

Establishing a Scale for Kullback-Leibler Divergence in Language Models Across Various Settings

Ryo Kishino¹, Yusuke Takase¹, Momose Oyama^{1,2}, Hiroaki Yamagiwa³, Hidetoshi Shimodaira^{1,2}

¹Kyoto University ²RIKEN ³SB Intuitions



Summary

- Representing each language model by its log-likelihood vector embeds models into a common space, where distances reflect KL divergence.
- This work establishes an interpretable absolute scale for KL divergence across various settings.
- We observed that training trajectories are strongly subdiffusive in log-likelihood space.

Log-Likelihood Vector [1]

- A log-likelihood vector represents an LM by its log probabilities over a fixed text set.

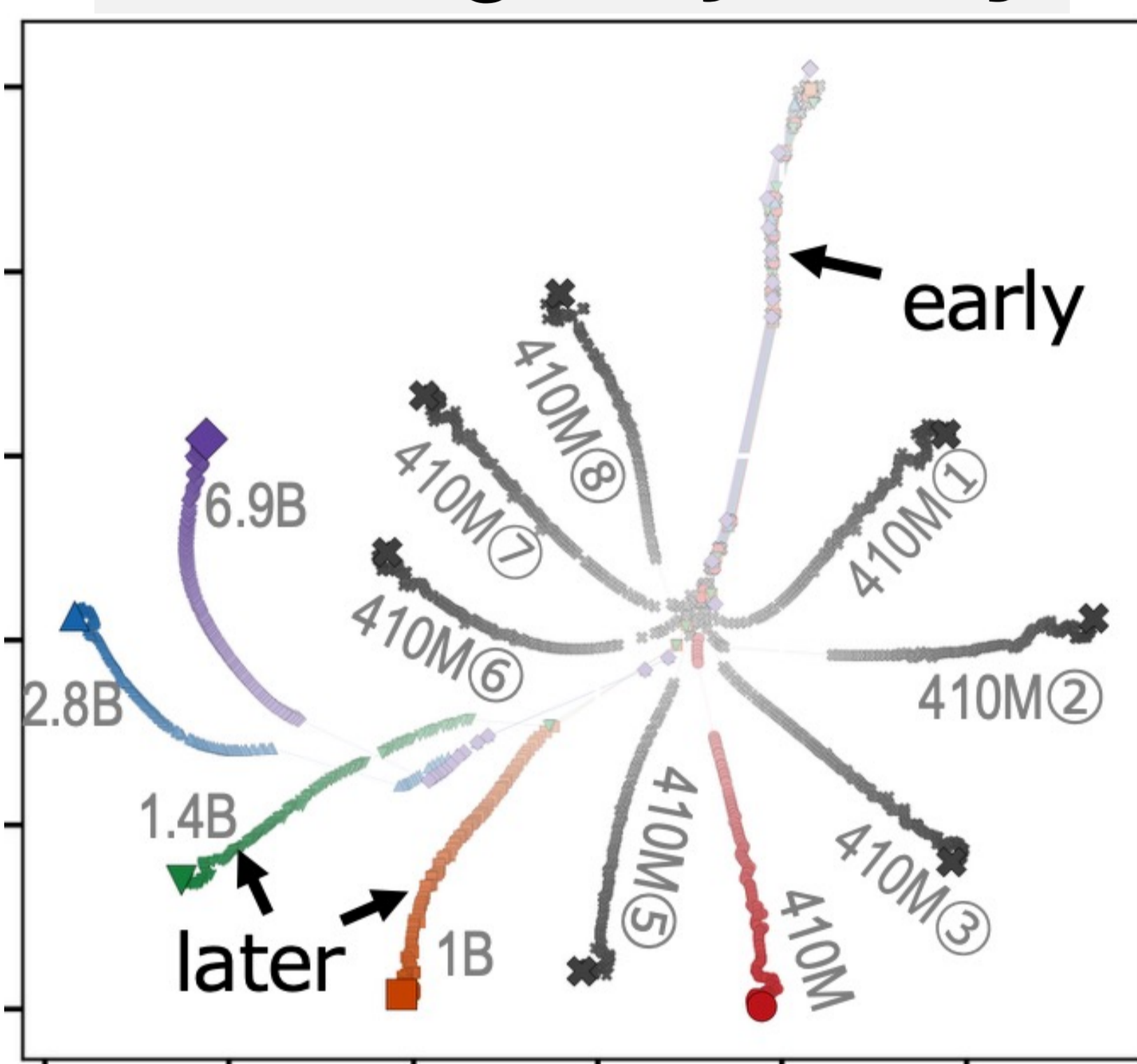
$$(\log p_i(x_1), \dots, \log p_i(x_N)) \in \mathbb{R}^N$$

KL Divergence

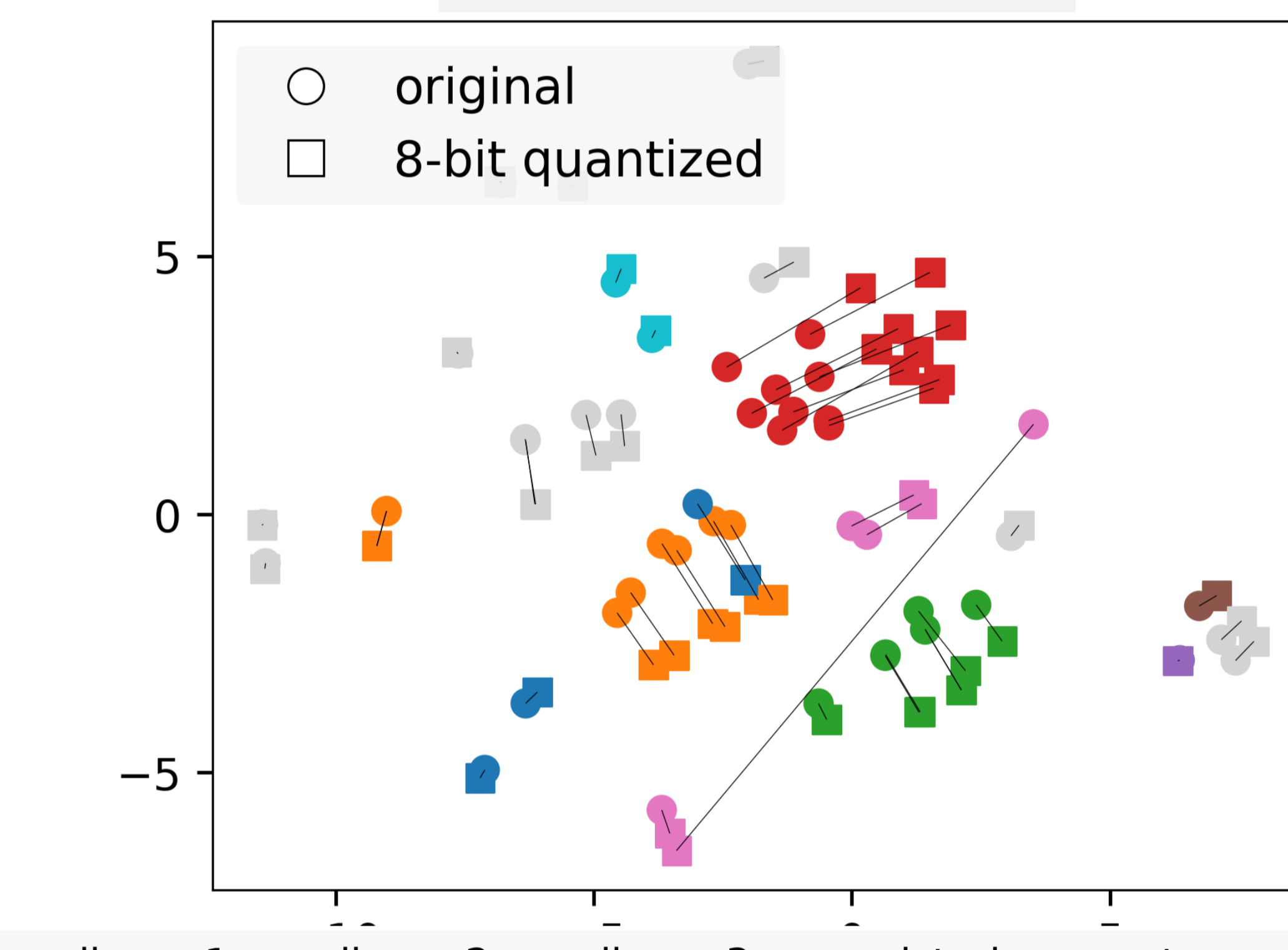
$$\text{KL}(p_i, p_j) \approx \frac{1}{2N} \|\mathbf{q}_i - \mathbf{q}_j\|^2$$

Model Map Across Various Settings

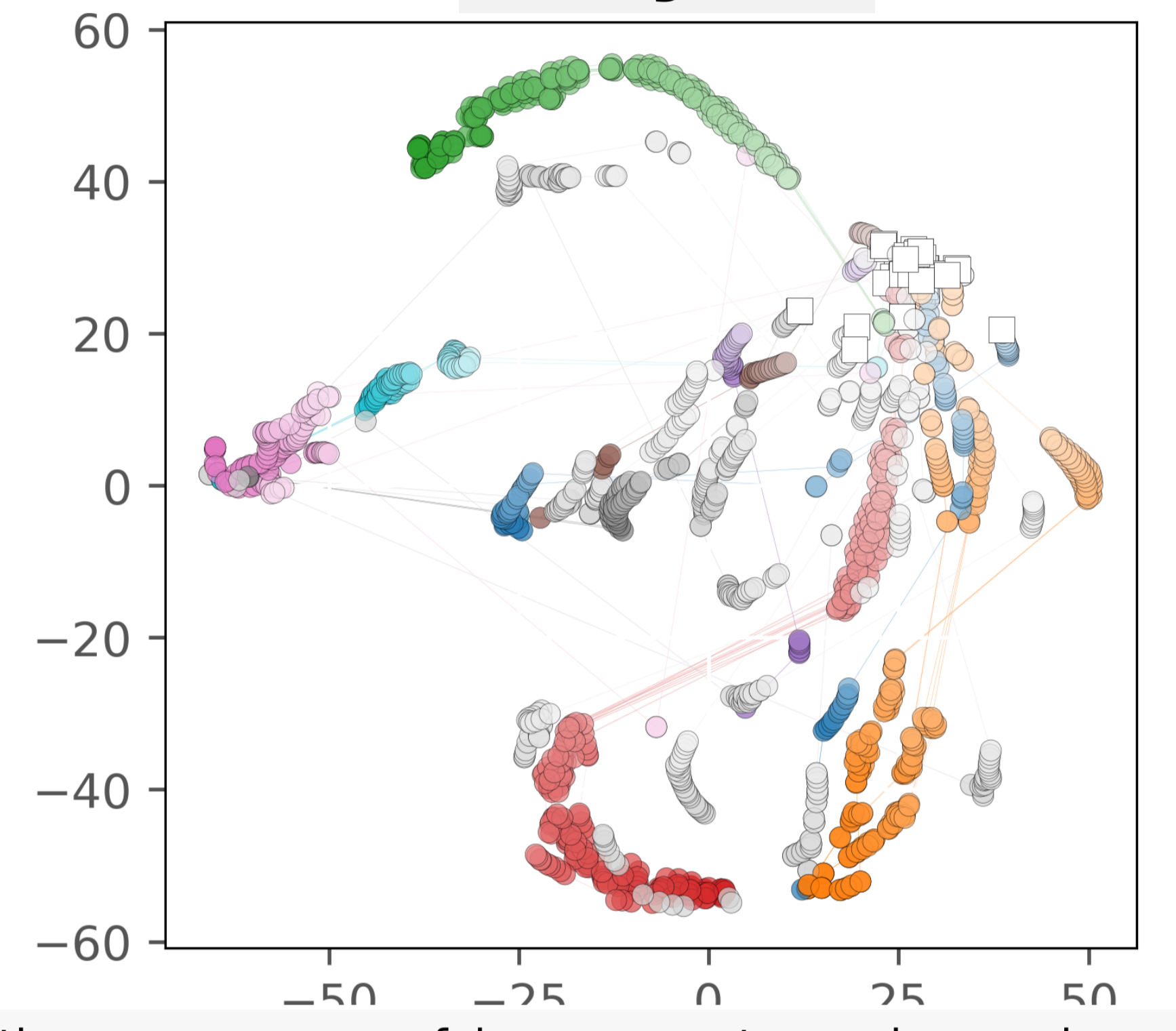
Training Trajectory



Quantization

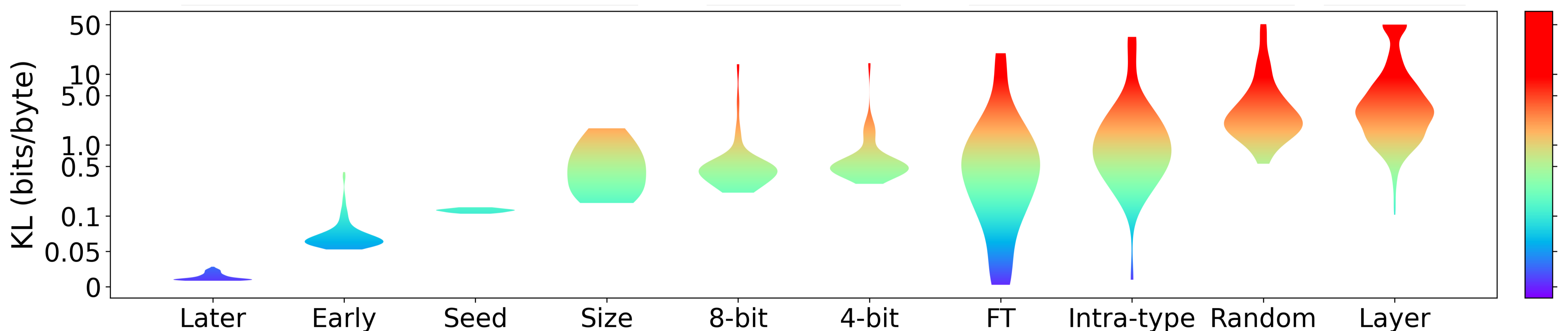


Layer



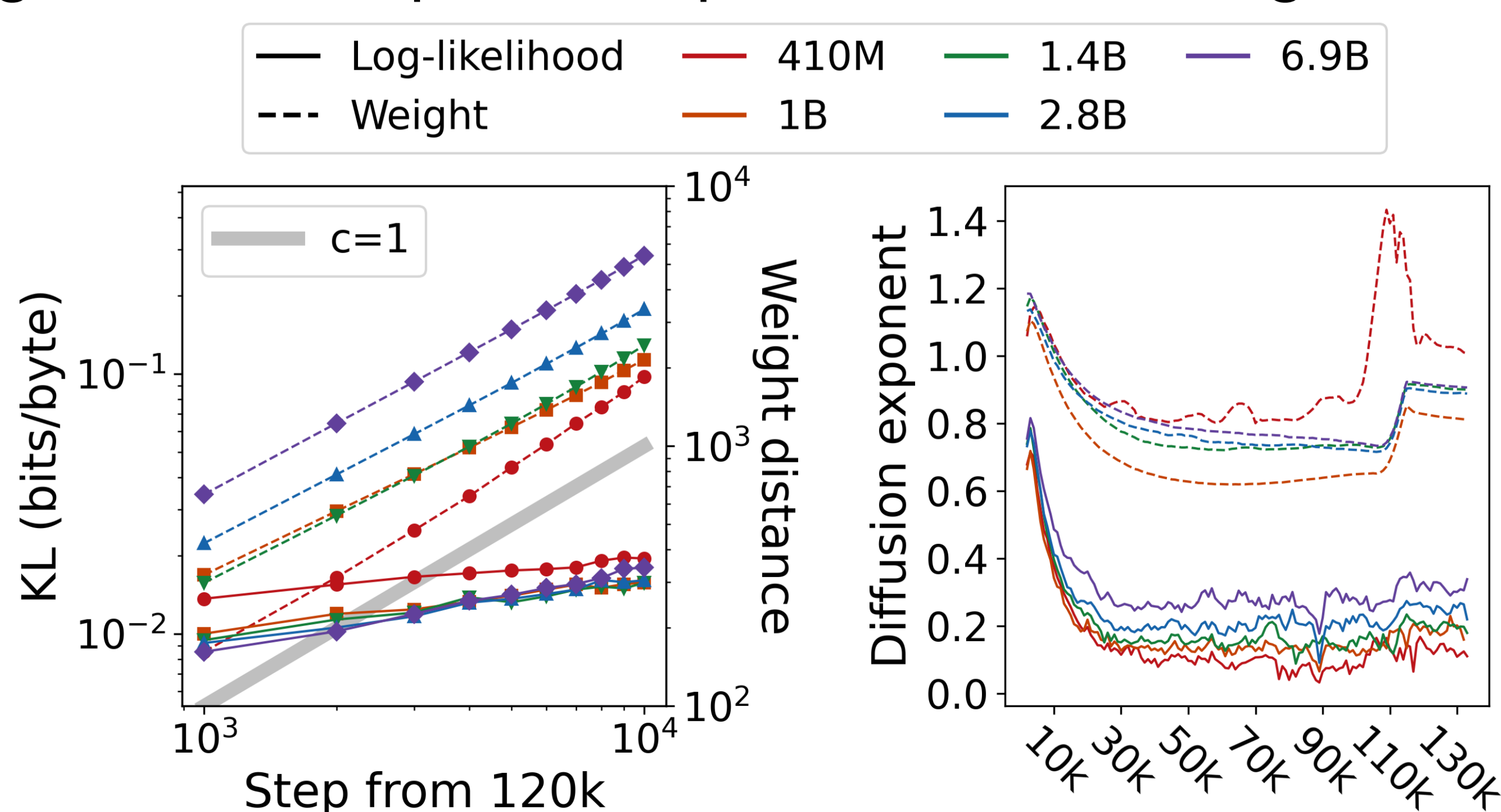
KL Divergence Scale

- KL divergence provides an interpretable scale for comparing model changes across checkpoints, seeds, sizes, quantization, fine-tuning, and layers.



Training Dynamics Analysis

- We analyze pretraining trajectories via KL divergence and find strongly subdiffusive behavior in log-likelihood space despite continued weight drift.



Weight Space vs. Log-Likelihood Space

- Anomalous diffusion:

$$\|\mathbf{W}_t - \mathbf{W}_{t_0}\|^2 \propto |t - t_0|^{c_w}, \quad c_w \approx 1$$

$$\|\mathbf{q}_t - \mathbf{q}_{t_0}\|^2 \propto |t - t_0|^{c_q}, \quad c_q \approx 0.2$$

- Hölder exponent of $f : \mathbf{W} \mapsto \mathbf{q}(\mathbf{W})$ can be calculated by

$$\alpha = c_q / c_w \approx 0.2$$