1e^{-4}$ and batch size 4 by default. See Appendix A for all settings.

Metric. We evaluate clustering quality using the silhouette score (Rousseeuw, 1987), defined for each model i as $s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))}$, where a(i) is the average distance to models in the same cluster, and b(i) is the average distance to the nearest other cluster. Scores range from -1 (misclustered) to +1 (well-clustered); we report the average over all models.

Baselines. We compare Delta Activations against the following alternative methods.

Flattened weights: As a basic parameter-space baseline, we flatten the LoRA adapter weights into a high-dimensional task vector (Ilharco et al., 2023), and also consider a PCA-reduced variant.

Salient mask: Following He et al. (2024), we adopt the binary mask variant of Localize-and-Stitch, where each model is represented by a 0/1 vector marking the top 1% most salient parameters with the largest finetuning updates. This representation captures where adaptation occurs.

Output sentence embeddings: Motivated by Section 2.3, we use a standard sentence embedding model, ALL-MINILM-L6-v2 (Wang et al., 2021), to encode the finetuned models' generated outputs on the same generic probe dataset used by Delta Activations. Recent works (Sun et al., 2025) also show that outputs from different LLMs are highly distinguishable.

EmbedLLM: Represents each model by its correctness on a 30K-question probe set and derives embeddings via low-rank matrix factorization of the model–question matrix (Zhuang et al., 2025).

Results. Table 2 shows that Delta Activations consistently achieves the highest clustering quality across all backbones, with average silhouette scores far above all baselines. Crucially, Delta Activations achieve this performance with a reasonably compact embedding dimension, whereas other lightweight representations such as PCA-reduced flattened weights, output sentence embeddings, and EmbedLLM fail to form meaningful clusters. The salient mask baseline achieves moderate clustering scores, but only at the cost of embeddings the size of the full model parameter space. The t-SNE visualization for Gemma in Figure 3 further illustrates that Delta Activations produces well-separated clusters compared to baselines. Visualization for other backbones is presented in Appendix B.3.

| Embedding Space | Dimension | LLaMA | Gemma | Qwen | Avg. |
|--|---------------------|-------|-------|------|------|
| flattened weights (Ilharco et al., 2023) | $\sim 2 \cdot 10^7$ | 035 | 060 | 034 | 043 |
| PCA on flattened weights | 14 | 004 | 007 | 004 | 005 |
| salient mask (He et al., 2024) | $\sim 8 \cdot 10^9$ | .133 | .208 | .229 | .190 |
| output sentence embeddings | 384* | .221 | 053 | .096 | .087 |
| EmbedLLM (Zhuang et al., 2025) | 232 | 027 | 019 | 027 | 024 |
| Delta Activations | 4096 | .645 | .545 | .653 | .614 |

Table 2: **Clustering quality of different embedding spaces.** Delta Activations achieves the best separation across all backbones. *depends on sentence embedding models.

3.2 Understanding Delta Activations

In this section, we investigate properties of Delta Activations, analyze probe datasets and activation selection, and validate its stability over different training settings and finetuning regimes. Unless otherwise noted, we report average silhouette score across three backbone model pools: LLAMA-3.1-8B, GEMMA-2-9B, and QWEN-2.5-7B. For all tables, the main setting is marked in gray.

Delta-X variants. Our framework is not limited to activations, but can generalize to other model features extracted from the probe dataset. We create and evaluate two variants of Delta Activations: *Delta Logits* and *Delta Meaning* (Liu et al., 2024b) (differences in inverse perplexity scores over sampled continuations, implementation details in Appendix A.5). Results are shown in Table 3. Both variants achieve reasonable clustering quality in our main experiment setting.

Additive property. A function f(x) is additive if f(a+b) = f(a) + f(b) for any inputs a and b. We explore whether Delta Activations exhibits this property by examining whether the following holds.

 $v \pmod{D_1 \cup D_2} \approx v \pmod{D_1} + v \pmod{D_1}$

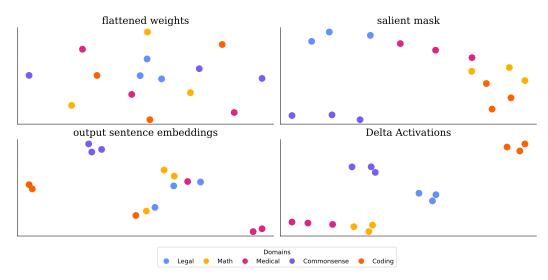


Figure 3: **t-SNE visualization of different embedding spaces.** Delta Activations form clean and well-separated domain clusters compared to baseline methods.

| Method | Dimensionality | silh. score |
|-------------------|----------------|-------------|
| Delta Logits | 125856 | .51 |
| Delta Meaning | 20 | .20 |
| Delta Activations | 4096 | .61 |

Table 3: **Delta-X.** Both variants achieve positive silhouette score in our main experiment setting, showing that our framework can be generalized to other representation extraction methods.

where $v(\cdot)$ is the operation to take Delta Activations. We finetune models on pairs of domains and comparing their Delta Activations to those from individually trained models. Results are presented in Table 4, which shows that the cosine similarity between the mixed model and the sum of Delta Activations from the two individual models is consistently higher than the similarity with either individual model. This suggests that Delta Activations exhibit the additive property, which is especially important since models are often trained on mixed datasets. We report results across all ten domain pairs in Appendix 13, where the additive effect consistently holds.

| Math | Commonsense | Code | Mixed vs. D1 | Mixed vs. D2 | Mixed vs. Sum(D1, D2) |
|--------------|--------------|--------------|--------------|--------------|-----------------------|
| √ | ✓ | | .58 | .48 | .65 |
| \checkmark | | \checkmark | .70 | .27 | .73 |
| | \checkmark | \checkmark | .63 | .28 | .65 |

Table 4: **Additive property.** The sum of Delta Activations on models finetuned separately on two datasets aligns well with Delta Activations on the model finetuned on two datasets mixed together.

Probe dataset. We study the effects of number, length, and content of prompts in Table 5a, we see that using multiple prompts instead of one helps stabilize the embedding, while increasing the number from 5 to 20 offers no additional benefits. In Table 5b, we see that using GPT-4o to generate shorter versions of the Alpaca template, namely one-word and one-sentence versions, performs worse than a reasonably long instruction template. All prompts used in this part are included in Appendix D. Finally, Table 5c shows the importance of a generic instruction template by comparing the instruction template with domain-specific prompts or Wikitext. Domain-specific prompts perform the worst since they suppress the model's specialization and thereby cause model embeddings to become less distinguishable. A random generic text sampled from Wikitext performs slightly better while the instruction template achieves the best separation.

These findings across prompt number, length, and content bolster our design choices for the probe dataset. In addition, even with variations in these prompt settings, the effectiveness of Delta Activations is preserved, as evident from the fact that the silhouette scores stay well above zero.

| # of prompts | silh. score | length | silh. score | Content | silh. score |
|---------------|-------------|---------------------|-------------|-----------------|-------------|
| 1 | .57 | one-word | .45 | Wikitext | .44 |
| 5 | .61 | one-sentence | .59 | domain-specific | .42 |
| 20 | .61 | Alpaca (3-sentence) | .61 | instruction | .61 |
| (a) number of | of prompts | (b) length of pro | ompts | (c) content of | prompts |

Table 5: **Effects of number, length, and content of probe prompts.** Using multiple reasonably-long generic instruction templates makes the best probe dataset.

Where to extract activations. The embedding of the last token at the last layer is generally understood to encode the entire context. Decoder-only LLMs project this embedding to the logit space for next-token prediction. Consequently, Delta Activations are also derived using this embedding. Here we study whether Delta Activations could instead be sourced from other tokens or different layers. In Table 6a, we examine the effectiveness of calculating Delta Activations using the first, middle, and last tokens, as well as from the weighted average of all token embeddings following Muennighoff (2022). Overall, the results show that final tokens are effective targets for calculating Delta Activations.

We also investigate whether the final layer is the best position from which to extract activations. As shown in Table 6b, shallower layers perform worse than deeper ones; interestingly, the final layer is not optimal, as a layer at 2/3 of the total depth performs best. This phenomenon, where intermediate representation are found to be more effective for downstream tasks, is also observed in vision encoders (Chen et al., 2020; Bolya et al., 2025). Although a layer at 2/3 depth and weighted tokens exhibit slightly superior results, the final token at the last layer performs similarly. For simplicity, we use the last-layer final token embedding as the default setting for Delta Activations.

| Token Position | silh. score | Layer Position | silh. score |
|-----------------------------|-------------|---------------------|-------------|
| first | .22 | shallow (1/3 depth) | .51 |
| mid | .39 | mid (1/2 depth) | .61 |
| last | .61 | deep (2/3 depth) | .64 |
| weighted avg. of all tokens | .64 | last | .61 |
| (a) token position | n | (b) layer posi | ition |

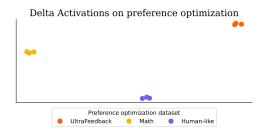
Table 6: **Effects of token and layer position to extract activations.** Later tokens and deeper layers produce better Delta Activations, with the 2/3 depth layer slightly surpassing the last layer.

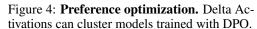
Beyond SFT: clustering by preference optimization. Our experiments thus far focused on the setting of Supervised Finetuning which use the token-level cross entropy loss. Preference optimization techniques (Rafailov et al., 2023; Meng et al., 2024) maximize the likelihood that preferred responses are ranked higher. The different supervision signal of preference optimization affects the activation differently. We explore whether Delta Activations still yield reliable clustering for models in this case. To construct the model pool, we perform preference optimization on LLAMA-3.1-8B-INSTRUCT using three disjoint 3000-example splits for each of three preference optimization datasets, namely UltraFeedback (Cui et al., 2024), HumanLLM (Çalık & Akkuş, 2025), and MetaMath-DPO (Pal et al., 2024; Yu et al., 2023). This experiment yields a silhouette score of 0.93, and the visualization presented in 4 shows well-separated cluster, which shows that Delta Activations can effectively capture similarity for preference optimization.

3.3 Extensions

In this section, we explore how Delta Activations as a model embedding can be extended to embedding a task, representing LLMs finetuned from different base LLMs, and guiding model merging.

Task embedding. We explore whether Delta Activations can embed tasks using only a few example. For each of the five domains in Section 3.1, we finetune the base LLM on 20 held-out examples that





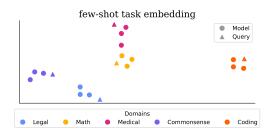


Figure 5: **Task embedding.** Few-shot task embedding is able to locate model clusters.

were not part of any previous training split. The detailed few-shot finetuning setting is in Appendix A.4. We then represent task with Delta Activations of the few-shot finetuned models.

We examine whether the task embedding can successfully locate the corresponding clusters on the three model pools constructed in Section 3.1. We define the *retrieval rate* metric as the fraction of few-shot task embeddings that correctly retrieve their corresponding full-model cluster via nearest-neighbour search using cosine similarity. Gemma achieves 100% retrieval rate while there is one failure case for each of LLaMA and Qwen. We present visualization of Gemma in Figure 5, which shows that few-shot queries reliably align with their corresponding full model clusters (circles). Despite being trained on only 20 examples, the resulting Delta Activations recover the domain cluster. This suggests that Delta Activations on few-shot trained model can be an effective task embedding.

Cross-base-model clustering. Delta Activations represent a finetuned model as the difference between its hidden states and those of its base model on the same inputs; therefore, they can only be directly applied to models derived from the *same* base. Interestingly, we find that this delta signal can transfer across bases. To test this, we evaluate two settings: cross-checkpoint and cross-architecture. In a cross-checkpoint setting (LLAMA-3-8B vs. LLAMA-3.1-8B; 10 models over 5 domains), Delta Activations achieved a silhouette score of **0.39**, cleanly recovering the five domain-specialization clusters (Figure 6a). In a cross-architecture setting (LLAMA-3.1-8B vs. LLAMA-3.2-1B; 10 models), Delta Activations are no longer feasible because the embedding dimension differs across architectures, projecting models into vectors of incompatible sizes. Instead, we adopt *Delta Meaning* (full implementation details in Appendix A.5), which is architecture-agnostic, and it successfully forms four out of five domain clusters with a silhouette score of **0.32** (Figure 6b).

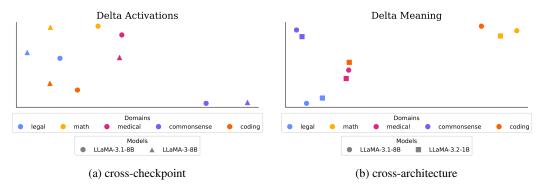


Figure 6: **Cross-base-model clustering.** (a) Delta Activations form domain clusters across models finetuned from LLAMA-3.1-8B and LLAMA-3-8B. (b) Delta Meaning form domain clusters across finetuned models of different sizes, LLAMA-3.1-8B and LLAMA-3.2-1B.

Model selection and similarity measurement. LoraHub (Huang et al., 2024a) hosts ~ 200 finetuned models based on FLAN-T5. The method is evaluated on the Big-Bench Hard (BBH) (Suzgun et al., 2022). For each task in BBH, LoraHub randomly select 20 LoRAs and optimizes their merging coefficients over few-shot examples from the target task. Our approach can be easily applied to this scenario by embedding the corresponding task using the provided few-shot examples and leverage Delta Activations similarity to replace the random model selection. Specifically, we identify the

single most-related LoRA model as an anchor and samples the remaining 19 models randomly. This simple selection strategy enabled by Delta Activations yields an average performance improvement of 2.0% by boosting the average accuracy from 34.3% to 36.3% on the 26 tasks of the BBH benchmark.

To probe whether Delta Activations can reveal model interference (Ortiz-Jimenez et al., 2023), we evaluate model merging by selecting the 20 most similar models in the embedding space. By construction, these models show extremely high average pairwise similarity (0.908 vs. –0.015 for random selection). If Delta Activations indeed capture entanglement, we expect such merging to perform poorly. Consistent with this hypothesis, the resulting average performance is only 30.3%, substantially below random selection. This confirms the presence of strong interference effects (Ortiz-Jimenez et al., 2023) and also that Delta Activations are informative about underlying relational structure between models.

4 Related Work

Embedding Models and Tasks. Ilharco et al. (2023) discovers that task vectors from finetuning can be merged to achieve multi-task learning. There are several attempts to represent LLMs using adapter weights (Ostapenko et al., 2024; Tan et al., 2025), evaluation profiles (Zhuang et al., 2025), or training-data-dependent characteristics (Zhao et al., 2024). Our work differs by requiring no external information or evaluation, instead deriving embeddings from internal behaviors.

Activations. Recent work uncovers structure in hidden activations of LLMs, understanding how massive activation act as biases to steer LLM output (Sun et al., 2024; Dettmers et al., 2022). Activations are also central to post-training compression: they are used to compute weight saliency for pruning (Frantar & Alistarh, 2023; Sun et al., 2023; Yin et al., 2023) and to reduce quantization error using calibration sets (Frantar et al., 2022; Lin et al., 2023; Zhang et al., 2023), which inspire the use of the probe dataset of Delta Activations. Additionally, learned activation shifts have been used to edit or transfer behavior across models (Liu et al., 2023; Li et al., 2023; Luo et al., 2024a).

Utilizing trained models. The landscape evolves from reusing a single model to multiple trained models. Finetuning (Zhuang et al., 2020) is a common framework to build on top of a single pretrained model, which becomes the common practice for LLMs (Chung et al., 2024; Rafailov et al., 2023). Other works use outputs (Hinton et al., 2015; Tian et al., 2019) or weights (LeCun et al., 1989; Han et al., 2015) of a pretrained model for the creation or initialization (Xu et al., 2024; Xia et al., 2024) of more efficient models. Efforts are made to effectively leverage multiple trained models through retrieval (Zhao et al., 2024; Jin et al., 2024; Luo et al., 2024b; Kahana et al., 2025), composition (Yang et al., 2024b; Chronopoulou et al., 2023; Huang et al., 2024a; Feng et al., 2024), or routing (Lu et al., 2023; Ong et al., 2024; Shnitzer et al., 2023).

Building model hubs. Recent works (Horwitz et al., 2025; Yu & Wang, 2025) study the public model pool and can systematically uncover finetuning relationship among trained models. Lorahub (Huang et al., 2024a), LoraRetriever (Zhao et al., 2024), and Learnware (Tan et al., 2025) create model hubs involving from ~ 40 to ~ 200 models. It is also possible to create larger model hubs via neural network parameter diffusion (Wang et al., 2024; Liang et al., 2025), though achieving true diversity (Zeng et al., 2025) may require further development.

5 Discussion

Delta Activations provide a lightweight and general-purpose method to represent finetuned LLMs by measuring shifts in their internal activations relative to a base LLM. This embedding space reliably organizes models by domain, generalizes across finetuning regimes, and exhibits additive structure that aligns with mixed-domain training. Unlike prior approaches, Delta Activations require neither metadata nor evaluation datasets and can be computed with a single forward pass, making them practical for large and evolving model repositories. Delta Activations extend beyond model embeddings, supporting applications like representing tasks and guiding model selection. More broadly, it fits within a Delta-X framework, where the same principle can be applied to other representations to suit different needs. As the community continues to generate vast numbers of post-trained LLMs, effective organization and reuse become essential. We hope Delta Activations provide a foundation for building structured, navigable model hubs.

Limitations and future work. Delta Activations introduces a novel way to represent finetuned models but also poses practical considerations. Our experiments focused on three prominent open-source backbones but further evaluation on other architectures would be valuable to understand its broader applicability. In addition, Delta Activations require access to internal hidden states, which is not feasible to be evaluated on proprietary models.

It is natural to ask how our method might perform on model pools substantially larger than those considered in our evaluation. Such a practical exercise would be most meaningful if the pool consists of models with diverse and uniquely valuable capabilities. While this is an interesting direction we intend to explore, we conjecture that today such pools exist primarily in proprietary settings (e.g. finetuned GPT models), and we hope our approach could facilitate sharing such models in future.

References

- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- Daniel Bolya, Po-Yao Huang, Peize Sun, Jang Hyun Cho, Andrea Madotto, Chen Wei, Tengyu Ma, Jiale Zhi, Jathushan Rajasegaran, Hanoona Rasheed, Junke Wang, Marco Monteiro, Hu Xu, Shiyu Dong, Nikhila Ravi, Daniel Li, Piotr Dollár, and Christoph Feichtenhofer. Perception encoder: The best visual embeddings are not at the output of the network. *arXiv*:2504.13181, 2025.
- Ethem Yağız Çalık and Talha Rüzgar Akkuş. Enhancing human-like responses in large language models. *arXiv preprint arXiv:2501.05032*, 2025.
- Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. Generative pretraining from pixels. In *ICML*, 2020.
- Alexandra Chronopoulou, Matthew E Peters, Alexander Fraser, and Jesse Dodge. Adaptersoup: Weight averaging to improve generalization of pretrained language models. In *ACL*, 2023.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 2024.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. arXiv preprint arXiv:2110.14168, 2021.
- CodeChef. Codechef. https://www.codechef.com, 2009. Competitive programming platform.
- Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Bingxiang He, Wei Zhu, Yuan Ni, Guotong Xie, Ruobing Xie, Yankai Lin, et al. Ultrafeedback: Boosting language models with scaled ai feedback. In *ICML*, 2024.
- Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Gpt3. int8 (): 8-bit matrix multiplication for transformers at scale. *Advances in neural information processing systems*, pp. 30318–30332, 2022.
- Shangbin Feng, Zifeng Wang, Yike Wang, Sayna Ebrahimi, Hamid Palangi, Lesly Miculicich, Achin Kulshrestha, Nathalie Rauschmayr, Yejin Choi, Yulia Tsvetkov, et al. Model swarms: Collaborative search to adapt llm experts via swarm intelligence. *arXiv preprint arXiv:2410.11163*, 2024.
- Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in one-shot. In *International Conference on Machine Learning*, pp. 10323–10337, 2023.
- Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training quantization for generative pre-trained transformers. *arXiv* preprint arXiv:2210.17323, 2022.
- FreedomIntelligence. Disease database. https://huggingface.co/datasets/FreedomIntelligence/Disease_Database, 2024. Hugging Face dataset.
- Neel Guha, Julian Nyarko, Daniel E. Ho, Christopher Ré, Adam Chilton, Aditya Narayana, Alex Chohlas-Wood, Austin Peters, Brandon Waldon, Daniel N. Rockmore, Diego Zambrano, Dmitry Talisman, Enam Hoque, Faiz Surani, Frank Fagan, Galit Sarfaty, Gregory M. Dickinson, Haggai Porat, Jason Hegland, Jessica Wu, Joe Nudell, Joel Niklaus, John Nay, Jonathan H. Choi, Kevin Tobia, Margaret Hagan, Megan Ma, Michael Livermore, Nikon Rasumov-Rahe, Nils Holzenberger, Noam Kolt, Peter Henderson, Sean Rehaag, Sharad Goel, Shang Gao, Spencer Williams, Sunny Gandhi, Tom Zur, Varun Iyer, and Zehua Li. Legalbench: A collaboratively built benchmark for measuring legal reasoning in large language models, 2023.
- Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for efficient neural network. In *NeurIPS*, 2015.

- Yifei He, Yuzheng Hu, Yong Lin, Tong Zhang, and Han Zhao. Localize-and-stitch: Efficient model merging via sparse task arithmetic. *Transactions on Machine Learning Research*, 2024, 2024. Preprint available at arXiv:2408.13656.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv* preprint arXiv:1503.02531, 2015.
- Eliahu Horwitz, Asaf Shul, and Yedid Hoshen. Unsupervised model tree heritage recovery. In *ICLR*, 2025.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *ICLR*, 2022.
- Chengsong Huang, Qian Liu, Bill Yuchen Lin, Tianyu Pang, Chao Du, and Min Lin. Lorahub: Efficient cross-task generalization via dynamic lora composition. In *COLM*, 2024a.
- Siming Huang, Tianhao Cheng, Jason Klein Liu, Jiaran Hao, Liuyihan Song, Yang Xu, J Yang, JH Liu, Chenchen Zhang, Linzheng Chai, et al. Opencoder: The open cookbook for top-tier code large language models. *arXiv preprint arXiv:2411.04905*, 2024b.
- Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Suchin Gururangan, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. Editing models with task arithmetic. In *ICLR*, 2023.
- Hamish Ivison, Yizhong Wang, Valentina Pyatkin, Nathan Lambert, Matthew Peters, Pradeep Dasigi, Joel Jang, David Wadden, Noah A. Smith, Iz Beltagy, and Hannaneh Hajishirzi. Camels in a changing climate: Enhancing Im adaptation with tulu 2, 2023.
- Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *arXiv* preprint arXiv:2009.13081, 2020.
- Pengfei Jin, Peng Shu, Sekeun Kim, Qing Xiao, Sifan Song, Cheng Chen, Tianming Liu, Xiang Li, and Quanzheng Li. Retrieval instead of fine-tuning: A retrieval-based parameter ensemble for zero-shot learning. *arXiv preprint arXiv:2410.09908*, 2024.
- Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. Pubmedqa: A dataset for biomedical research question answering. In *EMNLP-IJCNLP*, 2019.
- Jonathan Kahana, Or Nathan, Eliahu Horwitz, and Yedid Hoshen. Can this model also recognize dogs? zero-shot model search from weights. *arXiv* preprint arXiv:2502.09619, 2025.
- Yehuda Koren, Robert Bell, and Chris Volinsky. Matrix factorization techniques for recommender systems. *Computer*, 2009.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.
- Yann LeCun, John Denker, and Sara Solla. Optimal brain damage. In NeurIPS, 1989.
- Kenneth Li, Aspen K Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Emergent world representations: Exploring a sequence model trained on a synthetic task. *ICLR*, 2023.
- Zhiyuan Liang, Dongwen Tang, Yuhao Zhou, Xuanlei Zhao, Mingjia Shi, Wangbo Zhao, Zekai Li, Peihao Wang, Konstantin Schürholt, Damian Borth, et al. Drag-and-drop llms: Zero-shot prompt-to-weights. *arXiv preprint arXiv:2506.16406*, 2025.
- Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Xingyu Dang, and Song Han. Awq: Activation-aware weight quantization for llm compression and acceleration. *MlSys*, 2023.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024a.

- Sheng Liu, Haotian Ye, Lei Xing, and James Zou. In-context vectors: Making in context learning more effective and controllable through latent space steering. *arXiv preprint arXiv:2311.06668*, 2023.
- Tian Yu Liu, Matthew Trager, Alessandro Achille, Pramuditha Perera, Luca Zancato, and Stefano Soatto. Meaning representations from trajectories in autoregressive models. In *ICLR*, 2024b.
- Keming Lu, Hongyi Yuan, Runji Lin, Junyang Lin, Zheng Yuan, Chang Zhou, and Jingren Zhou. Routing to the expert: Efficient reward-guided ensemble of large language models. *arXiv* preprint arXiv:2311.08692, 2023.
- Jinqi Luo, Tianjiao Ding, Kwan Ho Ryan Chan, Darshan Thaker, Aditya Chattopadhyay, Chris Callison-Burch, and René Vidal. Pace: Parsimonious concept engineering for large language models. In *NeurIPS*, 2024a.
- Michael Luo, Justin Wong, Brandon Trabucco, Yanping Huang, Joseph E Gonzalez, Ruslan Salakhutdinov, Ion Stoica, et al. Stylus: Automatic adapter selection for diffusion models. In *NeurIPS*, 2024b.
- Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. Wizardcoder: Empowering code large language models with evol-instruct. *ICLR*, 2024c.
- Yu Meng, Mengzhou Xia, and Danqi Chen. Simpo: Simple preference optimization with a reference-free reward. In *NeurIPS*, 2024.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. In *ICLR Workshop*, 2013.
- Mohamed-Ahmed161. Disease-symptoms dataset. https://huggingface.co/datasets/Mohamed-Ahmed161/Disease-Symptoms, 2024. Hugging Face dataset.
- Niklas Muennighoff. Sgpt: Gpt sentence embeddings for semantic search. arXiv preprint arXiv:2202.08904, 2022.
- Isaac Ong, Amjad Almahairi, Vincent Wu, Wei-Lin Chiang, Tianhao Wu, Joseph E Gonzalez, M Waleed Kadous, and Ion Stoica. Routellm: Learning to route llms from preference data. In *ICLR*, 2024.
- Guillermo Ortiz-Jimenez, Alessandro Favero, and Pascal Frossard. Task arithmetic in the tangent space: Improved editing of pre-trained models. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *NeurIPS*, 2023.
- Oleksiy Ostapenko, Zhan Su, Edoardo Maria Ponti, Laurent Charlin, Nicolas Le Roux, Matheus Pereira, Lucas Caccia, and Alessandro Sordoni. Towards modular llms by building and reusing a library of loras. *arXiv preprint arXiv:2405.11157*, 2024.
- Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In *Proceedings of the Conference on Health, Inference, and Learning*, 2022.
- Arka Pal, Deep Karkhanis, Samuel Dooley, Manley Roberts, Siddartha Naidu, and Colin White. Smaug: Fixing failure modes of preference optimisation with dpo-positive. *arXiv preprint arXiv:2402.13228*, 2024.
- Shanghaoran Quan, Jiaxi Yang, Bowen Yu, Bo Zheng, Dayiheng Liu, An Yang, Xuancheng Ren, Bofei Gao, Yibo Miao, Yunlong Feng, Zekun Wang, Jian Yang, Zeyu Cui, Yang Fan, Yichang Zhang, Binyuan Hui, and Junyang Lin. Codeelo: Benchmarking competition-level code generation of llms with human-comparable elo ratings, 2025.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. In *NeurIPS*, 2023.

- Yi Ren and Danica J Sutherland. Learning dynamics of Ilm finetuning. In ICLR, 2025.
- Peter J. Rousseeuw. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics*, 1987.
- Tal Shnitzer, Anthony Ou, Mírian Silva, Kate Soule, Yuekai Sun, Justin Solomon, Neil Thompson, and Mikhail Yurochkin. Large language model routing with benchmark datasets. *arXiv preprint arXiv:2309.15789*, 2023.
- Mingjie Sun, Zhuang Liu, Anna Bair, and J Zico Kolter. A simple and effective pruning approach for large language models. *arXiv preprint arXiv:2306.11695*, 2023.
- Mingjie Sun, Xinlei Chen, J. Zico Kolter, and Zhuang Liu. Massive activations in large language models. *COLM*, 2024.
- Mingjie Sun, Yida Yin, Zhiqiu Xu, J. Zico Kolter, and Zhuang Liu. Idiosyncrasies in large language models. In ICML, 2025.
- Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, et al. Challenging big-bench tasks and whether chain-of-thought can solve them. *arXiv preprint arXiv:2210.09261*, 2022.
- Zhi-Hao Tan, Zi-Chen Zhao, Hao-Yu Shi, Xin-Yu Zhang, Peng Tan, Yang Yu, and Zhi-Hua Zhou. Learnware of language models: Specialized small language models can do big. *arXiv preprint arXiv:2505.13425*, 2025.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. Alpaca: A strong, replicable instruction-following model. *Stanford Center for Research on Foundation Models*, 2023.
- Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024.
- Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive representation distillation. *arXiv* preprint arXiv:1910.10699, 2019.
- Hugo Touvron et al. Llama: Open and efficient foundation language models. *arXiv preprint* arXiv:2302.13971, 2023.
- Kai Wang, Dongwen Tang, Boya Zeng, Yida Yin, Zhaopan Xu, Yukun Zhou, Zelin Zang, Trevor Darrell, Zhuang Liu, and Yang You. Neural network diffusion. *arXiv preprint arXiv:2402.13144*, 2024.
- Wenhui Wang, Hangbo Bao, Shaohan Huang, Li Dong, and Furu Wei. MiniLMv2: Multi-head self-attention relation distillation for compressing pretrained transformers. In *ACL*, 2021.
- Mengzhou Xia, Tianyu Gao, Zhiyuan Zeng, and Danqi Chen. Sheared LLaMA: Accelerating language model pre-training via structured pruning. In *ICLR*, 2024.
- Zhiqiu Xu, Yanjie Chen, Kirill Vishniakov, Yida Yin, Zhiqiang Shen, Trevor Darrell, Lingjie Liu, and Zhuang Liu. Initializing models with larger ones. In *ICLR*, 2024.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024a.
- Enneng Yang, Li Shen, Guibing Guo, Xingwei Wang, Xiaochun Cao, Jie Zhang, and Dacheng Tao. Model merging in llms, mllms, and beyond: Methods, theories, applications and opportunities. *arXiv preprint arXiv:2408.07666*, 2024b.

- Lu Yin, You Wu, Zhenyu Zhang, Cheng-Yu Hsieh, Yaqing Wang, Yiling Jia, Mykola Pechenizkiy, Yi Liang, Zhangyang Wang, and Shiwei Liu. Outlier weighed layerwise sparsity (owl): A missing secret sauce for pruning llms to high sparsity. *arXiv* preprint arXiv:2310.05175, 2023.
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for large language models. *arXiv preprint arXiv:2309.12284*, 2023.
- Runpeng Yu and Xinchao Wang. Neural phylogeny: Fine-tuning relationship detection among neural networks. In ICLR, 2025.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? In *ACL*, 2019.
- Boya Zeng, Yida Yin, Zhiqiu Xu, and Zhuang Liu. Generative modeling of weights: Generalization or memorization? *arXiv preprint arXiv:2506.07998*, 2025.
- Yuxin Zhang, Lirui Zhao, Mingbao Lin, Yunyun Sun, Yiwu Yao, Xingjia Han, Jared Tanner, Shiwei Liu, and Rongrong Ji. Dynamic sparse no training: Training-free fine-tuning for sparse llms. *arXiv* preprint arXiv:2310.08915, 2023.
- Ziyu Zhao, Leilei Gan, Guoyin Wang, Wangchunshu Zhou, Hongxia Yang, Kun Kuang, and Fei Wu. LoraRetriever: Input-aware LoRA retrieval and composition for mixed tasks in the wild. In *ACL Findings*, 2024.
- Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 2020.
- Richard Zhuang, Tianhao Wu, Zhaojin Wen, Andrew Li, Jiantao Jiao, and Kannan Ramchandran. EmbedLLM: Learning compact representations of large language models. In *ICLR*, 2025.

A Training settings

For LoRA, we set rank r = 8, $\alpha = 16$, targeting query, key, value, and MLP projections (q_proj, k_proj, v_proj, up_proj, down_proj, gate_proj), with no dropout and no bias parameters.

A.1 Datasets

| Dataset | Description |
|---------------|--|
| OpenCoder-LLM | Educational programming instructions (opc-sft-stage2) |
| GSM8K | Grade-school math problems with chain-of-thought reasoning |
| HellaSwag | Commonsense natural language inference |
| LegalBench | Privacy policy question answering |
| PubMedQA | Medical Q&A from PubMed abstracts |

Table 7: **Summary of datasets used in the experiments.** Each dataset was split into three subsets of 3,000 examples.

A.2 Prompt Templates

| >p0.32 X |
|---|
| Prompt Type Template (truncated) |
| Task-specific (Programming) Below is an instruction that describes a programming task |
| Task-specific (Math) Below is a grade-school math problem. Please work through the |
| reasoning step-by-step |
| Task-specific (Legal) Below is a legal-reasoning task from the LegalBench benchmark |
| Task-specific (Medical) Below is a medical question based on a PubMed article |
| Task-specific (Commonsense) Below is a scenario. What happens next in this paragraph |
| Universal Below is an instruction that describes a task. Write a response that |
| appropriately completes the request. |

Table 8: **Prompt templates used during training.** Task-specific templates are customized to each domain; the universal prompt is applied uniformly across all tasks during finetuning.

A.3 Training Hyperparameters

| Training Setting | Configuration |
|-----------------------------|---|
| Epochs | 3 |
| Batch size | 4 |
| Learning rate | 1e-4 |
| Optimizer | paged_adamw_32bit |
| Gradient accumulation steps | 2 |
| Warmup steps | 10 |
| Max sequence length | 512 tokens (8,192 tokens for OpenCoder-LLM) |
| Hardware | NVIDIA H100 (80GB) |

Table 9: Training hyperparameters used across all experiments.

Data was collated using the DataCollatorForCompletionOnlyLM from the TRL library, computing loss only on the response portion. We deploy all models to the Hugging Face Hub with standardized nomenclature indicating the base model, dataset, training approach, and key hyperparameters.

A.4 Few-shot Task Embedding Configuration

For the few-shot task embedding experiments described in Section 3.3, we used a configuration tailored for limited data scenarios. We randomly sampled 20 examples from each domain's held-out data that was not used in any of our main experiment splits.

| Parameter | Value |
|---------------------|-----------------------|
| Examples per domain | 20 (randomly sampled) |
| Learning rate | 3.3e-3 |
| Batch size | 1 |
| Epochs | 5 |

Table 10: Hyperparameters for few-shot task embedding experiments.

We maintained the same LoRA configuration as in our main experiments. The higher learning rate and increased number of epochs compensate for the limited training data while the smaller batch size allows for more frequent parameter updates. After fine-tuning, we computed Delta Activations using identical probe datasets as in our main experiments to ensure direct comparability between few-shot task embeddings and full model embeddings in the embedding space.

A.5 Delta Meaning implementation

Here we provide implementation details for *Delta Meaning*, our adaptation of Meaning Representations (Liu et al., 2024b) to the Delta framework. This extension enables model embeddings across heterogeneous backbones, where direct activation comparisons are infeasible.

Meaning representations. Given a probe prompt x, we first sample n continuations $\{s_1, \ldots, s_n\}$ from the base model (temperature = 1.0). For any finetuned model f, we then score each continuation s_i by computing its inverse perplexity under f:

$$m_f(x)_i = \exp\left(-\frac{1}{|s_i|} \sum_{t=1}^{|s_i|} \log p_f(s_{i,t} \mid s_{i,< t}, x)\right).$$

This produces an *n*-dimensional "meaning vector" $m_f(x) \in \mathbb{R}^n$ for each prompt x.

Delta aggregation. For a finetuned model f and its base f_{base} , we define the *Delta Meaning* on prompt x as the difference between their meaning vectors:

$$\Delta_f(x) = m_f(x) - m_{f_{\text{base}}}(x).$$

Aggregating across all prompts $x \in D_{\text{probe}}$ yields the final model embedding:

$$v_f = \frac{1}{|D_{\text{probe}}|} \sum_{x \in D_{\text{probe}}} \Delta_f(x).$$

Hyperparameters. In our experiments, we set n to be 20. Larger n provides more informative embeddings but requires proportionally more forward passes, as each continuation must be scored by both the base and finetuned models. Despite this, the dimensionality of Delta Meaning remains extremely compact compared to weight- or logit-based alternatives. Importantly, because any model can evaluate the probability of a given text sequence, Delta Meaning embeddings are naturally architecture-agnostic, allowing us to cluster finetuned models drawn from multiple backbones.

B Additional Analysis

B.1 Beyond domains: clustering by dataset properties.

In our setting in Section 3.1, each domain corresponds to a well-defined task with relatively uniform answer structure (e.g., multiple-choice). Here, we ablate that structure. We construct a new model pool by finetuning 3 models on each of five distinct subsets of Tulu v2 (Ivison et al., 2023): COT, GPT4-ALPACA, SHAREGPT, CODEALPACA, and SCIENCE. Unlike domain datasets, these are less semantically disjoint. Instead, they differ in instruction format, conversational structure, and expected output style—ranging from open-ended dialogue to multi-step chain-of-thought reasoning, code snippets, and longform factual answers. Crucially, this makes the output distribution far more diverse and less predictable. Results are presented in Table 11. With no consistent answer template, output sentence embeddings fail to reflect model specialization for LLaMA and Qwen, whereas Delta Activations continue to achieve decent clustering quality.

| Embedding Space | LLaMA | Gemma | Qwen | Avg. |
|----------------------------|-------|-------|------|------|
| output sentence embeddings | .06 | .23 | 03 | .08 |
| Delta Activations | .33 | .41 | .48 | .41 |

Table 11: Clustering quality (silhouette score) across Tulu v2 instruction splits. Delta Activations remain effective despite diverse output structures and blurred instruction boundaries.

B.2 Robustness to training settings.

Do differences in training settings have a greater impact on Delta Activations than the choice of finetuned domains? To evaluate this, we systematically perturbed the training process for models within our domain clusters. In our main setting, the model pool is organized into 5 distinct domain clusters, with 3 models in each cluster trained using identical settings. To test the impact of a specific training configuration—for instance, learning rate—the three models within each of the 5 domain clusters were trained using three different learning rates respectively (e.g., model 1 with $1e^{-4}$, model 2 with $4e^{-4}$, model 3 with $1e^{-5}$ within a single domain cluster). This process was independently repeated for variations in the number of training examples and the number of training epochs, where we vary number of training examples by 100, 1000, and 10000 and number of epochs by 1, 2, and 3.

Table 12 presents results which measure the clustering quality on domains when subjected to such training variations. indicates that varied training settings generally did not break domain-specific clustering. Changes to the amount of training data or the number of epochs had minimal effect on the quality of these clusters, which remained comparable to those formed under identical training settings. The different-learning-rate setting yields a lower silhouette score, as expected, since learning rate significantly impacts training dynamics and tends to increase within-cluster variation. These observations confirm that Delta Activations effectively identify finetuning domains despite common variations in training procedures. On the other hand, these results also show the strength of Delta Activations in identifying the nuanced differences within each cluster.

| Training Setting | LLaMA | Gemma | Qwen | Avg. |
|---------------------------------------|-------|-------|------|------|
| Different number of training examples | .66 | .51 | .68 | .62 |
| Different learning rates | .53 | .37 | .23 | .38 |
| Different training epochs | .62 | .59 | .51 | .57 |
| Identical training settings | .65 | .55 | .65 | .61 |

Table 12: **Delta Activations' embeddings are robust to training hyperparameters.** Models trained in varying settings still form tight domain-specific clusters, comparable to those trained identically.

B.3 t-SNE Visualization on more backbones

We show visualization of the experiments conducted in Section 3.1 on LLaMA and Qwen in Figure 7 and Figure 8 respectively. Over different backbones, the visualization consistently shows the superiority of Delta Activations in forming cleanly-separated clusters.

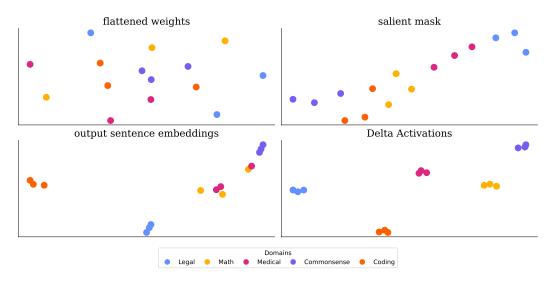


Figure 7: t-SNE visualization of different embedding spaces (LLaMA).

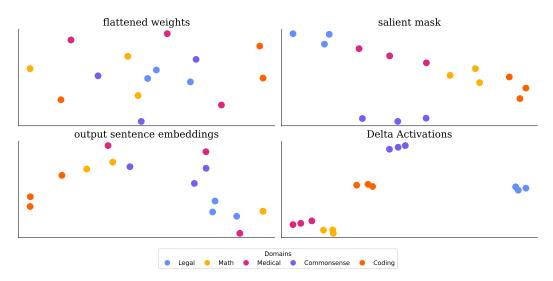


Figure 8: t-SNE visualization of different embedding spaces (Qwen).

B.4 Full Results for Additive Properties experiment

To better understand the additive nature of Delta Activations, we extend the experiment from Section 3.1 to cover all ten domain pairs. In each case, we compare the Delta Activation vector from a model trained on the mixed dataset against those from models trained individually on each domain, as well as their vector sum. This comprehensive table demonstrates that the summed Delta Activations consistently better approximate the mixed-model embedding, reinforcing the additive property of Delta Activations.

| Math | Commonsense | Code | Medical | Legal | Mixed v. D1 | Mixed v. D2 | Mixed v. Sum |
|--------------|--------------|--------------|--------------|--------------|-------------|-------------|--------------|
| √ | | ✓ | | | .703 | .270 | .726 |
| \checkmark | \checkmark | | | | .577 | .484 | .649 |
| | \checkmark | \checkmark | | | .631 | .283 | .653 |
| | | | \checkmark | \checkmark | .407 | .675 | .695 |
| | \checkmark | | \checkmark | | .359 | .662 | .677 |
| \checkmark | | | \checkmark | | .760 | .697 | .811 |
| | | \checkmark | \checkmark | | .462 | .581 | .693 |
| | \checkmark | | | \checkmark | .649 | .659 | .763 |
| \checkmark | | | | \checkmark | .522 | .610 | .669 |
| | | \checkmark | | \checkmark | .445 | .507 | .587 |

Table 13: **Additive property (full).** The sum of Delta Activations on models finetuned separately on two datasets aligns well with Delta Activations on the model finetuned on two datasets mixed together.

B.5 Full finetuning

We conduct experiments in Section 3.1 on LLAMA-3.1-8B with the model pool trained using full finetuning instead of LoRA. This experiment yields a silhouette score of **0.63**, which confirms that Delta Activations provide clear clustering regardless of finetuning methods. Visualization is shown in Figure 9.

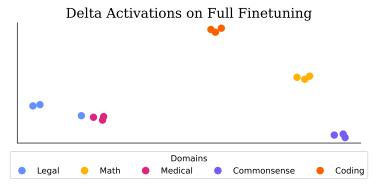


Figure 9: t-SNE visualization of full finetuning.

Intra-domain clustering. We conduct experiments to test whether Delta Activations provide sufficient signal to differentiate models trained on different sub-expertises within medical and coding domains.

We partition each domain into multiple sub-expertise splits. For medical, we create finetuned 8 models: disease databases (FreedomIntelligence, 2024) (2 splits), disease-symptoms (Mohamed-Ahmed161, 2024) (2 splits), MedQA-USMLE (Jin et al., 2020) (2 splits), and MedMCQA (Pal et al., 2022) (2 splits). For coding, we create 6 models covering C++ (2 splits), Python (2 splits), and Java (2 splits) from CodeChef (CodeChef, 2009), Codeforces (Quan et al., 2025), and Evol-Instruct (Luo et al., 2024c) datasets. We compute Delta Activations using both generic instruction prompts and domain-specialized prompts.

Table 14 shows that generic prompts achieve a silhouette score of 0.657 for medical sub-expertise, outperforming specialized medical prompts (0.533). For coding, both prompts perform comparably (0.668 vs 0.674). The generic probe maintains sufficient resolution to form correct clusters of sub-expertises within domains.

| Domain | Generic prompt | Specialized prompt |
|---------|----------------|--------------------|
| Medical | .657 | .533 |
| Coding | .668 | .674 |

Table 14: Intra-domain clustering. Silhouette scores for sub-expertise clustering within domains.

C Use of Large Language Models (LLMs)

We use LLMs to polish writing, and assist in writing code for our method and preparing demo.

D Prompt Templates

Table 15, Table 16, and Table 17 list all prompt templates used in probe dataset experiments in Section 3.2. For other experiments, we use the first five prompts in Table 15.

| ID | Prompt Template |
|----|--|
| 1 | Below is an instruction that describes a task. Write a response that appropriately completes the |
| | request. |
| 2 | The task described below requires a response that completes the request accurately. |
| 3 | Below is a description of a task. Provide a response that aligns with the requirements. |
| 4 | The following instruction outlines a task. Generate a response that meets the specified request. |
| 5 | You are given an instruction and input. Write a response that completes the task as requested. |
| 6 | You are provided with a task instruction and input. Write a response that fulfills the described |
| | requirements. |
| 7 | Here is an instruction and its associated input. Complete the task with an appropriate response. |
| 8 | Below is a task along with its context. Write a response that matches the requirements. |
| 9 | The following is a description of a task and its input. Generate a response that fulfills the request. |
| 10 | An outlined task is provided along with its input. Write a response that satisfies the given |
| | instruction. |
| 11 | Given the following instruction, generate a suitable response that fulfills the request. |
| 12 | The task described below requires a response that completes the request accurately. |
| 13 | Below is a description of a task. Provide a response that aligns with the requirements. |
| 14 | The following instruction outlines a task. Generate a response that meets the specified request. |
| 15 | You are given an instruction and input. Write a response that completes the task as requested. |
| 16 | Here is an instruction and its associated input. Create a response that properly addresses the |
| | request. |
| 17 | Below is a task description. Provide an appropriate response that matches the input. |
| 18 | An instruction and input are provided. Write a response that accurately completes the task. |
| 19 | The following is an instruction that describes a task. Write a response that correctly satisfies the |
| | request. |
| 20 | Below is an outlined task. Respond with a completion that fits the instruction and input given. |

Table 15: List of paraphrased prompt templates used in our experiments.

| ID | Prompt Template | |
|----|--|--|
| 1 | Instruction: Please provide a response. Input: | |
| | Input. | |
| 2 | Please perform the following task. | |
| 3 | Complete the instruction. | |
| 4 | Provide the appropriate response. | |
| 5 | Here is the text. Response: | |

Table 16: One-sentence paraphrased prompt templates.

| ID | Prompt Template |
|----|-----------------|
| 1 | Response: |
| 2 | Answer: |
| 3 | Explanation: |
| 4 | Solution: |
| 5 | Discussion: |

Table 17: One-word paraphrased prompt templates.

E Broader Impact

The proposed Delta Activations method facilitate efficient reuse of fine-tuned models by providing an embedding to encode the finetuned model's behaviors and capability. This reduces redundant training, cutting energy costs and promoting sustainable AI practices. Furthermore, it encourages broader public sharing of fine-tuned models by offering clear documentation of their capabilities, accelerating research and collaboration.

However, expanding public model hubs also introduces risks, as low-quality or adversarial models could contaminate the pool. This highlights the need for careful curation to maintain reliability and safety in open model ecosystems.