

Fresh in memory:
Training-order recency is linearly
encoded in LLM activations

Dmitrii Krasheninnikov

March 2026

FRESH IN MEMORY: TRAINING-ORDER RECENCY IS LINEARLY ENCODED IN LANGUAGE MODEL ACTIVATIONS

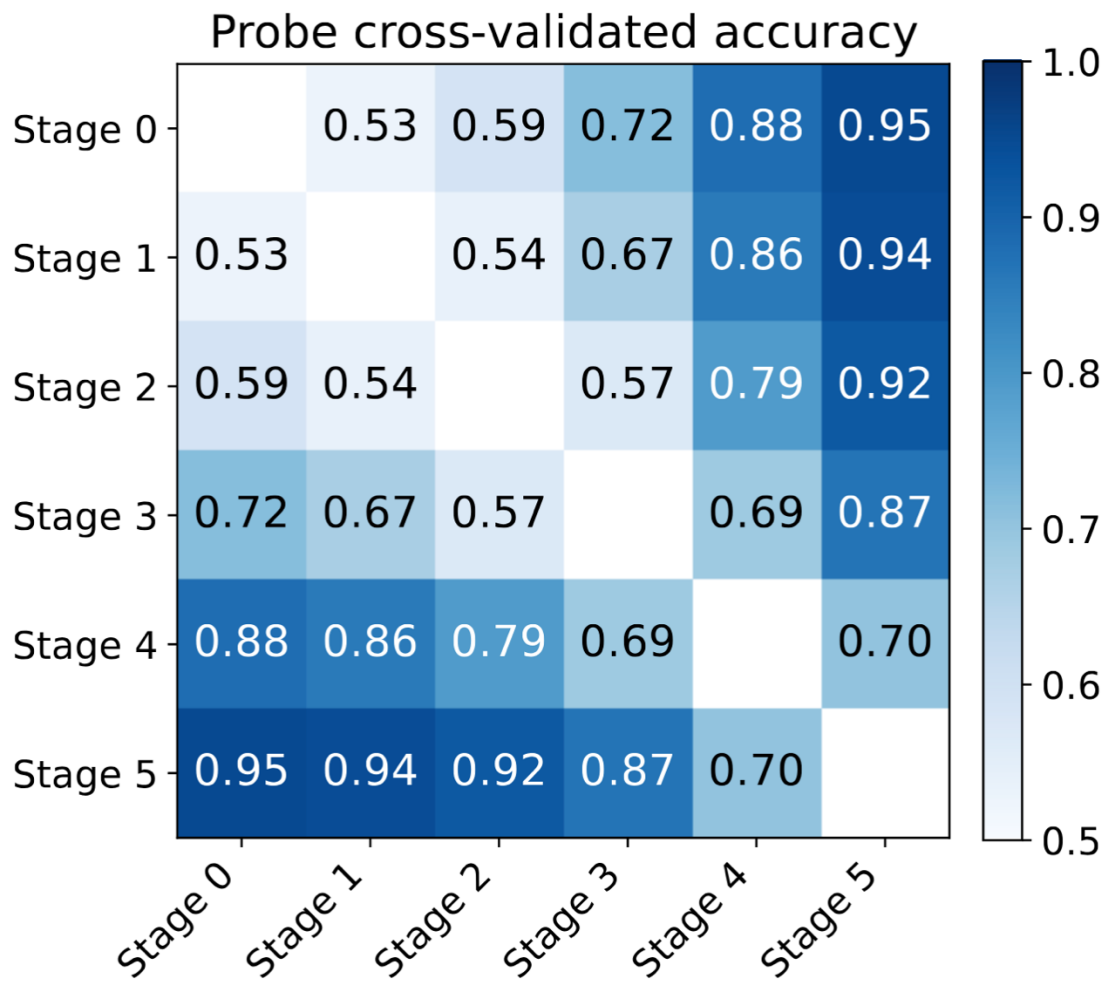
Dmitrii Krasheninnikov¹

Richard E. Turner¹

David Krueger²

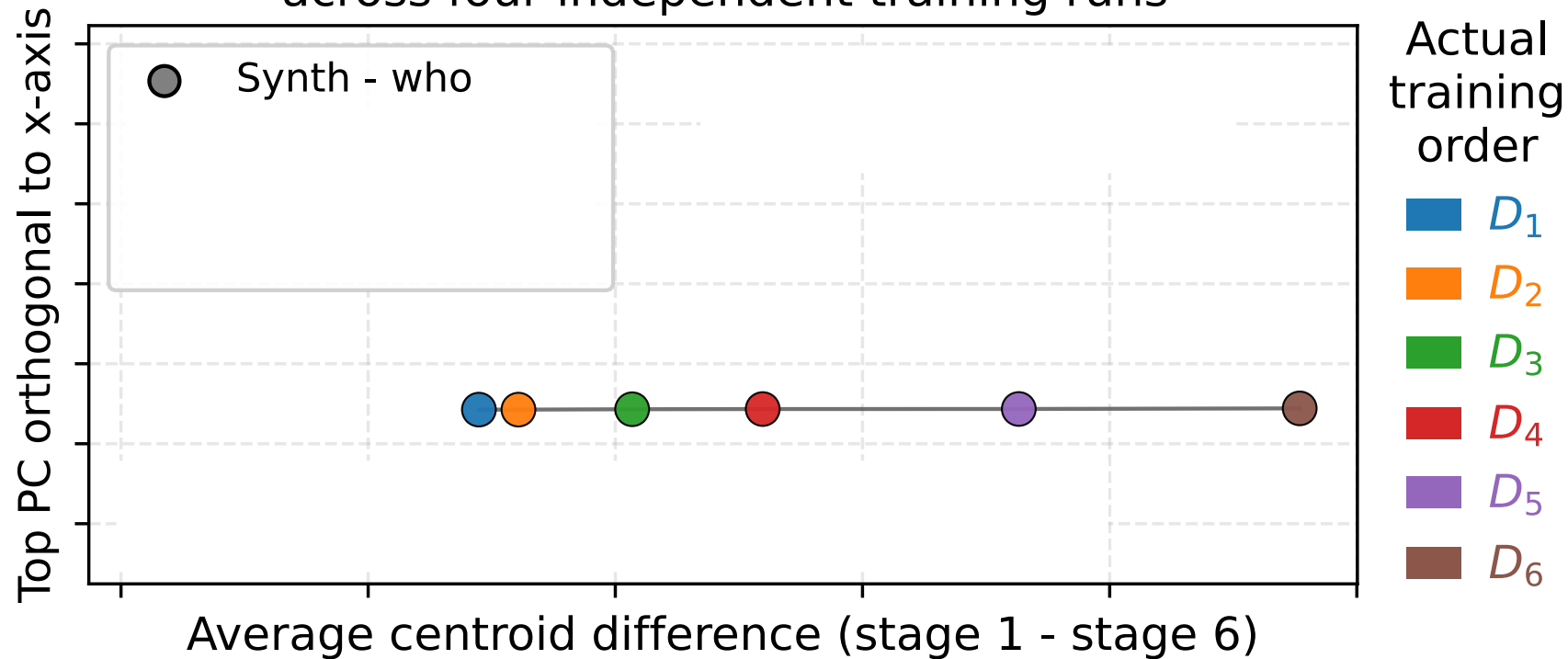
ABSTRACT

We show that language models’ activations linearly encode when information was learned during training. Our setup involves creating a model with a known training order by sequentially fine-tuning Llama-3.2-1B on six disjoint but otherwise similar datasets about named entities. We find that the average activations of test samples corresponding to the six training datasets encode the training order: when projected into a 2D subspace, these centroids are arranged exactly in the order of training and lie on a straight line. Further, we show that linear probes can accurately ($\sim 90\%$) distinguish “early” vs. “late” entities, generalizing to entities unseen during the probes’ own training. The model can also be fine-tuned to explicitly report an unseen entity’s training stage ($\sim 80\%$ accuracy). Notably, the training-order encoding does not seem attributable to simple differences in activation magnitudes, losses, or model confidence. Our paper shows that models can differentiate information by its acquisition time, and carries significant implications for how they might manage conflicting data and respond to knowledge modifications.



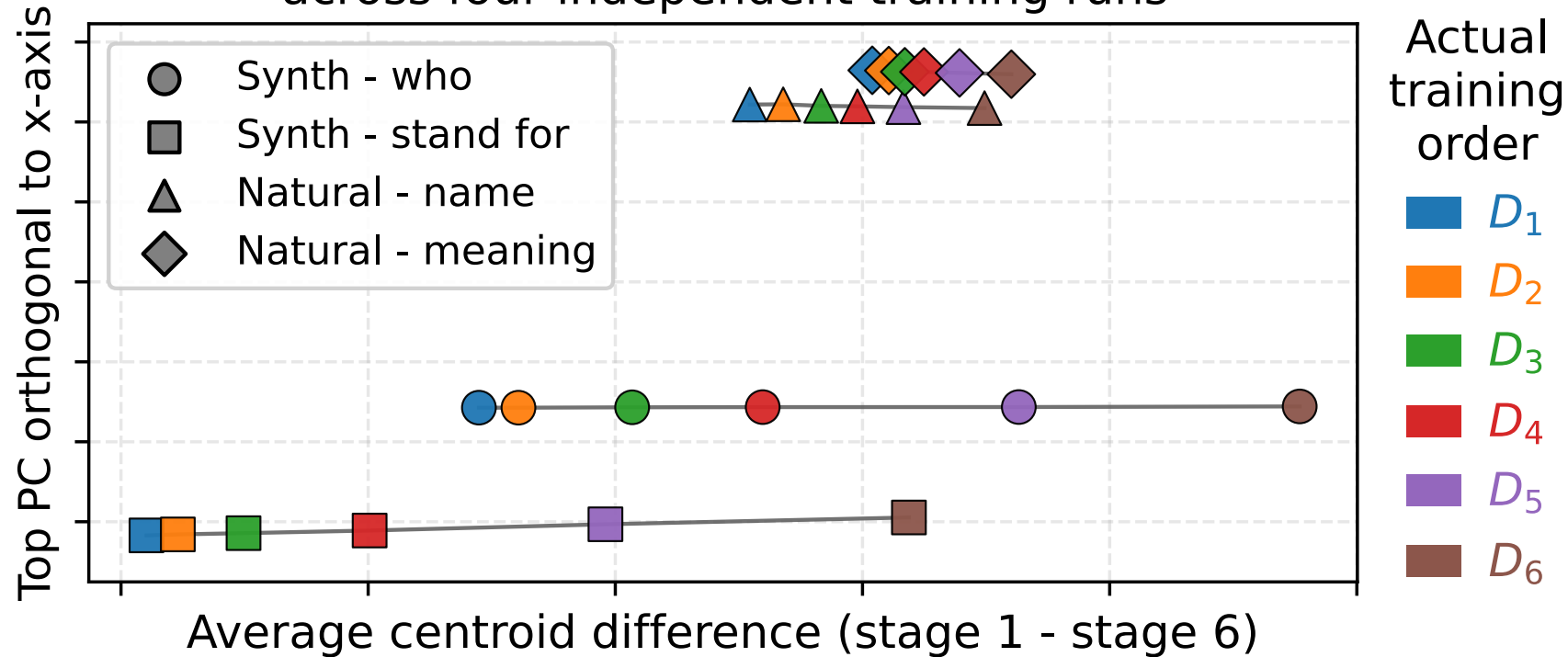
Training-order recency is encoded in activations

Activation centroids (averages) for the six *test* datasets, across four independent training runs



Training-order recency is encoded in activations

Activation centroids (averages) for the six *test* datasets, across four independent training runs



Check the paper for

- Results showing that “training-order recency” is a good interpretation
- Models can report the training stage in output tokens when finetuned to do so
- Mechanistic study of what the discovered axis encodes
- And more!