



## Do LLMs solve novel tasks?

**Out-of-distribution (OOD)** generalization measures the performance on novel tasks  $\mathcal{P}_{\text{train}} \neq \mathcal{P}_{\text{test}}$ . New challenges since advent of LLMs.

- Prompting, In-context learning.
- Compositional structure.
- Tasks that require “reasoning”.

**Goal:** In-depth empirical analysis to understand

- How composition is internally represented by LLMs;
- How critical geometric structure emerges from training;
- How they empower language & reasoning tasks across a wide range of models.

## A primer on Transformers: How concepts are represented

Transformers from circuits perspective. Let  $\mathbf{X} = [\mathbf{x}_1^{(\ell)}, \dots, \mathbf{x}_T^{(\ell)}]^\top \in \mathbb{R}^{T \times d}$  be the input vectors or the hidden states at a layer  $\ell$ .

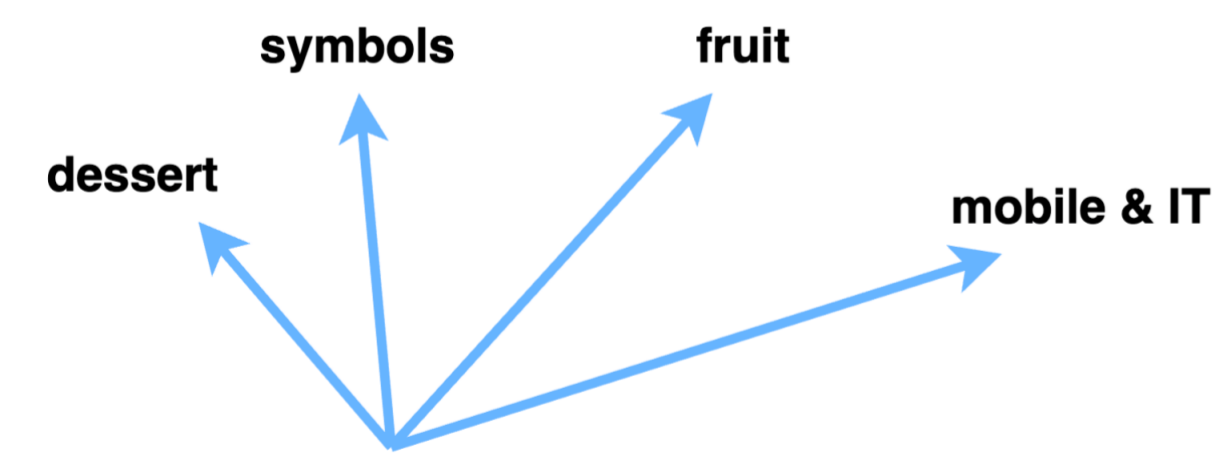
$$\mathbf{X} \leftarrow \mathbf{X} + \text{MSA}(\mathbf{X}; \mathbf{W}^{(\ell)}), \quad \mathbf{X} \leftarrow \mathbf{X} + \text{FFN}(\mathbf{X}; \widetilde{\mathbf{W}}^{(\ell)})$$

$$\text{MSA}(\mathbf{X}; \mathbf{W}) = \sum_{j=1}^H \text{Softmax} \left( \underbrace{\underbrace{(\mathbf{X}\mathbf{W}_{\text{QK},j}\mathbf{X}^\top)}_{\text{attention matrix}}}_{\text{QK circuit reads and matches info from stream}} \underbrace{\mathbf{X}\mathbf{W}_{\text{OV},j}^\top}_{\text{OV circuit writes and adds info to stream}} \right)$$

where  $\mathbf{W}_{\text{QK}}, \mathbf{W}_{\text{OV}} \in \mathbb{R}^{d \times d}$  are query-key, output-value matrices.

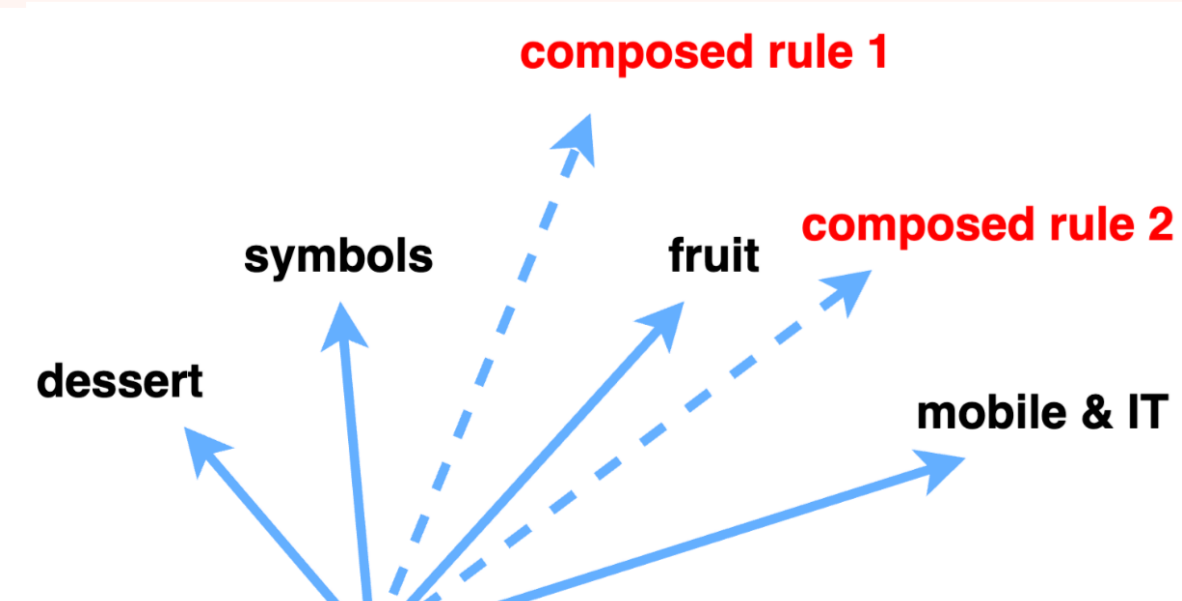
**Linear representation hypothesis:** concepts are encoded as linear subspaces within the embedding space.

**Feature superposition:** hidden states are sparse linear combinations of base concept vectors from a large dictionary.



$$\text{apple} = 0.09 \text{“dessert”} + 0.11 \text{“organism”} + 0.16 \text{“fruit”} + 0.22 \text{“mobile\&IT”} + 0.42 \text{“other”}$$

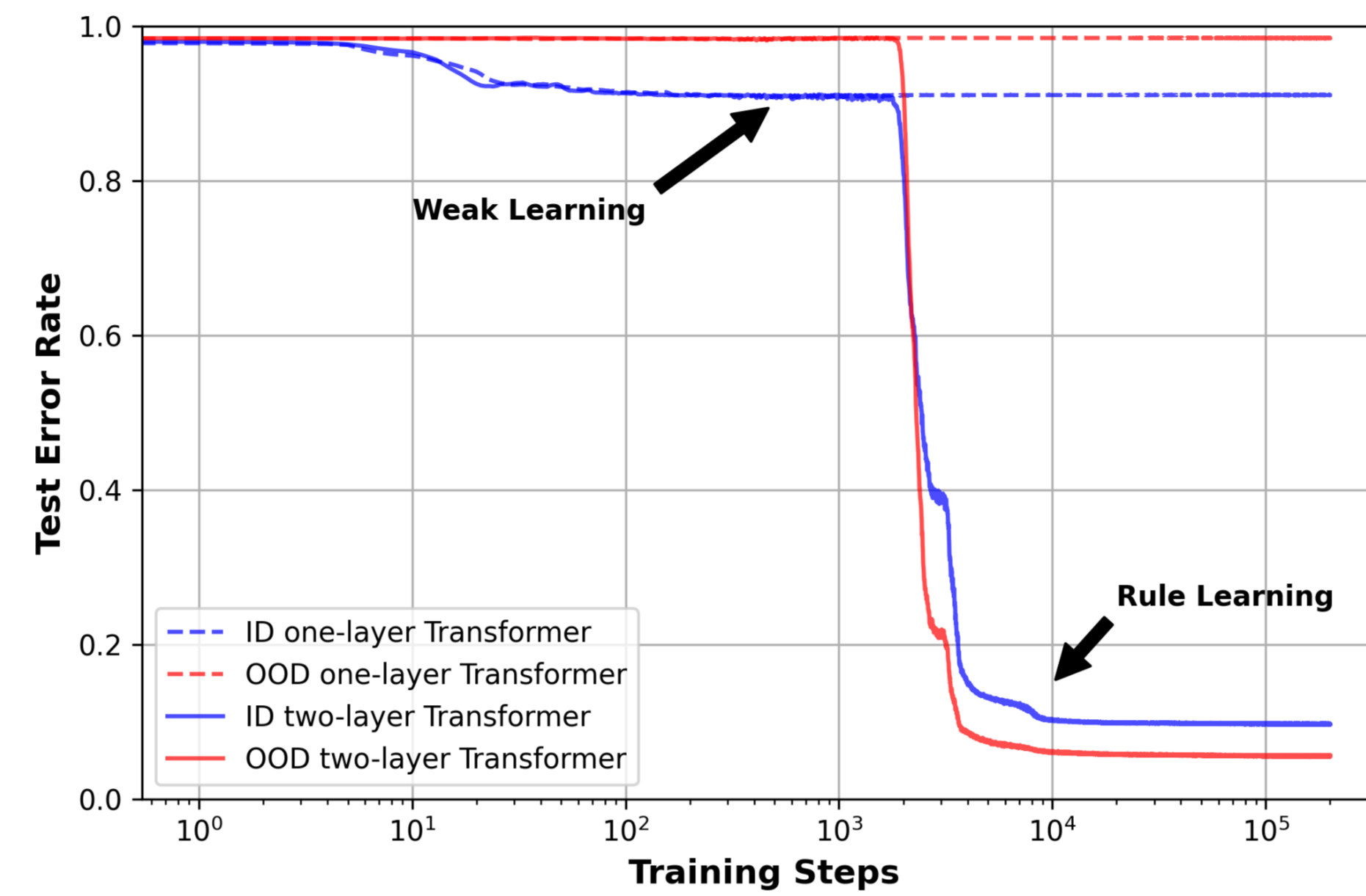
## Main message: composition through subspace matching empowers OOD generalization



- Concept subspaces + rule subspaces
- Composed rule 1 (e.g., copying), composed rule 2 ...
- Enables OOD generalization, esp. in novel context (ICL, CoT)

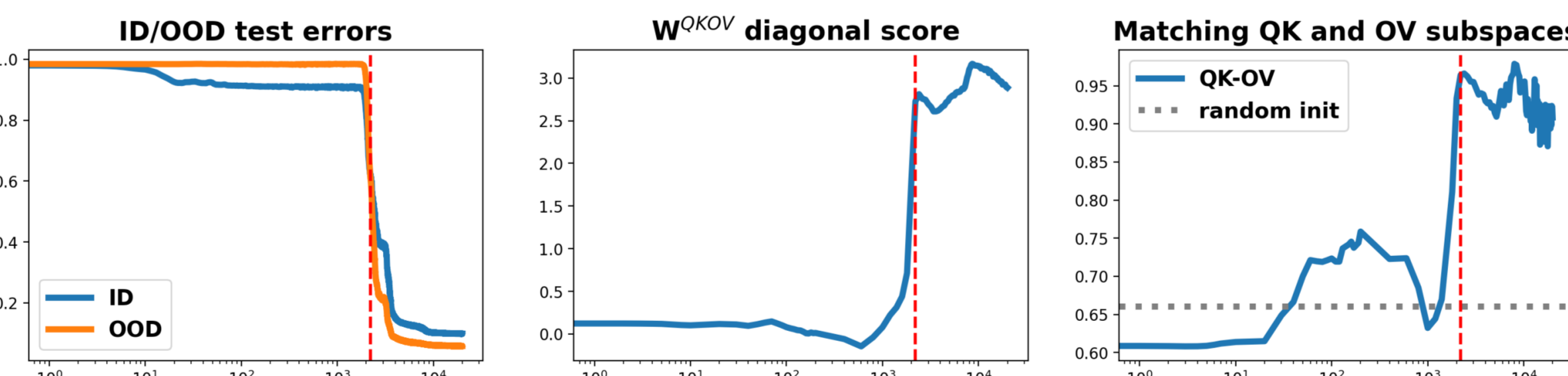
## Synthetic example: training dynamics on copying task

Copying task : ... [A], [B], [C] ... [A], [B]  $\xrightarrow{\text{next-token prediction}}$  ... [A], [B], [C] ... [A], [B], [C]



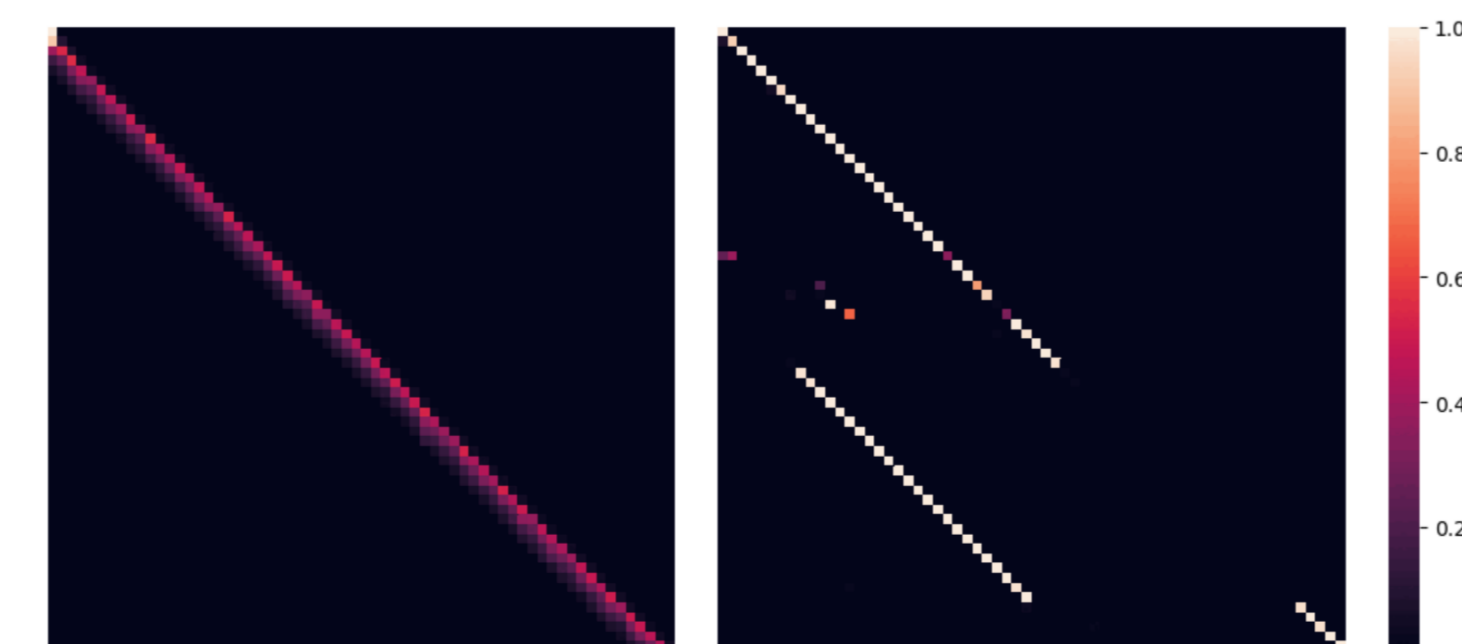
1. **Training data generation.** Vocabulary size 64, context length 64, i.i.d. tokens from power law distribution. Segment  $s^\#$  of random tokens with length  $\text{Unif}(\{10, 11, \dots, 19\})$ . Two copies of  $s^\#$  at random non-overlapping locations. Prompt format  $(*, s^\#, *, s^\#, *)$ .
2. **OOD data generation.** Token distribution uniform, segment length 25.
3. **Model.** 2-layer (1-layer) 1-head TF with no FFN, LayerNorm, RoPE, dropout.
4. **Training.** Fresh samples, autoregressive, AdamW.

## Low-dimensional subspace matching emerges abruptly



- **Diagonal score:** normalized average diagonal entries of  $\mathbf{W}_{\text{QK}}^{(2)}\mathbf{W}_{\text{OV}}^{(1)}$ .
- **Subspace matching:** generalized cosine sim between two principal subspaces ( $r = 10$ ).
- **Previous-token head (PTH) and induction head (IH).** Two types of attention heads. Follow similar sharp transition, complementary role (position vs. token matching).

PTH/IH attention: pool size None, step 20000



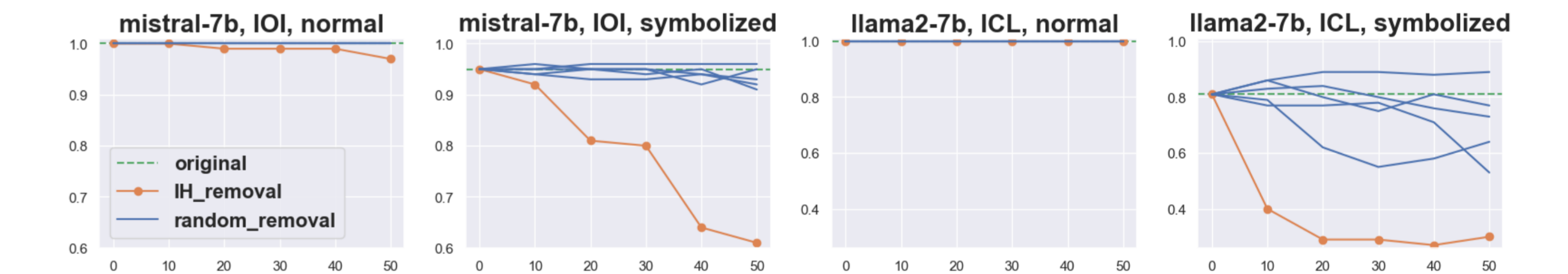
$$\text{score}^{\text{PTH}} = \text{Ave}_{i \leq N} \left( \frac{1}{T-1} \sum_{T \geq t \geq 2} (\mathbf{A}_i)_{t,t-1} \right)$$

$$\text{score}^{\text{IH}} = \text{Ave}_{i \leq N} \left( \frac{1}{|\mathcal{I}_i|} \sum_{t \in \mathcal{I}_i} (\mathbf{A}_i)_{t,t-L_i+1} \right)$$

$\mathcal{I}_i$  : index set of repeating tokens  
 $L_i$  : distance between two segments

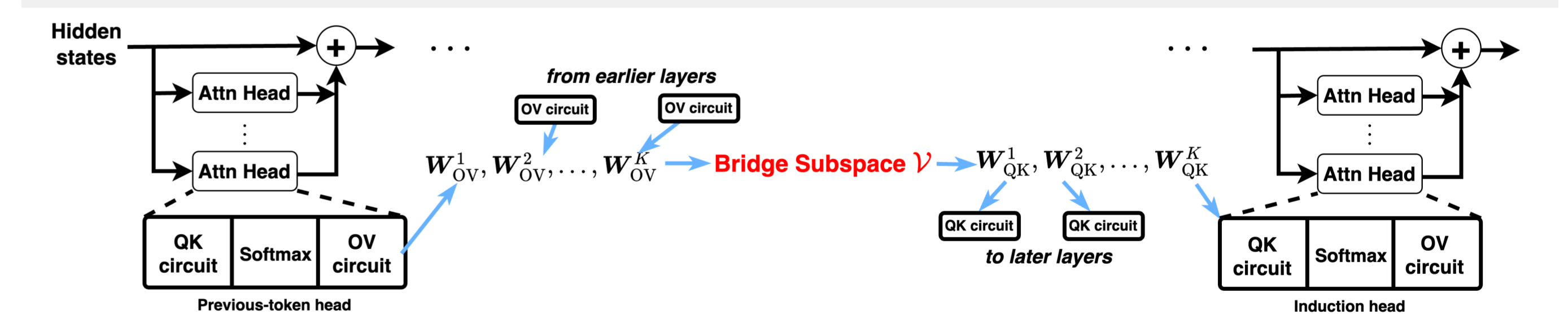
## Experiments on pretrained LLMs: symbolic & reasoning tasks

1. **Indirect object identification (IOI).** Normal (N) vs. Symbolized (S).
    - (N) “Then, Henry and Blake had a long argument. Afterwards Henry said to”  $\rightarrow$  Blake
    - (S) “Then, & and #\\$ had a long argument. Afterwards & said to”  $\rightarrow$  #\\$
  2. **In-context learning (ICL).**
    - (N) “baseball is sport, celery is plant, sheep is animal, volleyball is sport, lettuce is”  $\rightarrow$  plant
    - (S) “baseball is \\$#, celery is !%, sheep is !\*, volleyball is \\$#, lettuce is”  $\rightarrow$  !%
  3. **Math reasoning** with chain-of-thought (CoT) on GSM8K.
    - “Jerry is cutting up wood for his wood-burning stove. Each pine tree makes 80 logs, each maple tree makes 60 logs, and each walnut tree makes 100 logs. If Jerry cuts up 8 pine trees, 3 maple trees, and 4 walnut trees, how many logs does he get?” [...Deduction...] “### 1220”
- 100 test prompts, two versions (Normal as in-distribution, Symbolized as OOD).
  - Removal top- $K$  induction heads (ranked by attention scores) vs. removal of random heads,  $K = 0, 10, \dots, 50$ .



1. **Finding 1:** Normal prompts are insensitive to IH removal (likely memorization)
2. **Finding 2:** In contrast, OOD/reasoning prompts accuracy rely crucially on IHs as a component in composition.
3. **Verified extensively:** similar phenomena on 10+ LLMs, ranging from 36M to 70B.

## Common bridge representation hypothesis



**Hypothesis:** For a compositional task, there exists a low-dimensional subspace  $\mathcal{V} \subset \mathbb{R}^d$  s.t.

$$\mathcal{V} = \text{span}(\mathbf{W}_{\text{OV},j}) = \text{span}(\mathbf{W}_{\text{QK},k}^\top).$$

- Extension of linear representation hypothesis to compositional tasks.
- Key to OOD generalization.
- Supported by ablation experiments (projecting weights onto  $\mathcal{V}$  vs. onto  $\mathcal{V}^\perp$ )

## References

[1] Jiajun Song, Zhuoyan Xu, and Yiqiao Zhong. Out-of-distribution generalization via composition: a lens through induction heads in transformers. *Proceedings of National Academy of Sciences*, 122(6), 2025.