

DiffRACT

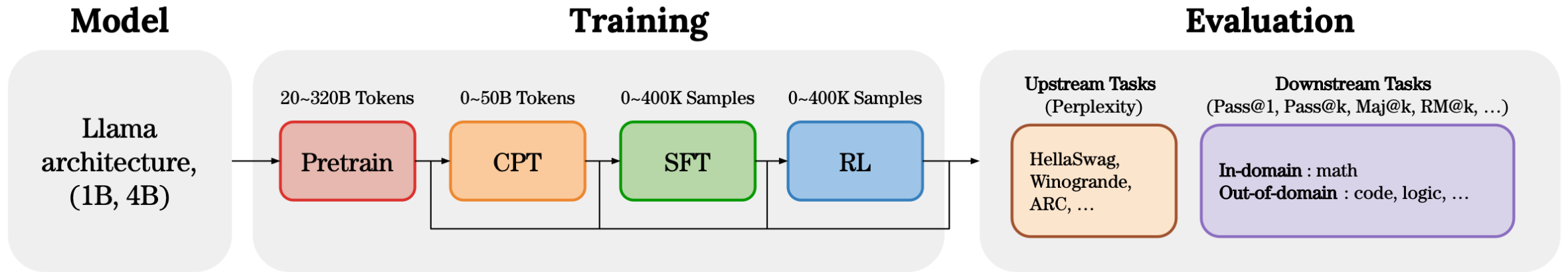
Spectral View of LLM Domain Adaptation

Nikita Borodin¹, Maria Krylova¹, Artem Zabolotnyi^{1,2}, Dmitry Aspisov¹, Egor Shikov¹,
Nikita Tyuplyaev¹, Oleg Travkin¹, Roman Alferov¹, Dmitry Vinichenko¹

¹ Risk AI Research Lab · ² Applied AI Institute

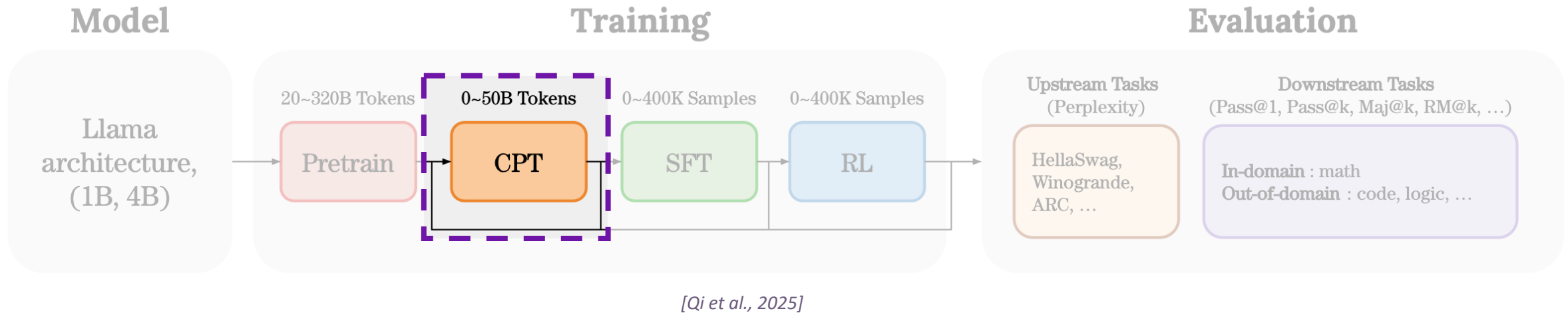


Continual pre-training – important stage for LLM domain adaptation



[Qi et al., 2025]

Continual pre-training – important stage for LLM domain adaptation



CENTRAL IN MODERN PIPELINES

Produces domain-specialized models like **DeepSeek-Math** and **LLEMMA**



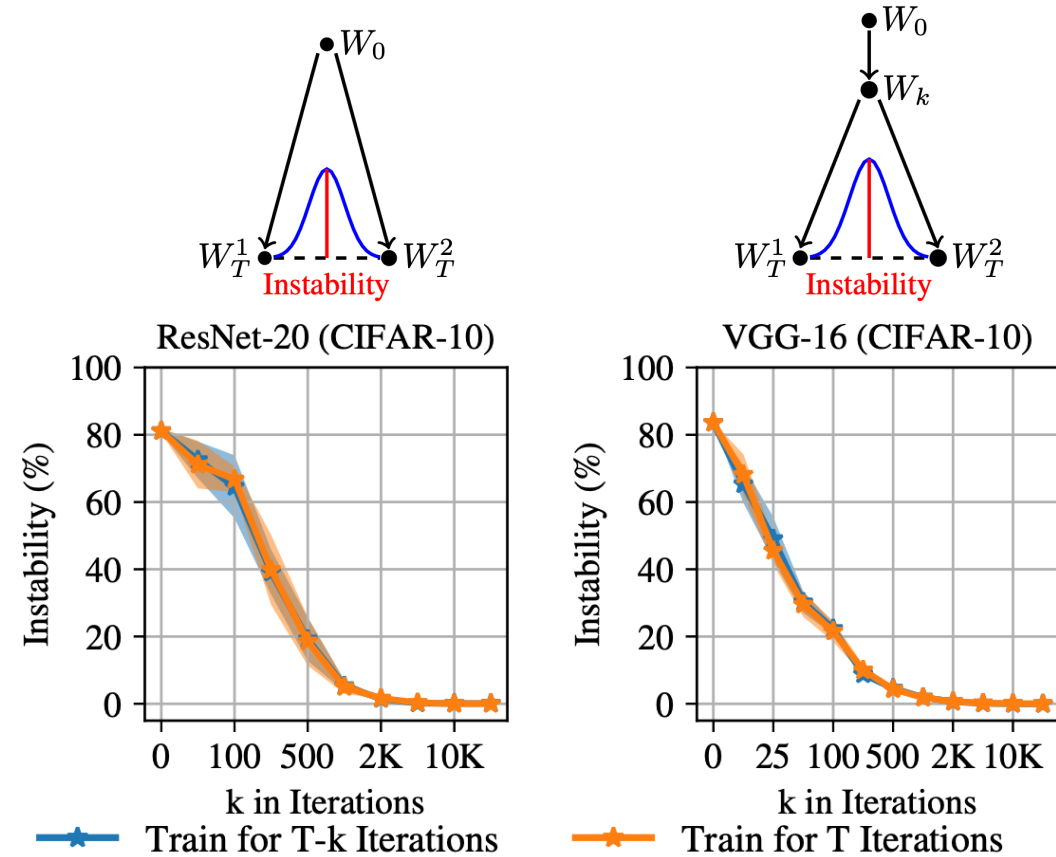
LLEMMA ■

A wide range of phenomena and properties is known for SFT stage

1

Linear mode connectivity

Linear Mode Connectivity



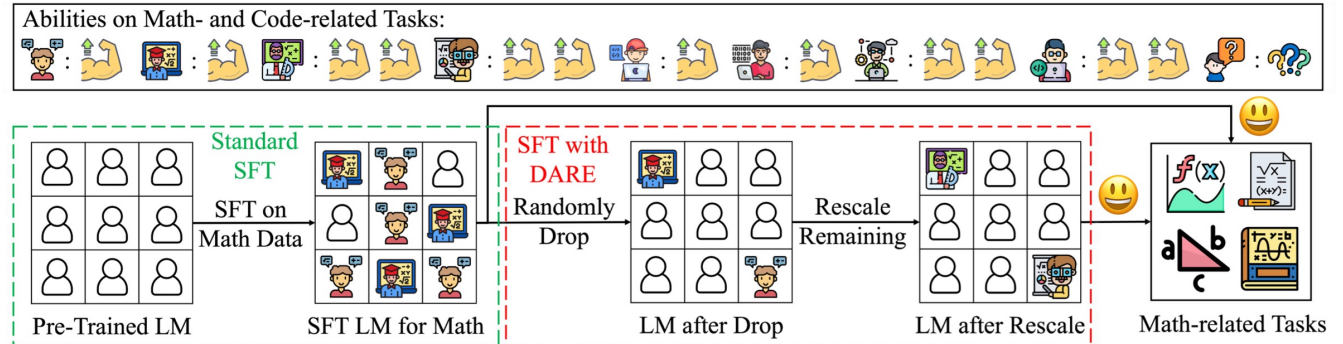
[Frankle et al., 2020]

A wide range of phenomena and properties is known for SFT stage

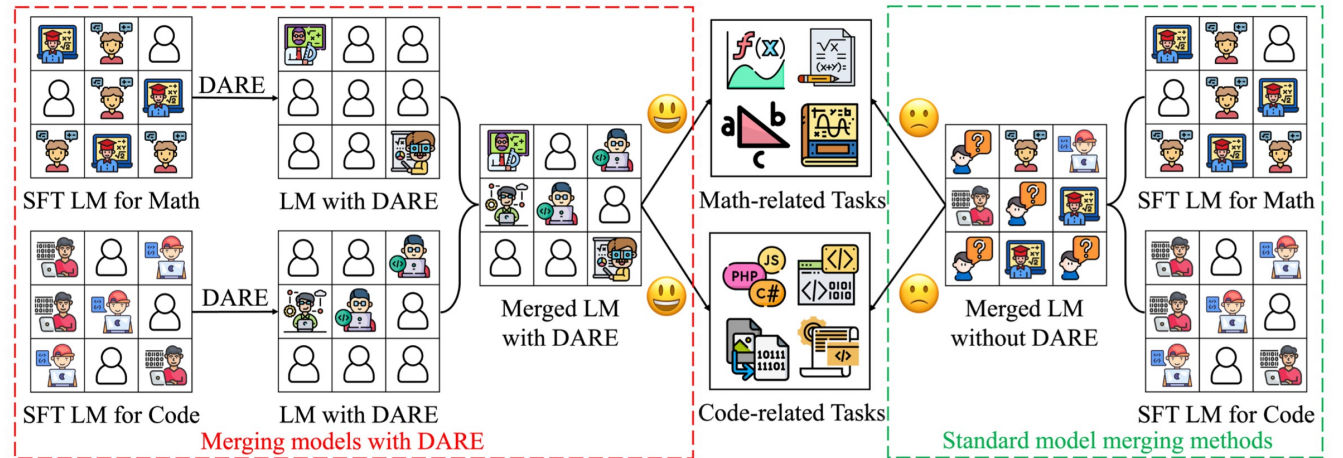
1 Linear mode connectivity

2 Task arithmetic, ability transfer

Language Models are Super Mario



(a) Standard SFT and SFT with DARE on math-related task.



(b) Merging models with and without DARE on math- and code-related tasks.

[Yu et al., 2024]

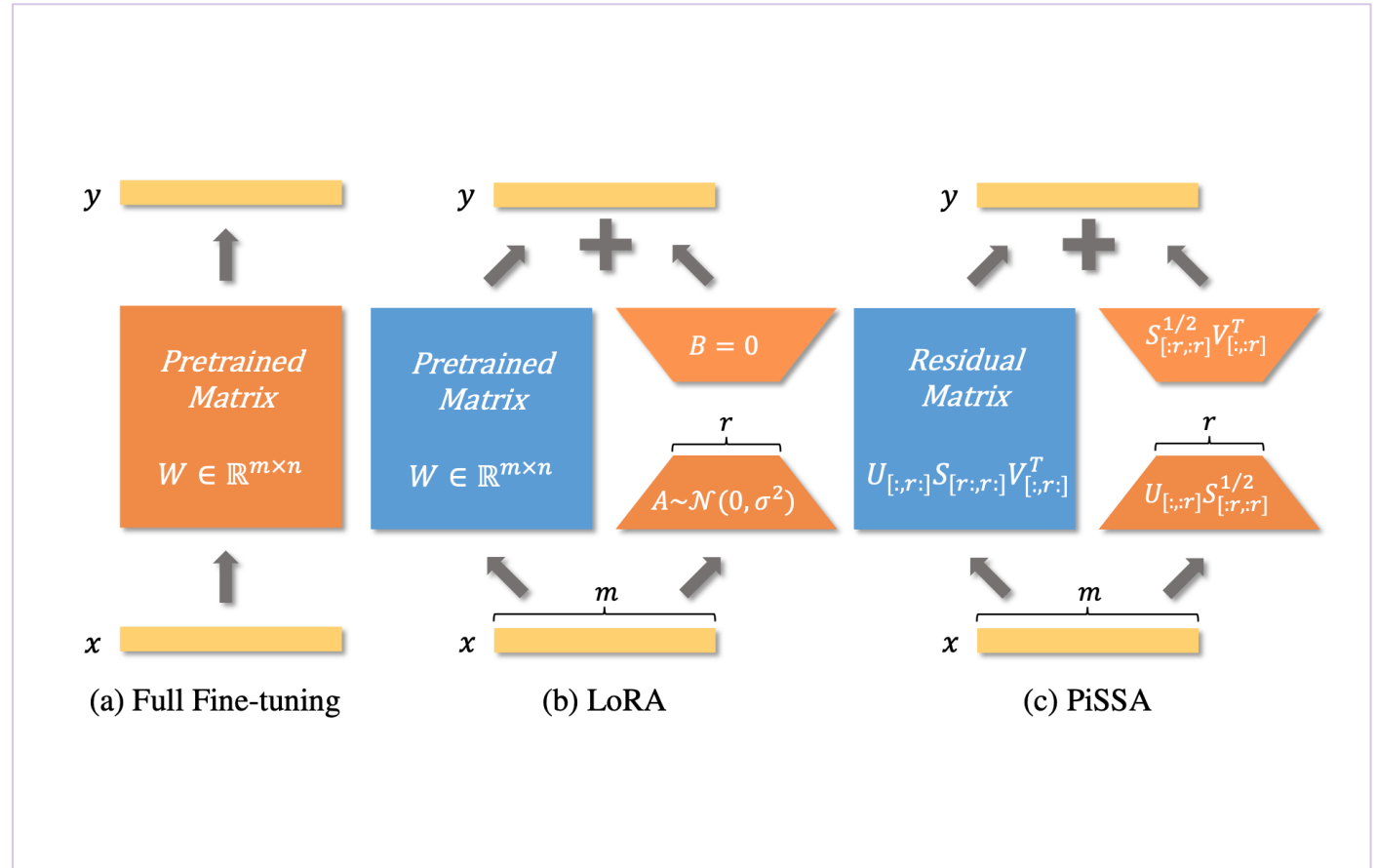
A wide range of phenomena and properties is known for SFT stage

1 Linear mode connectivity

2 Task arithmetic, ability transfer

3 Low-rank subspace modification

PiSSA



[Meng et al., 2024]

Objective: Investigate the CPT stage and identify its properties

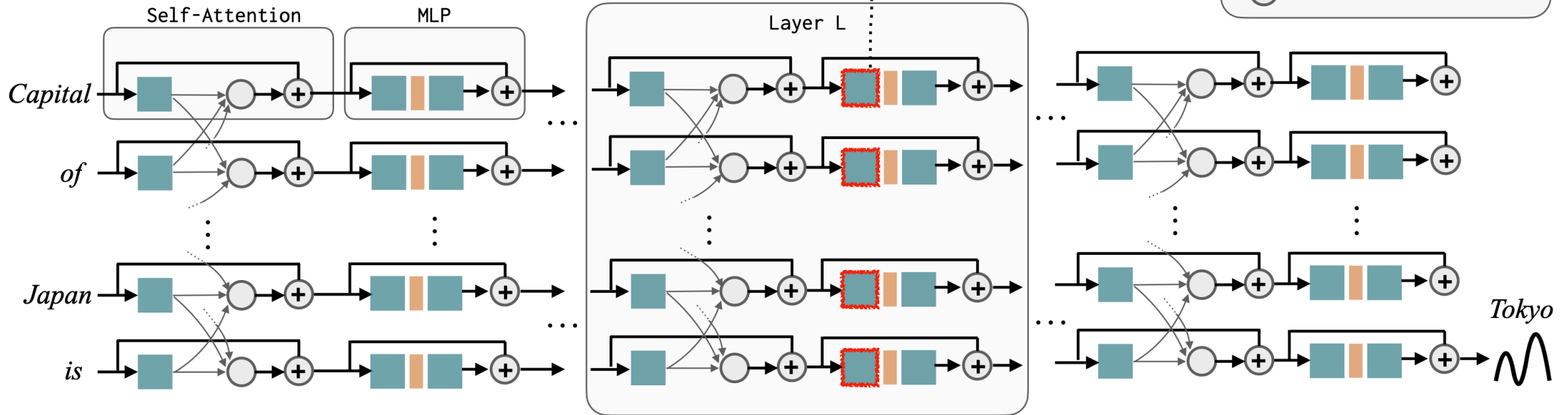
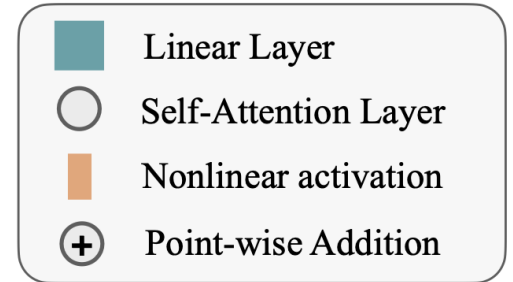
Investigation of weight matrices spectra can lead to new insights

LASER

LASER 

Replace W in specific layers with its low rank approximation W_{LR}

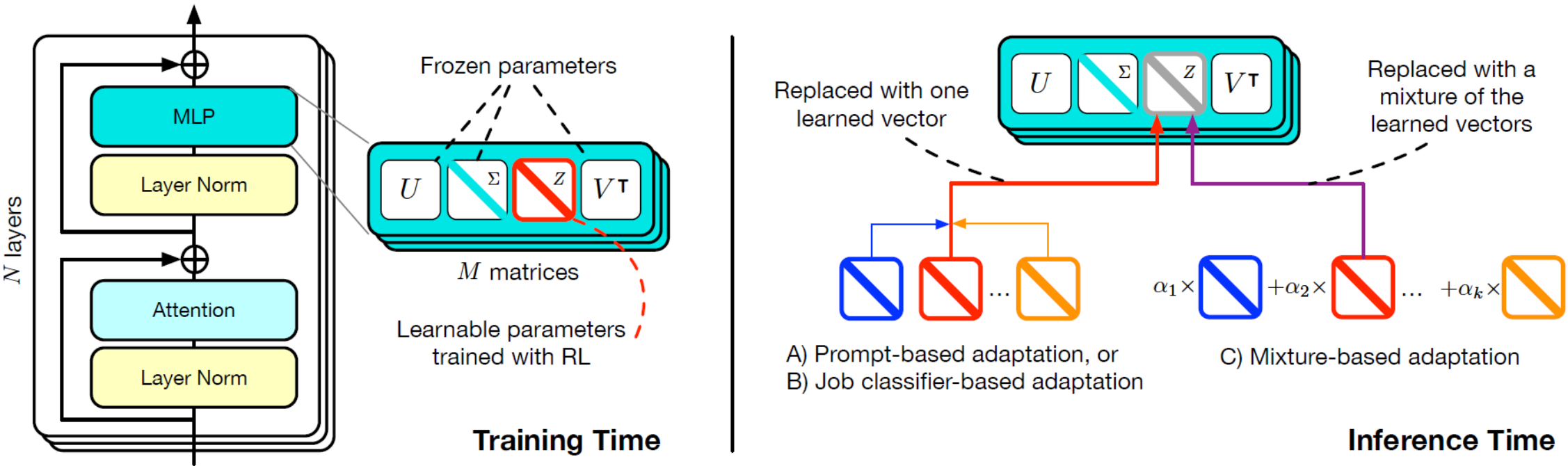
$$W \rightarrow W_{LR} = U \Sigma V^T$$



[Sharma et al., 2024]

Investigation of weight matrices spectra can lead to new insights

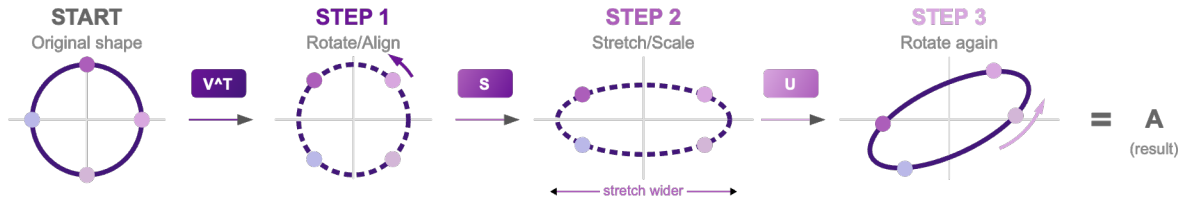
Transformer-Squared



[Sun et al., 2025]

Investigation of weight matrices spectra can lead to new insights

Heavy-tailed self-regularization theory

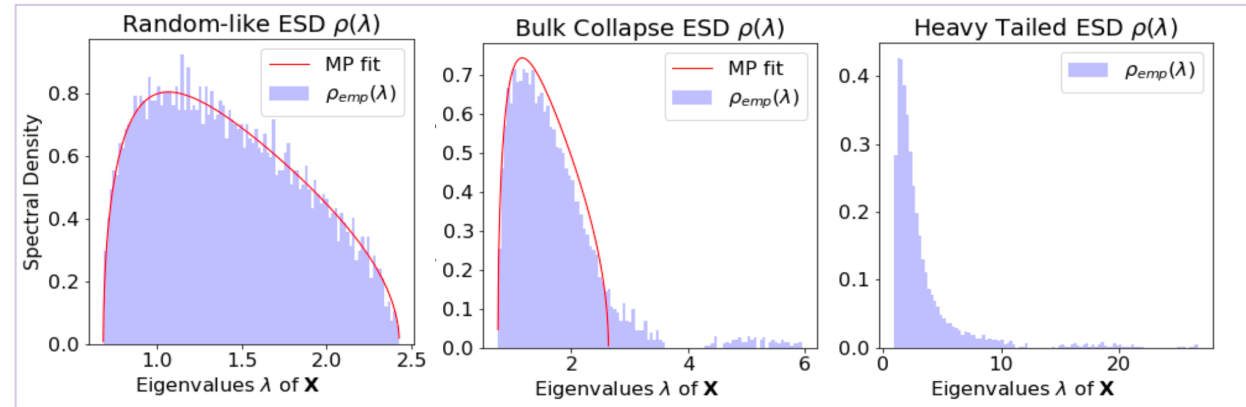


Marchenko-Pastur distribution for GOE random matrices

$$W_{i,j} \sim \mathcal{N}(0, \text{Var}(\mathbf{W}))$$

$$\rho^{\text{MP}}(\lambda; \mathbf{X}) = \begin{cases} \frac{m}{2\pi n \cdot \text{Var}(\mathbf{W})} \frac{\sqrt{(\lambda^{\max} - \lambda)(\lambda - \lambda^{\min})}}{\lambda}, & \lambda \in [\lambda^{\min}, \lambda^{\max}], \\ 0, & \text{otherwise,} \end{cases}$$

$$\lambda^{\max/\min}(\mathbf{X}) = \text{Var}(\mathbf{W})(1 \pm \sqrt{n/m})^2$$



[Martin & Mahoney, 2021]

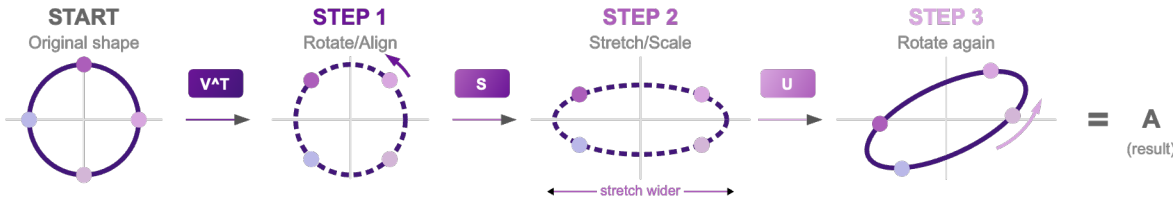
Power law tails develop as training progresses

$$\rho^{\text{PL}}(\lambda; \mathbf{X}) = c \cdot \lambda^{-\alpha}$$

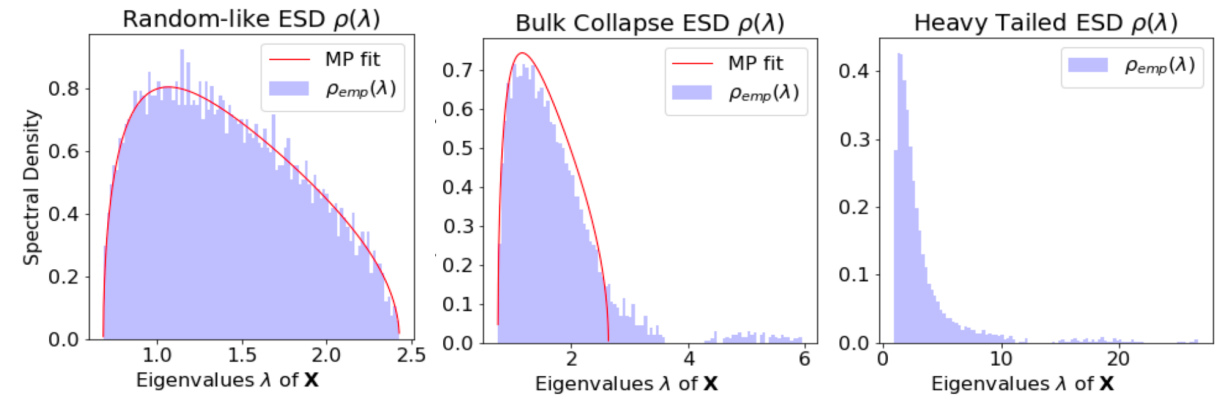
Hypothesis – power law parameter and other spectral properties correlate with model quality

Investigation of weight matrices spectra can lead to new insights

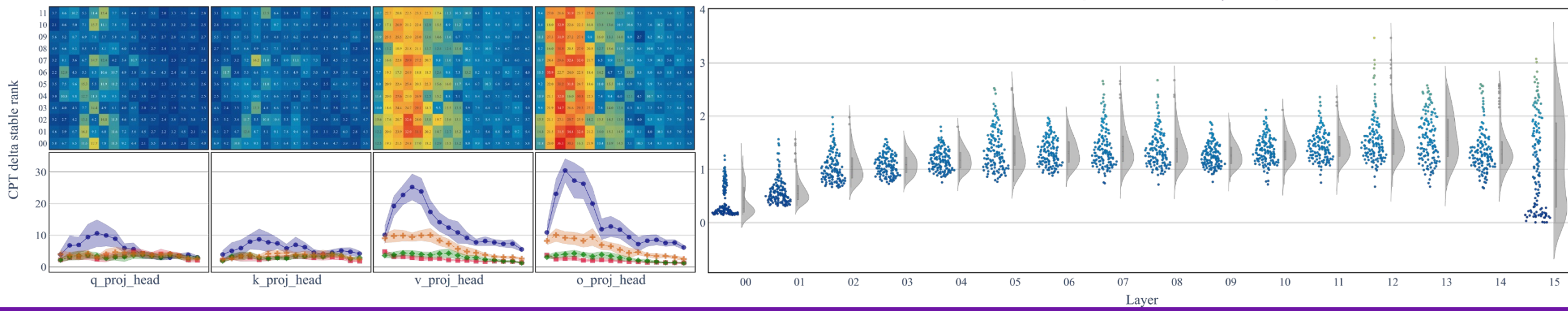
Heavy-tailed self-regularization theory



DIFFRACT · OPEN SOURCE PYTHON PACKAGE
<https://github.com/Risk-AI-Research/diffract>
SVD-based weight analysis for billion-parameter LLMs



[Martin & Mahoney, 2021]



We use DiffRACT package to investigate CPT stage through the spectral lens

Experimental setup: OLMo 2 pre-train + CPT



MODEL SIZE

PRE-TRAIN BUDGET

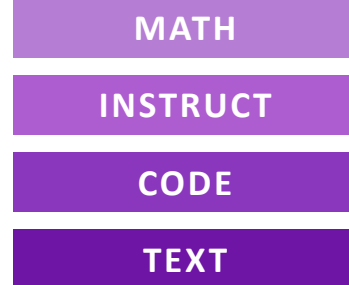
CPT DOMAINS

1B: 4T, 20 – 400B

7B: 4T

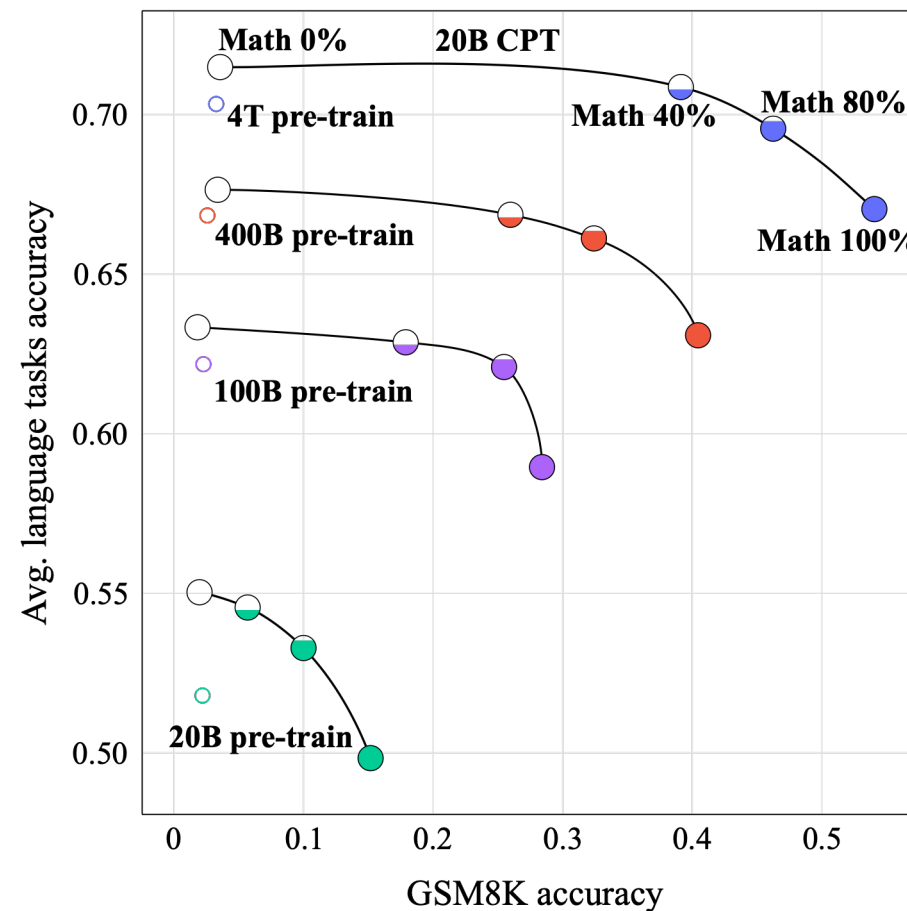
13B: 5T

32B: 6T



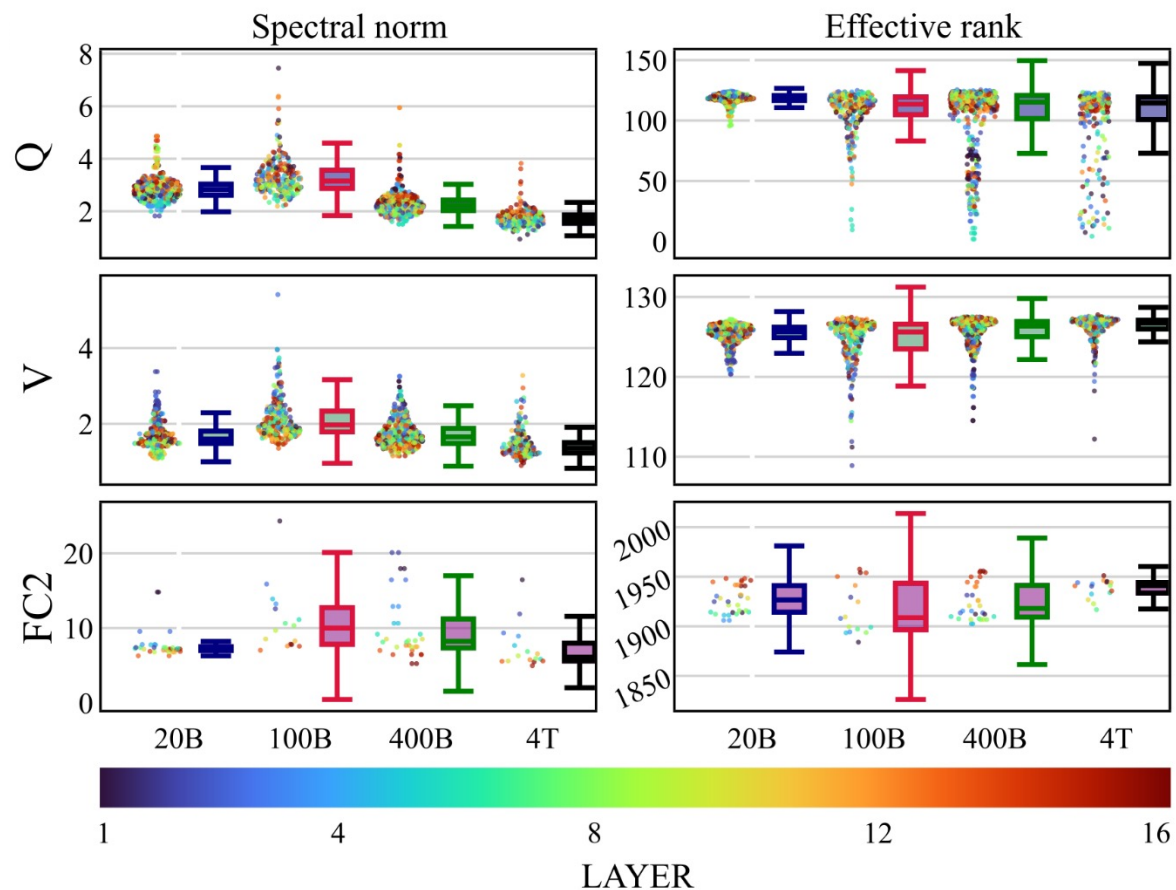
OLMo-2 CPT checkpoints

Impact of pre-train budget on math CPT quality



Takeaway: longer pre-train enables higher quality after CPT

Spectral evolution along the pre-train stage and CPT spectral properties



1B model

Spectral properties of attention heads:

Frobenius norm $\|\mathbf{W}\|_F = \sqrt{\sum_i \sigma_i(\mathbf{W})^2}$

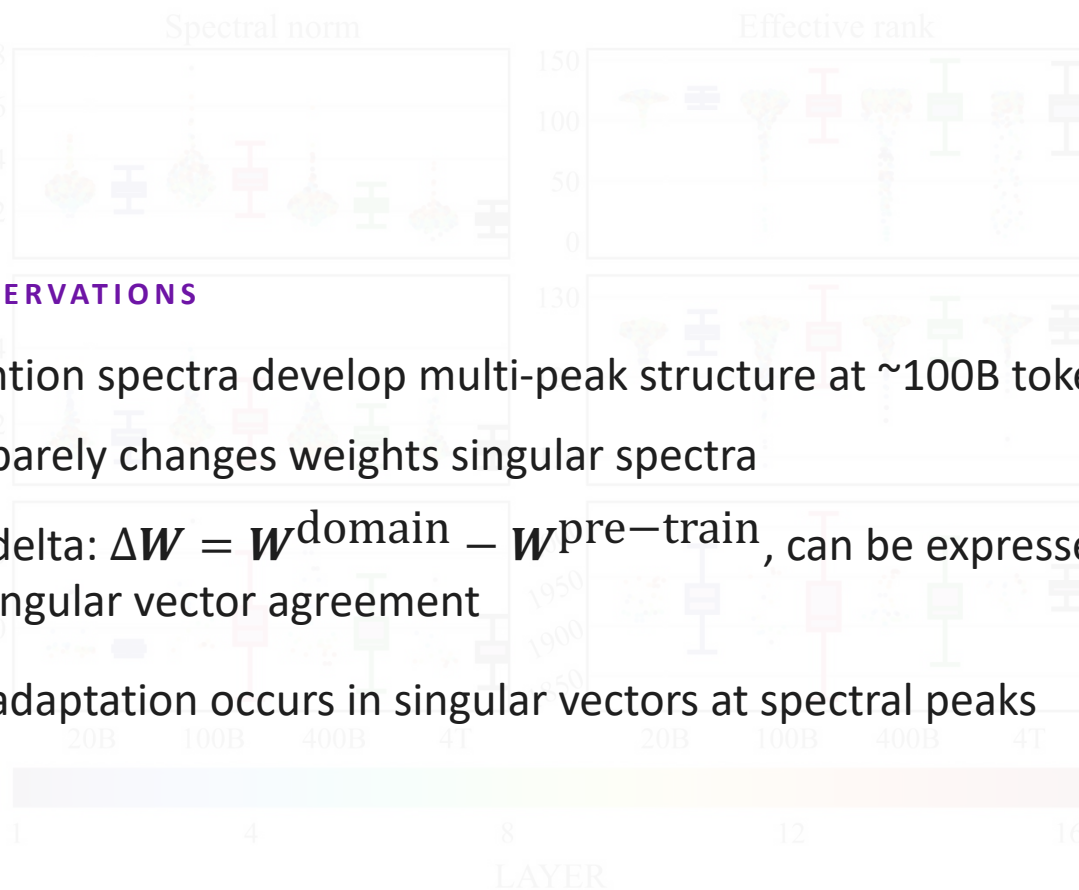
Spectral norm $\|\mathbf{W}\|_2 = \max_i \sigma_i(\mathbf{W})$

Stable rank $R^s(\mathbf{W}) = \frac{\|\mathbf{W}\|_F^2}{\|\mathbf{W}\|_2^2}$

Effective rank

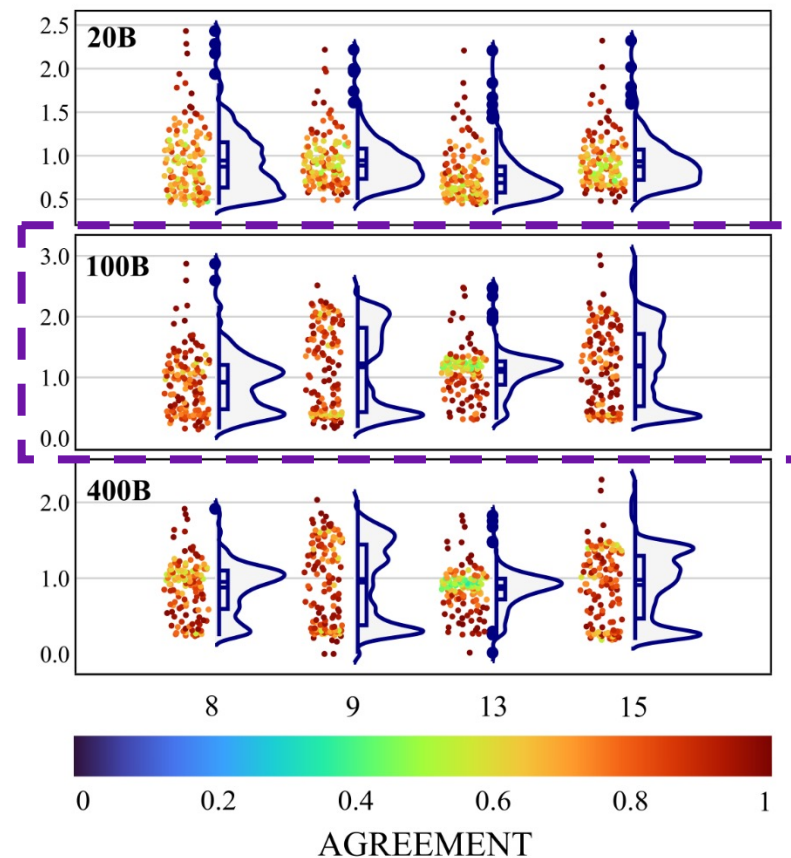
$$R^e(\mathbf{W}) = \exp\left(-\sum_{i=1}^{R(\mathbf{W})} \frac{\sigma_i(\mathbf{W})}{\sum_j \sigma_j(\mathbf{W})} \log\left(\frac{\sigma_i(\mathbf{W})}{\sum_j \sigma_j(\mathbf{W})}\right)\right)$$

Spectral evolution along the pre-train stage and CPT spectral properties



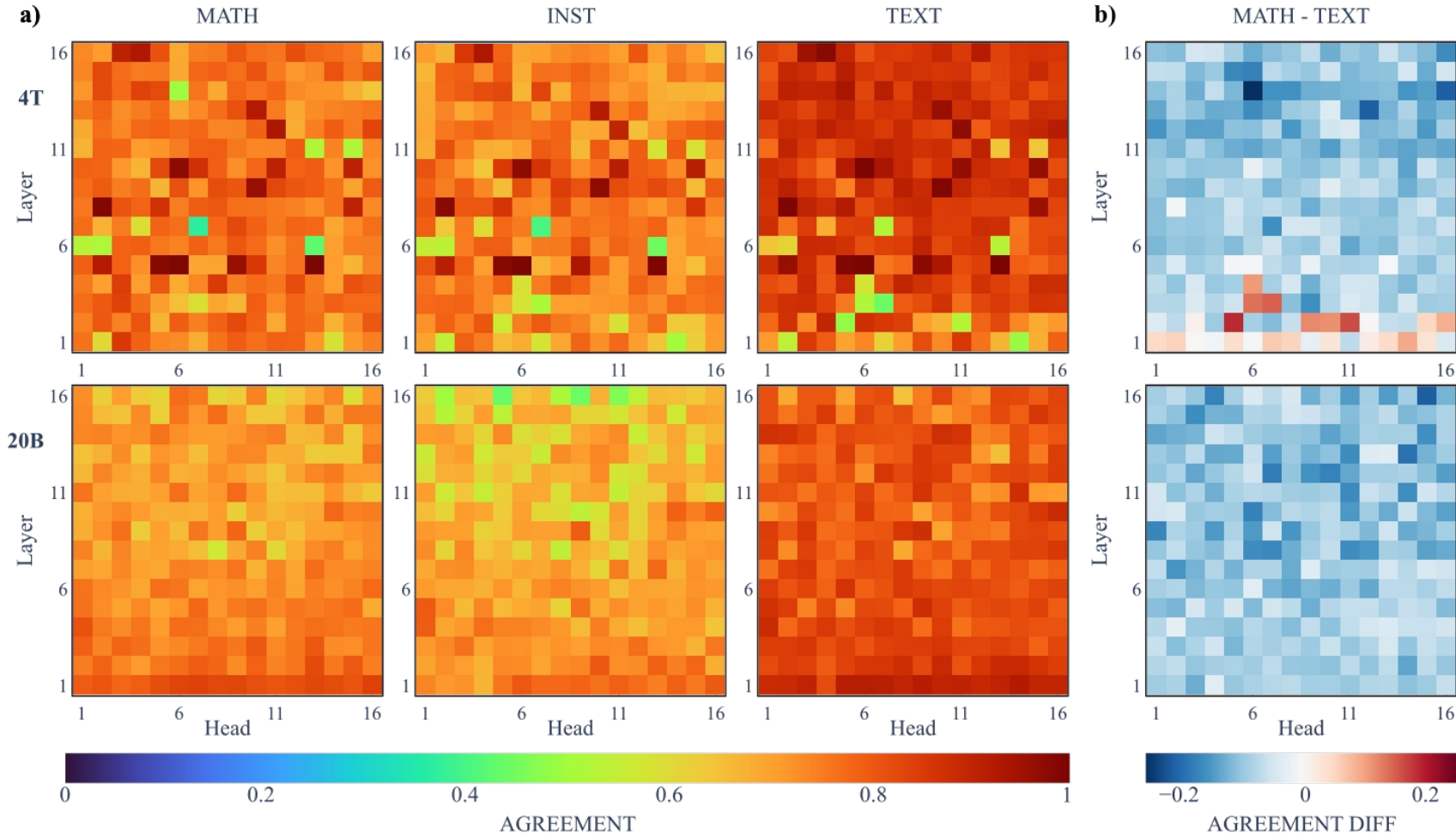
KEY OBSERVATIONS

- Attention spectra develop multi-peak structure at ~100B tokens
- CPT barely changes weights singular spectra
- CPT delta: $\Delta W = W^{\text{domain}} - W^{\text{pre-train}}$, can be expressed via singular vector agreement
- CPT adaptation occurs in singular vectors at spectral peaks

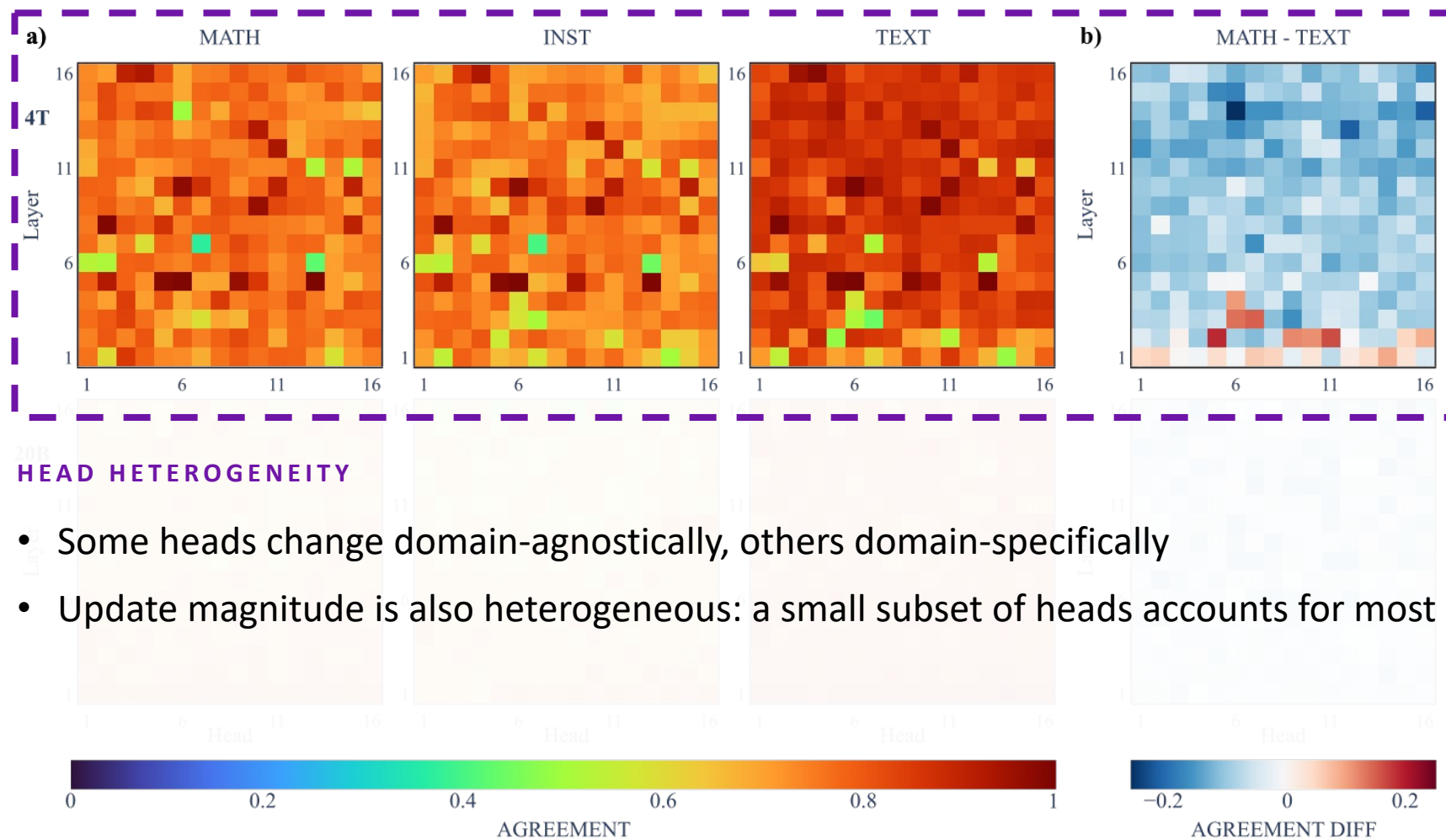


Takeaway: complex spectral shape enables finer CPT delta structure

Structure of CPT delta – head heterogeneity



Structure of CPT delta – head heterogeneity



HEAD HETEROGENEITY

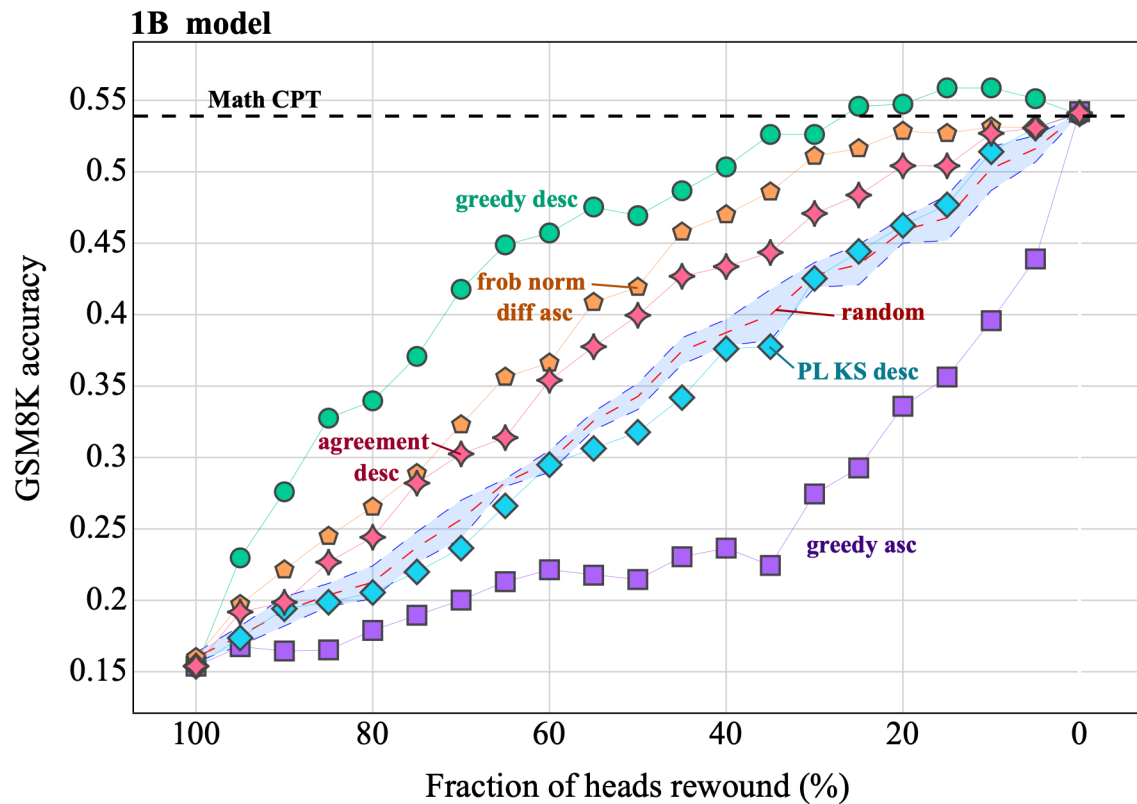
- Some heads change domain-agnostically, others domain-specifically
- Update magnitude is also heterogeneous: a small subset of heads accounts for most of the change

Takeaway: CPT changes are concentrated and domain-specific

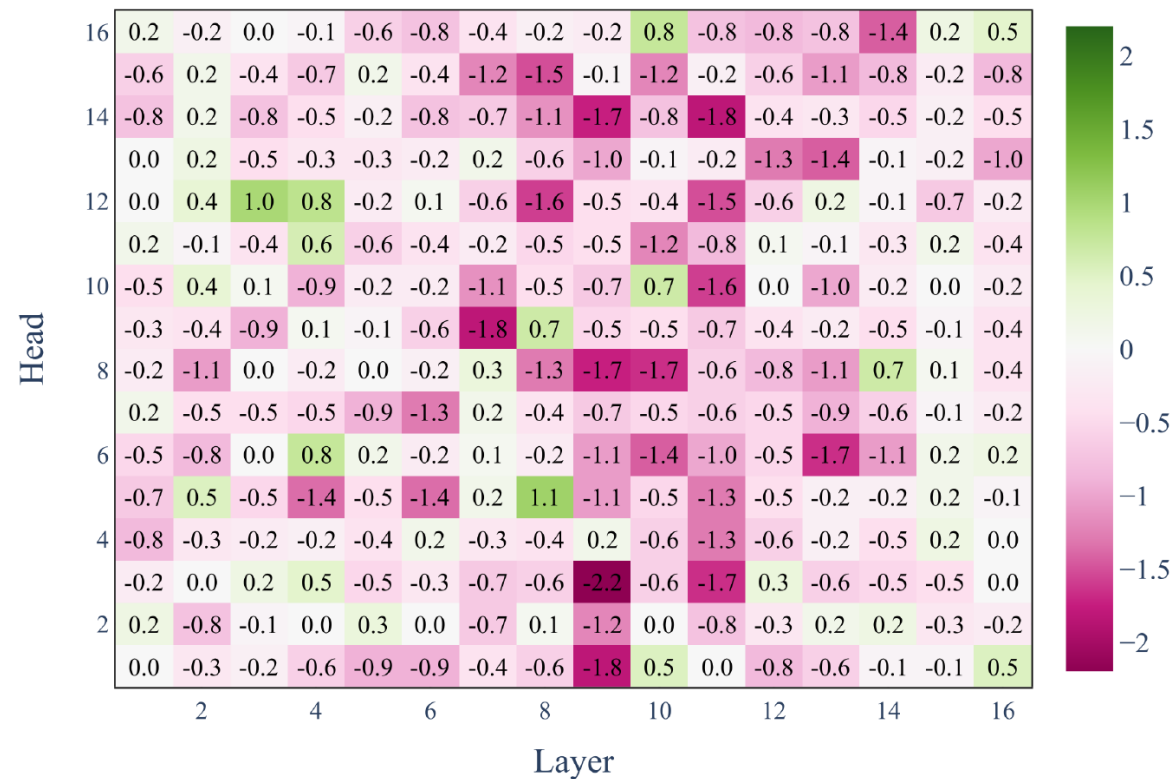
Structure of CPT delta – head-wise rewind

“Greedy” rewind pattern

20B math CPT after 4T pre-train



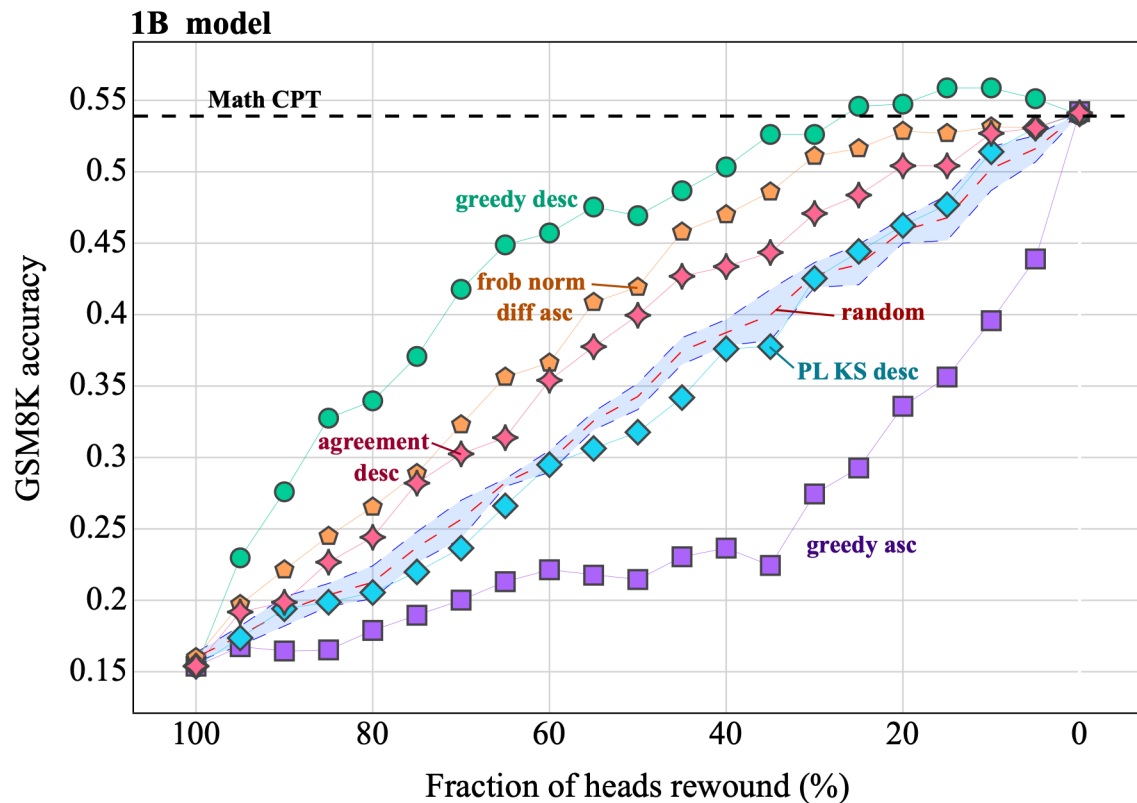
Heads are ranked by the effect of individual rewind on model GSM8K quality



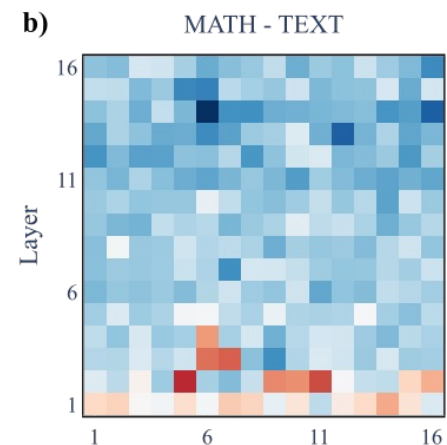
Structure of CPT delta – head-wise rewind

“Reference domain” rewind pattern

20B math CPT after 4T pre-train



Heads are ranked by difference between math and text CPTs

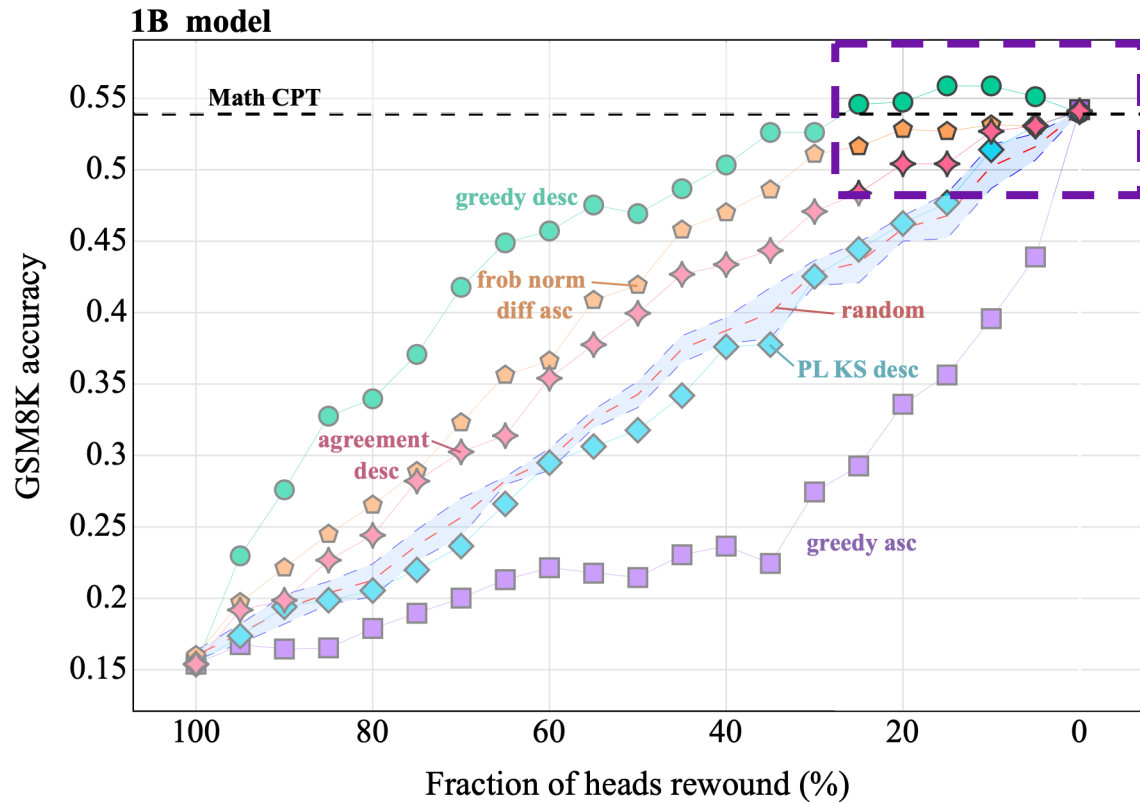


$$s_{l,h} = \text{scale}_{[0,1]} \left(\left\| \mathbf{W}_{l,h}^{\text{domain}} - \mathbf{W}_{l,h}^{\text{pre-train}} \right\|_F \right) - \text{scale}_{[0,1]} \left(\left\| \mathbf{W}_{l,h}^{\text{reference}} - \mathbf{W}_{l,h}^{\text{pre-train}} \right\|_F \right)$$

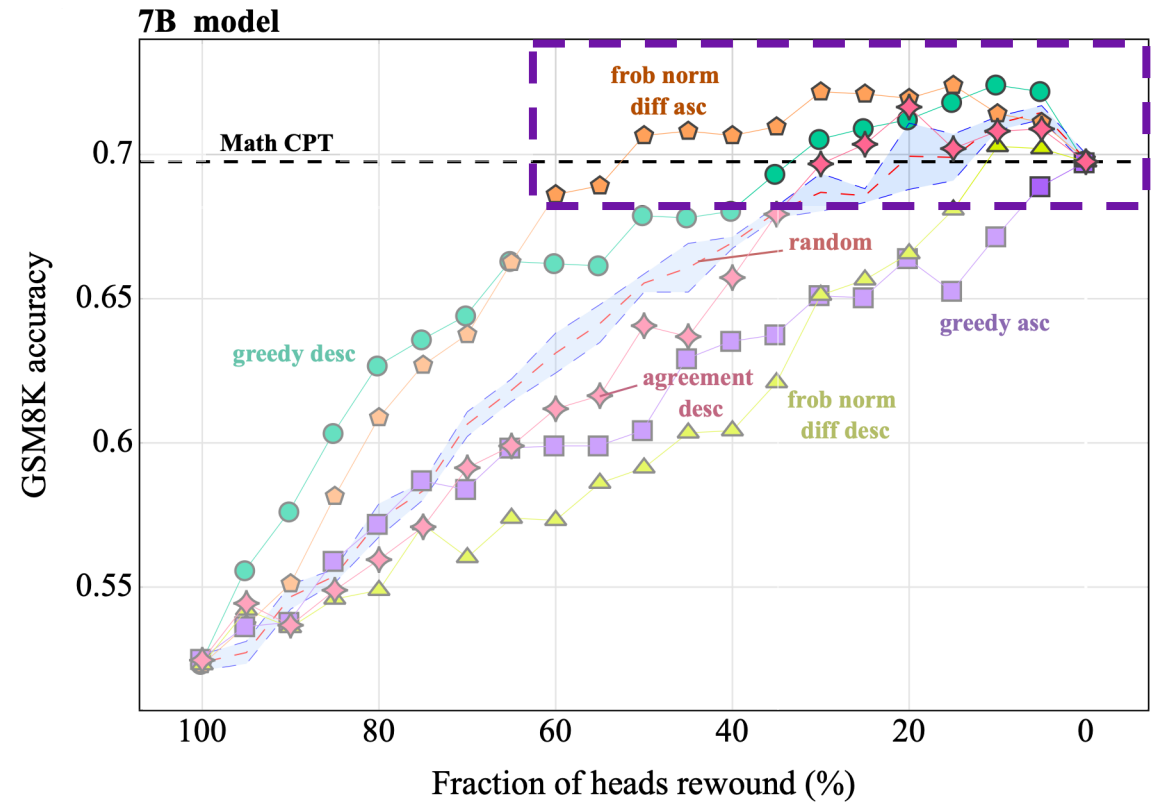
$$\text{scale}_{[0,1]}(\mathbf{X}) = (\mathbf{X} - X_{\min}) / (X_{\max} - X_{\min})$$

Structure of CPT delta – head-wise rewind

- Rewinding 15% heads yields **+2%** GSM8K
- Up to **25%** heads can be rewind without significant quality drop



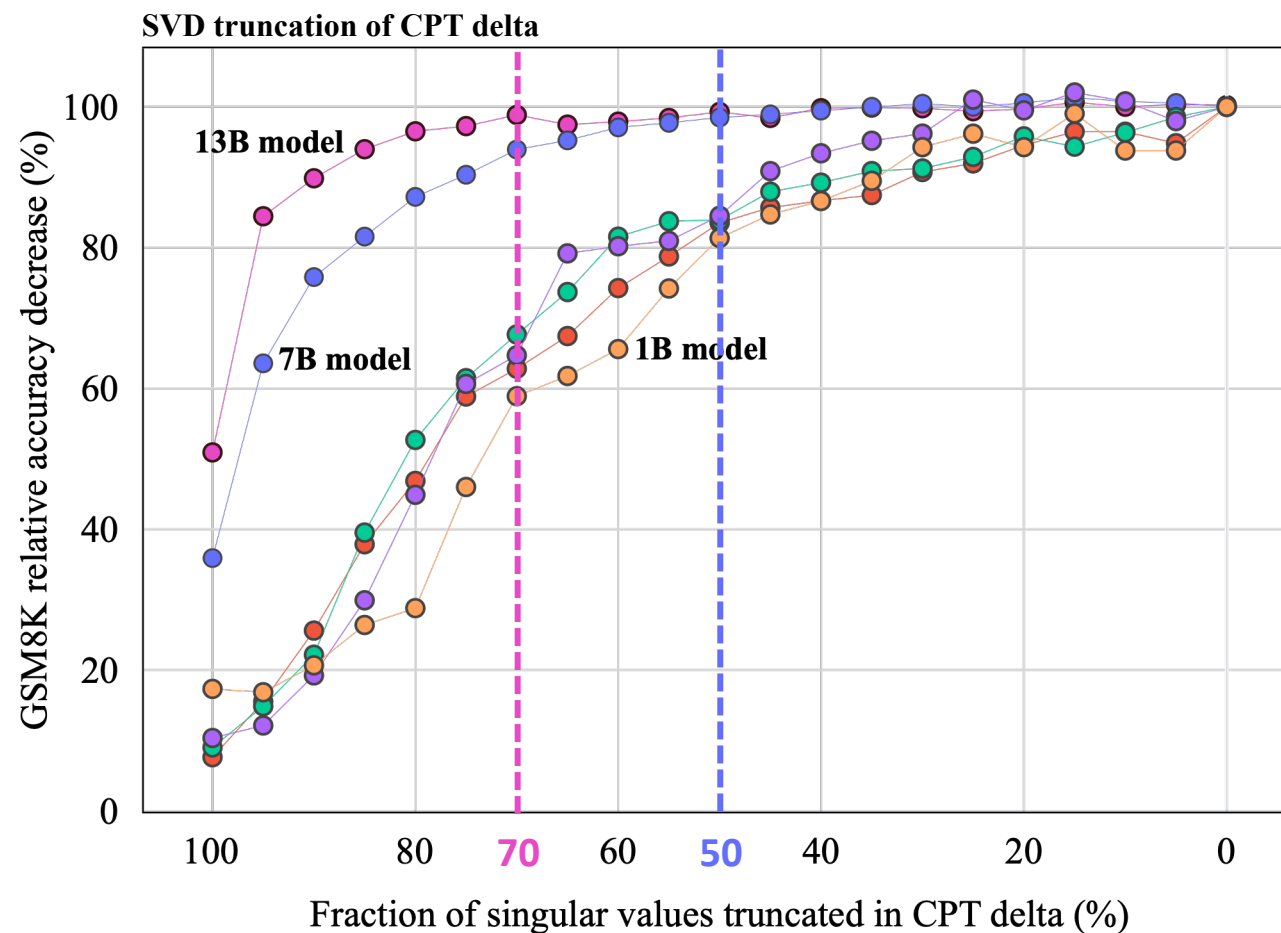
- Rewinding 15% heads yields **+4%** GSM8K
- Up to **60%** heads can be rewind without significant quality drop



Takeaway: head-wise rewind can consistently increase quality for larger models

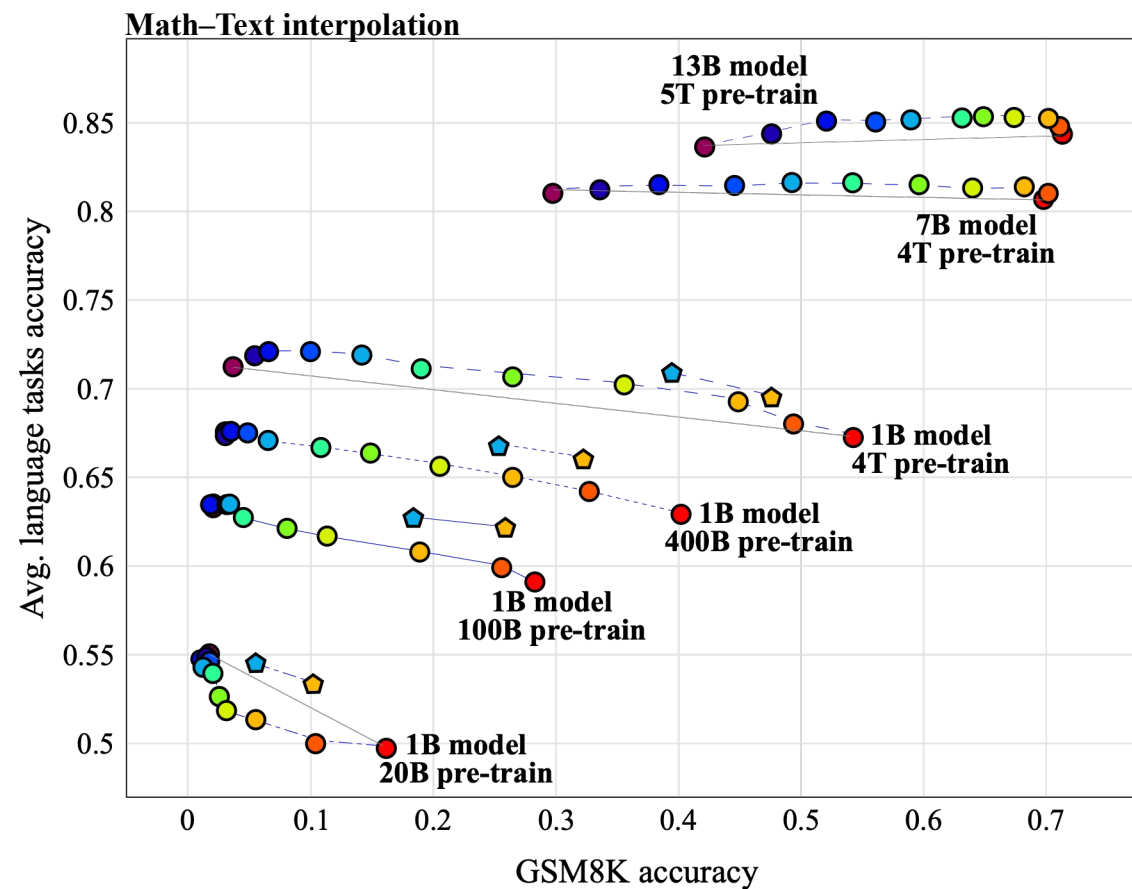
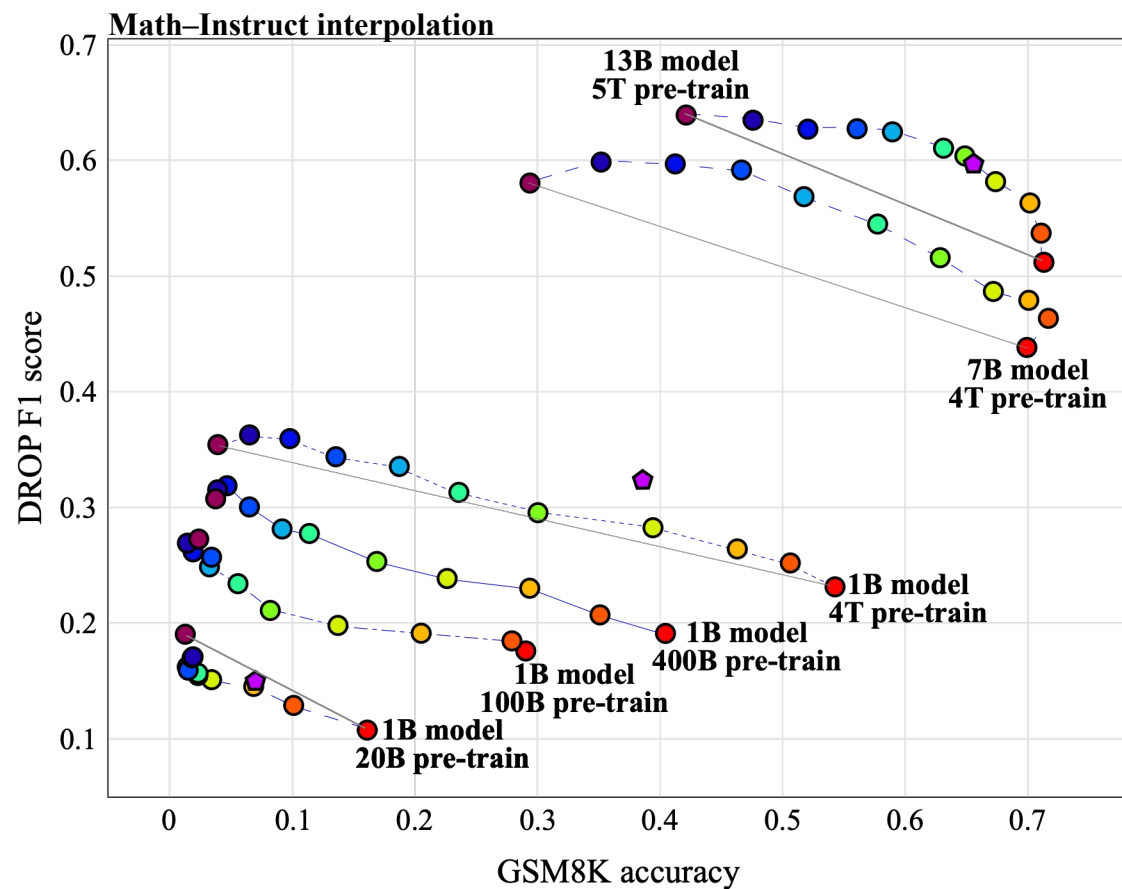
Structure of CPT delta – SVD truncation

- CPT delta is increasingly redundant with scale
- Larger model \Rightarrow lower-rank update
- 7B model: **50%** of singular values truncatable
- 13B model: up to **70%** truncatable

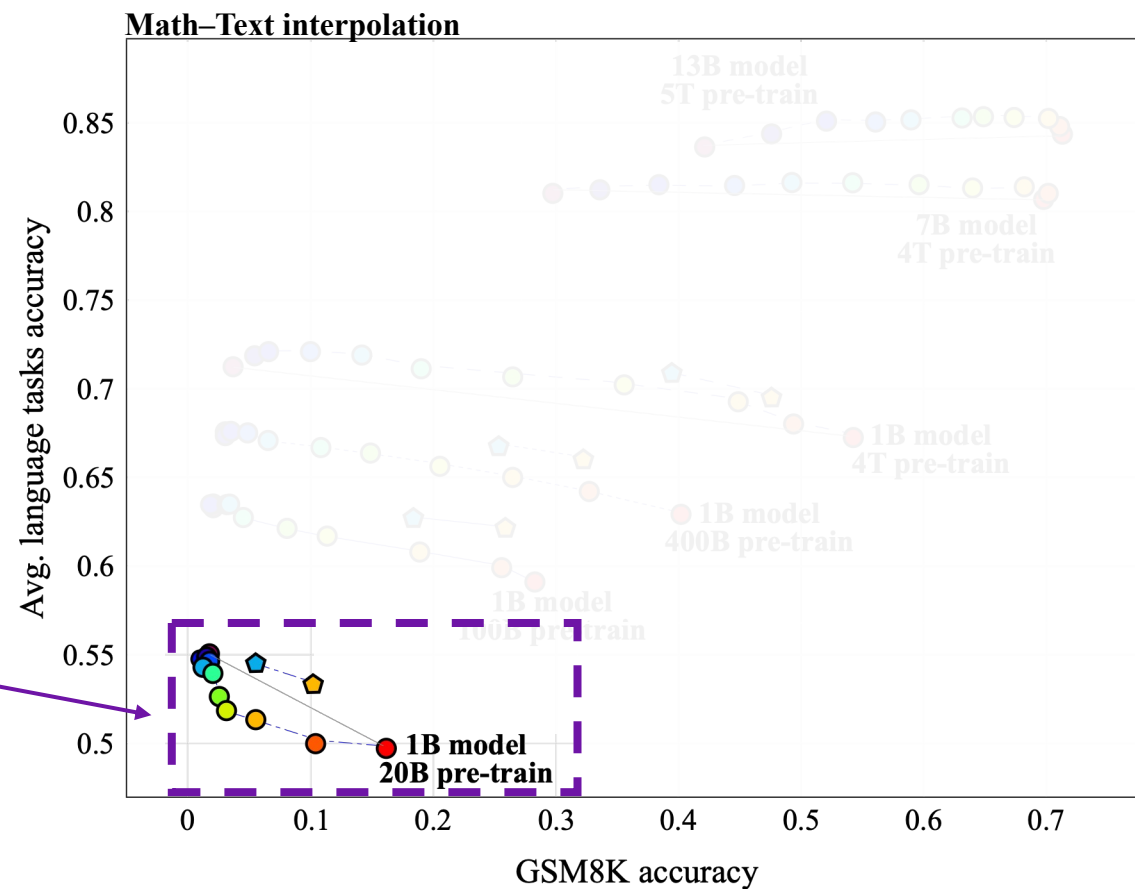
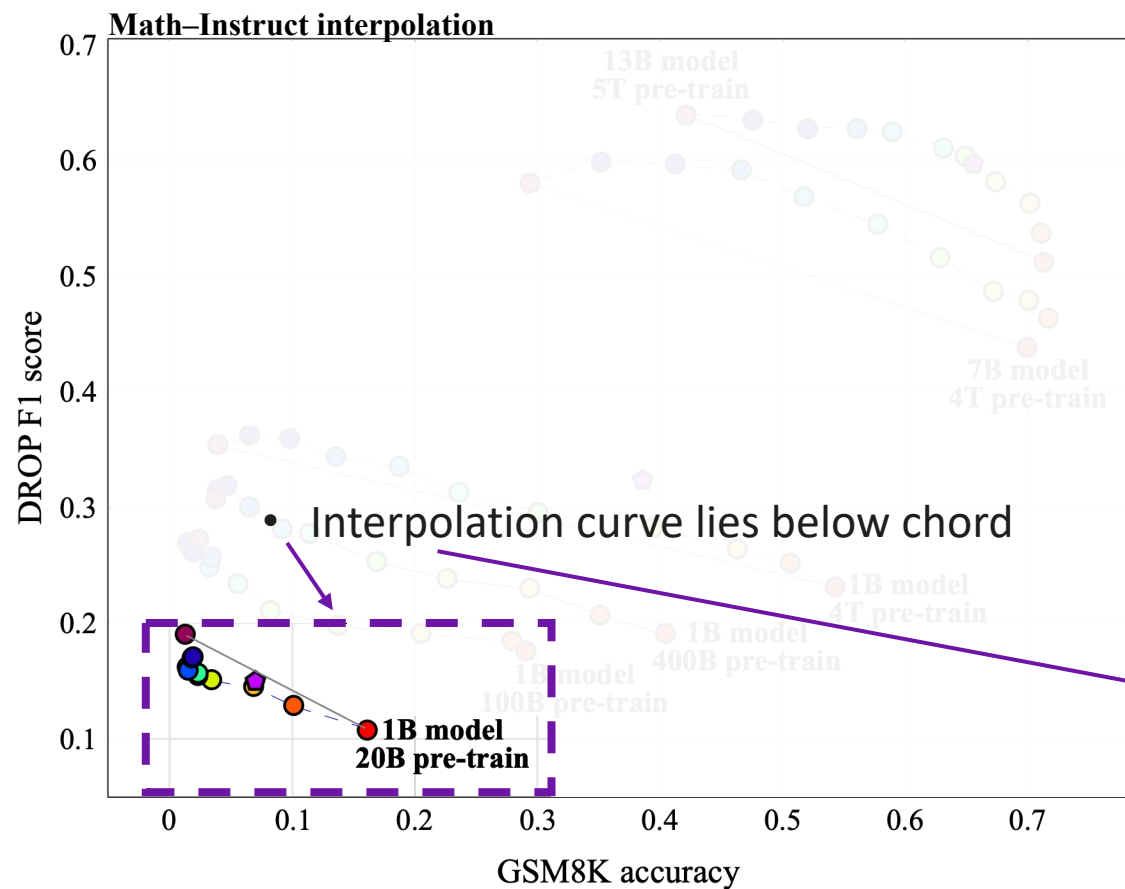


Takeaway: CPT update becomes more compressible for larger models, but not enough for DARE-like methods

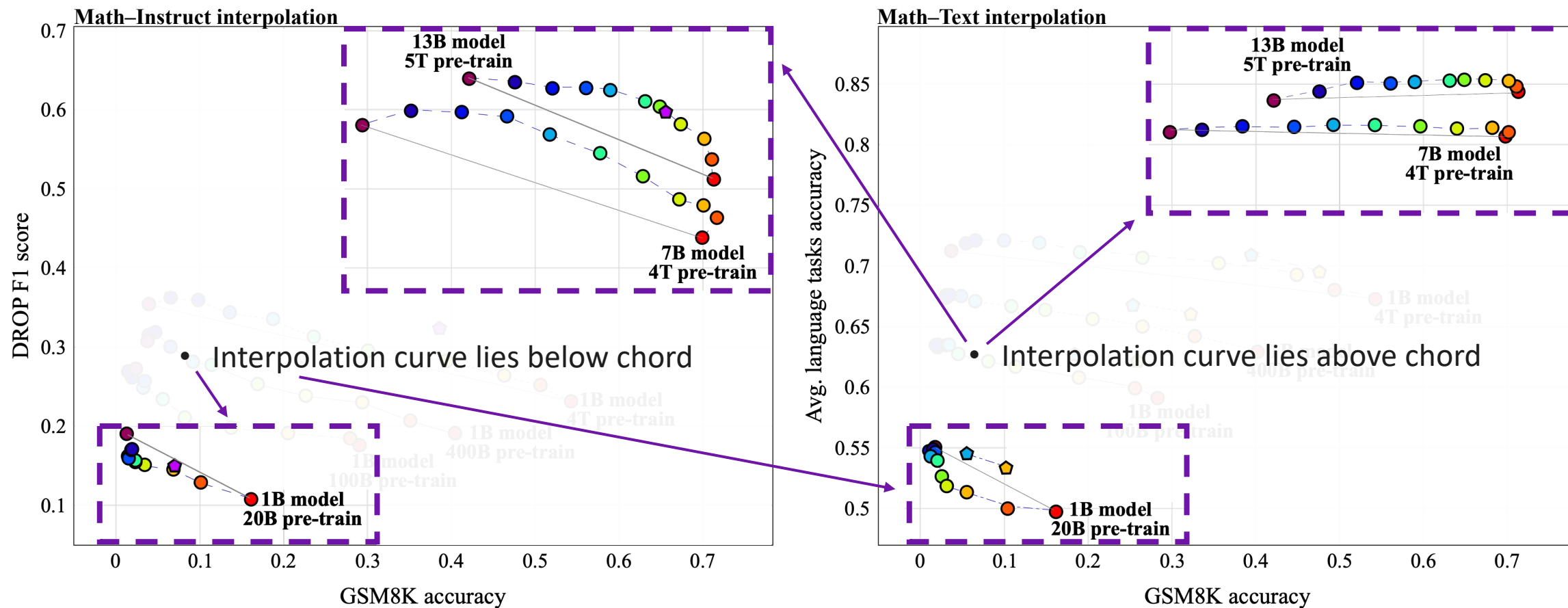
Structure of CPT delta – domain connectivity



Structure of CPT delta – domain connectivity



Structure of CPT delta – domain connectivity



Takeaway: larger pre-train budgets enable domain connectivity

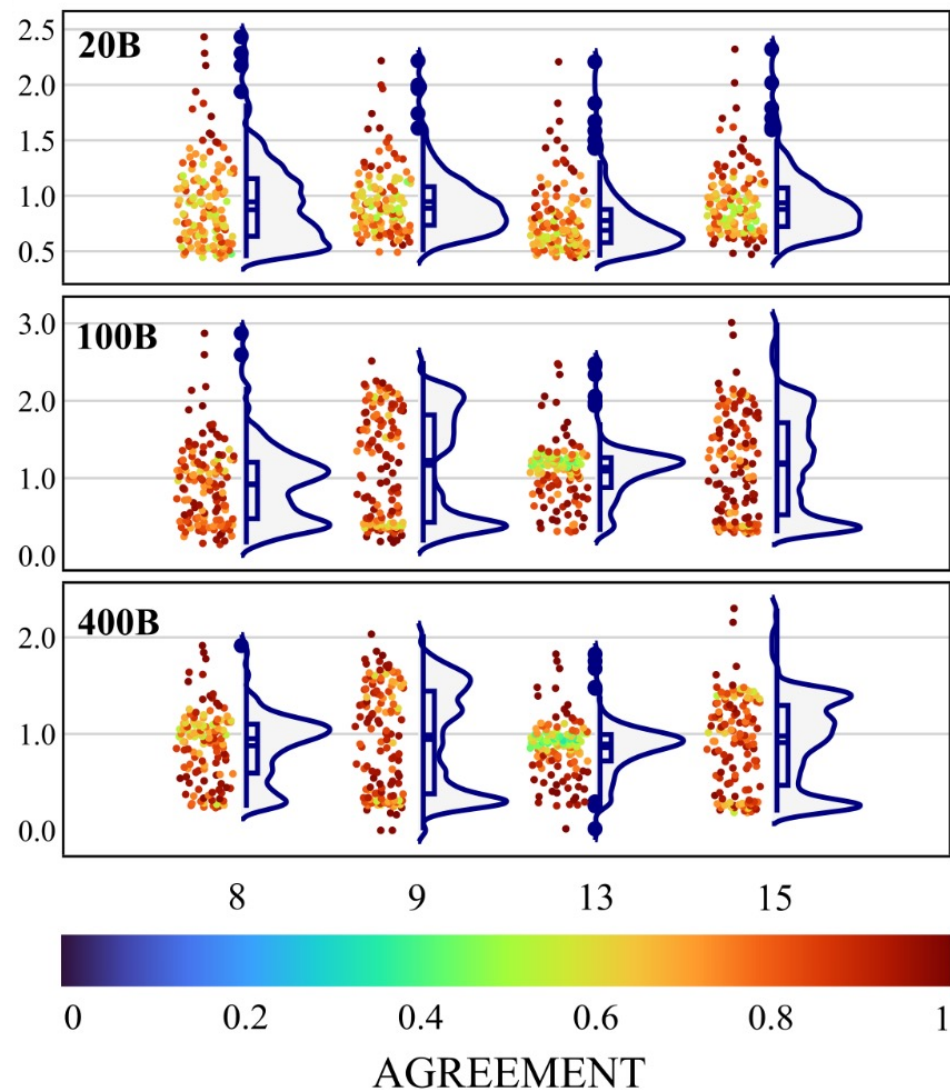
Further research directions

1 Theoretical framework for complex spectral shapes

2 Downstream effects of head-wise rewind

3 Spectral fingerprints of SFT and RL fine-tuning

4 Combination with analysis of activations



Thank you for your attention!

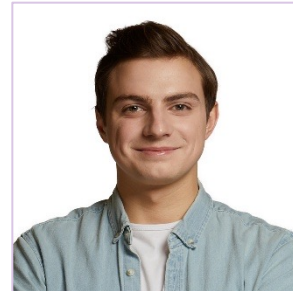
DIFFRACT · SPECTRAL VIEW OF LLM DOMAIN ADAPTATION



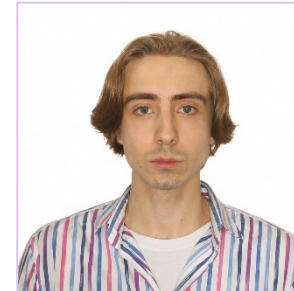
Nikita Borodin



Maria Krylova



Artem Zabolotnyi



Dmitry Aspisov



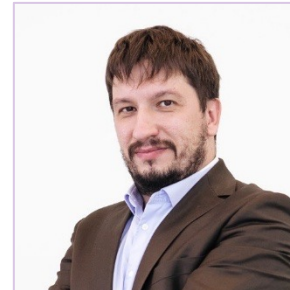
Egor Shikov



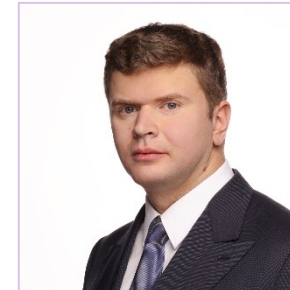
Nikita Tyuplyaev



Oleg Travkin



Roman Alferov



Dmitry Vinichenko



Code & analysis toolkit

<https://github.com/Risk-AI-Research/diffract-training>

<https://github.com/Risk-AI-Research/diffract>

Corresponding author: Dmitry Vinichenko — dmitry.vinichenko@risk-ai-research.tech

