

Dispersion loss counteracts embedding condensation and improves generalization in small language models

Chen Liu^{*1}, Xingzhi Sun^{*1}, Xi Xiao^{*2,3}, Alexandre Van Tassel^{*1}, Ke Xu¹, Kristof Reimann¹, Danqi Liao¹, Mark Gerstein¹, Tianyang Wang², Xiao Wang³, Smita Krishnaswamy¹
^{*} Equal contribution ¹ Yale University ² University of Alabama at Birmingham ³ Oak Ridge National Laboratory. Correspondence to: Smita Krishnaswamy <smita.krishnaswamy@yale.edu>.

One-liner summary: What makes LLMs better than small LMs? Data? Parameters? Geometry might play a role!

Motivation

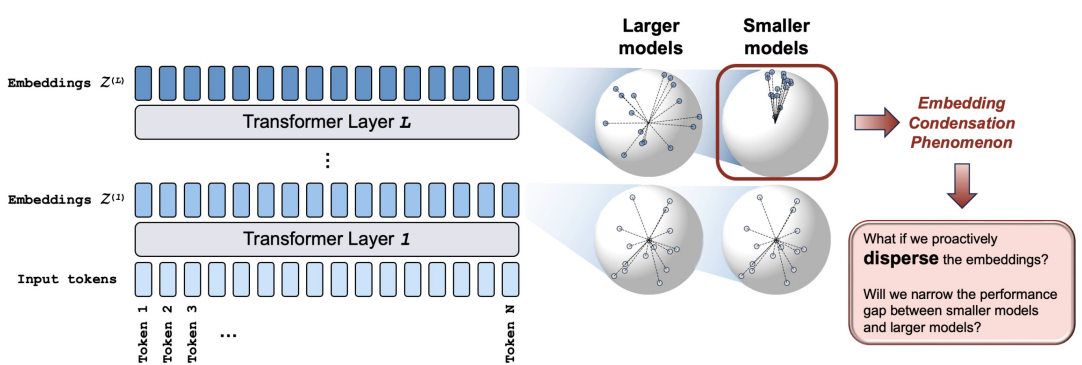


Figure 1. Illustration of the embedding condensation phenomenon. In pre-trained language models, embeddings of all tokens from the same input sequence condense into a narrow cone after being processed by many Transformer layers. This phenomenon is substantially more pronounced in smaller models than in larger models within the same family, which motivates our hypothesis in Section 3.3.

Features of embedding condensation

Feature 1. More severe in smaller models than in larger counterparts (Figure 2).

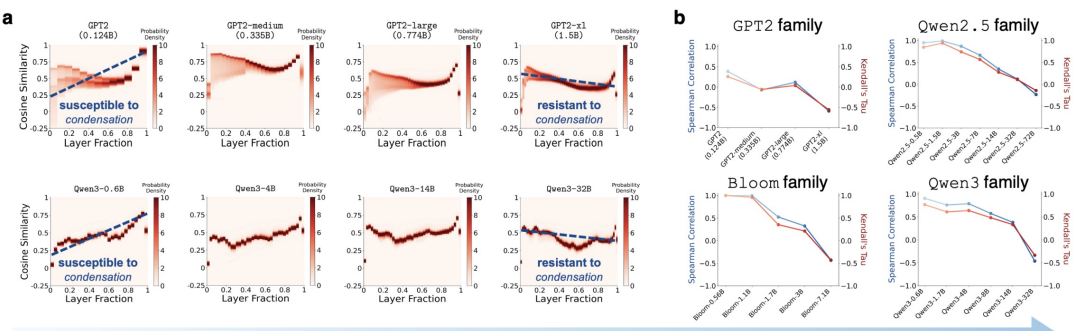


Figure 2. Qualitative and quantitative observations of the embedding condensation phenomenon. **a.** The cosine similarity heatmaps demonstrate that smaller models (e.g., GPT2, Qwen3-0.6B) are susceptible to condensation, since token cosine similarities become increasingly positive as the embeddings proceed to deeper layers. In contrast, larger models (e.g., GPT2-x1, Qwen3-32B) are more resistant to embedding condensation. **b.** Quantifications using Spearman correlation and Kendall's Tau demonstrate a consistent trend of "larger model, less condensation" across multiple families of language models. Additional results can be found in Figure S1.

Feature 2. Reproducible under confounder-controlled settings (Figure 3).

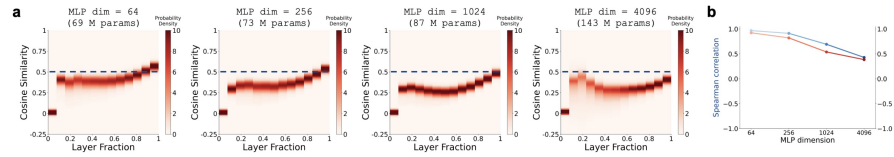


Figure 3. In a highly controlled experiment, we reproduced the observation of "larger model, less condensation". We pre-trained four GP2-like models of varying sizes that differ only in MLP dimension, while keeping all other factors fixed, including the number of layers, embedding dimension, dataset, and training configuration. The resulting models exhibit consistent trends in embedding condensation, shown qualitatively (panel a) and quantitatively (panel b). Horizontal dashed lines are added to panel a for easier visual comparison.

Feature 3. Emerging at model initialization (Figure 4).

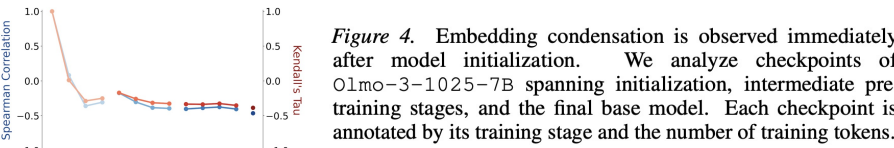


Figure 4. Embedding condensation is observed immediately after model initialization. We analyze checkpoints of O1mo-3-1025-7B spanning initialization, intermediate pre-training stages, and the final base model. Each checkpoint is annotated by its training stage and the number of training tokens.

Feature 4. Not resolved by distillation from a larger model (Figure 5).

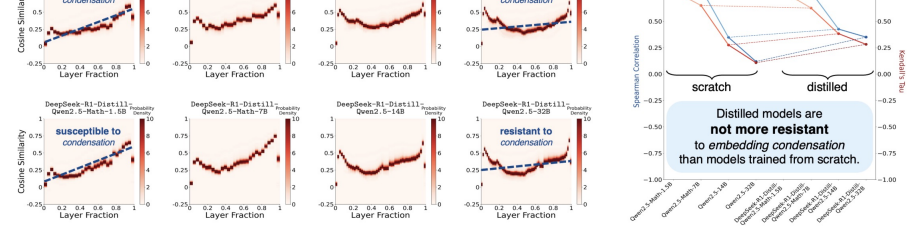


Figure 5. Knowledge distillation is not a remedy to embedding condensation, shown qualitatively (panel a) and quantitatively (panel b).

Solution: dispersion loss

Table 1. Our dispersion loss and its alternative formulations. Main implementation differences from (Wang & He, 2025) are highlighted in teal and magenta. Including or excluding diagonal terms yields identical gradients and is therefore cosmetic. For dispersion loss and ℓ_2 -repeL, we adopt the \log -sum-exp trick for numerical stability, which differs from $\log(\text{mean}(\exp(\cdot)))$ only by an additive constant. For ℓ_2 -repeL, we include a norm regularization term to prevent unbounded expansion of embeddings. For Orthogonalization, the distance margin is fixed to $\frac{1}{2}$ since we use angular distance, where $\frac{1}{2}$ corresponds to orthogonality and thus serves as the ideal margin.

| | For generative modeling in diffusion-based models (Wang & He, 2025) | | For improving generalization of Transformer-based language models (Ours) | |
|---------------------------------|---|--|---|--|
| | formulation | term definition | formulation | term definition |
| Dispersion loss | $\log E_{\mathbf{z}_i}[\exp(-D(\mathbf{z}_i, \mathbf{z}_j)/\tau)]$ | $D(\mathbf{z}_i, \mathbf{z}_j) = -\text{cossim}(\mathbf{z}_i, \mathbf{z}_j)$ | $\log \sum_{i \neq j} \exp(-D(\mathbf{z}_i, \mathbf{z}_j)/\tau)$ | $D(\mathbf{z}_i, \mathbf{z}_j) = \frac{\text{arccos}(\text{cossim}(\mathbf{z}_i, \mathbf{z}_j))}{2}$ |
| Alternative formulations | | | | |
| Decorrelation | $\sum_{m, n} \text{Cov}_{mn}^2$ | $\text{Cov}^2 = \frac{\mathbf{Z}^T \mathbf{Z}}{d} - \frac{\mathbf{Z} - \mu \mathbf{1}}{\sigma} \frac{\mathbf{Z} - \mu \mathbf{1}}{\sigma}$ | $\sum_{m, n} \text{Cov}_{mn}^2$ | $\text{Cov}^2 = \frac{\mathbf{Z}^T \mathbf{Z}}{d} - \frac{\mathbf{Z} - \mu \mathbf{1}}{\sigma} \frac{\mathbf{Z} - \mu \mathbf{1}}{\sigma}$ |
| ℓ_2 -repeL | $\log E_{\mathbf{z}_i}[\exp(-D(\mathbf{z}_i, \mathbf{z}_j)/\tau)]$ | $D(\mathbf{z}_i, \mathbf{z}_j) = \ \mathbf{z}_i - \mathbf{z}_j\ _2^2$ | $\log \sum_{i \neq j} \exp(-D(\mathbf{z}_i, \mathbf{z}_j)/\tau) + \lambda_{\text{norm}} \ \mathbf{Z}\ _F^2$ | $D(\mathbf{z}_i, \mathbf{z}_j) = \ \mathbf{z}_i - \mathbf{z}_j\ _2^2$ |
| Orthogonalization | $E_{\mathbf{z}_i}[\max(0, \tau - D(\mathbf{z}_i, \mathbf{z}_j))]^2$ | $D(\mathbf{z}_i, \mathbf{z}_j) = -\text{cossim}(\mathbf{z}_i, \mathbf{z}_j)$ | $E_{\mathbf{z}_i}[\max(0, \frac{1}{2} - D(\mathbf{z}_i, \mathbf{z}_j))]^2$ | $D(\mathbf{z}_i, \mathbf{z}_j) = \frac{\text{arccos}(\text{cossim}(\mathbf{z}_i, \mathbf{z}_j))}{2}$ |

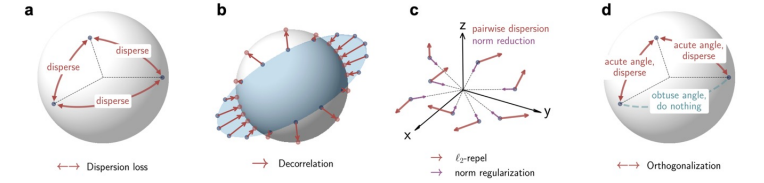


Figure 6. Illustration of how dispersion loss and its alternative formulations promote embedding dispersion. **a.** Dispersion loss enforces uniform angular dispersion by spreading out all pairs along the unit hypersphere. **b.** Decorrelation loss encourages different feature dimensions to remain uncorrelated. ℓ_2 -repeL loss increases pairwise Euclidean distance, while the norm regularization prevents unbounded expansion. **d.** Orthogonalization loss spreads out vectors forming acute angles while leaving obtuse ones unchanged.

Results: mid-training

Table 2. Using dispersion loss during mid-training improves performance on language tasks. For each base model, the best single-benchmark metrics are displayed in bold, whereas the best average ranks and best average performances are boxed and in bold. We also perform the Student's t -tests on the average performances and report the significance level with respect to " $\mathcal{L}_{\text{train}} + \text{Dispersion loss}$ ".

| Model | Mid-training | | Zero-shot | | | | | | | | | | Few-shot | | Average \uparrow | t-test | |
|------------|--------------|--------------|--------------|-----------------|--------------------|-----------------------|-----------------|-----------------------|-----------------------|----------------|----------------|-------------------|-----------------|-------------------|--------------------|-----------------|---------|
| | Train | Loss | A190hours | ANLI \uparrow | LAMBADA \uparrow | OpenbookQA \uparrow | PIQA \uparrow | TruthfulQA \uparrow | WinoGrande \uparrow | ARC \uparrow | ARC \uparrow | MeMCQA \uparrow | MMLU \uparrow | Rank \downarrow | | | |
| GPT2 | x | — | 34.4 | 30.8 | 15.6 | 61.6 | 40.3 | 62.2 \pm 0.7 | 43.6 \pm 0.2 | 52.8 \pm 3.1 | 41.6 \pm 0.0 | 17.4 \pm 0.6 | 24.2 \pm 3.9 | 24.8 \pm 0.4 | 32.6 | 34.95 \pm 0.1 | p<0.001 |
| | \checkmark | \checkmark | 1.122 (100x) | 340 \pm 0.3 | 32.4 \pm 0.6 | 16.4 \pm 0.1 | 62.2 \pm 0.7 | 43.6 \pm 0.2 | 52.8 \pm 3.1 | 41.6 \pm 0.0 | 17.4 \pm 0.6 | 24.2 \pm 3.9 | 24.8 \pm 0.4 | 32.6 | 34.95 \pm 0.1 | p<0.001 | |
| | \checkmark | \checkmark | 1.122 (100x) | 343 \pm 0.1 | 33.0 \pm 0.3 | 16.9 \pm 0.4 | 61.1 \pm 0.7 | 43.9 \pm 0.2 | 53.6 \pm 2.0 | 43.7 \pm 0.4 | 18.0 \pm 0.7 | 21.6 \pm 2.5 | 25.1 \pm 0.1 | 4.3 | 35.15 \pm 0.06 | p<0.001 | |
| | \checkmark | \checkmark | 1.127 (100x) | 350 \pm 0.3 | 33.8 \pm 0.5 | 16.6 \pm 0.3 | 61.0 \pm 0.0 | 44.0 \pm 0.1 | 53.2 \pm 0.8 | 44.3 \pm 0.4 | 18.6 \pm 1.1 | 22.4 \pm 1.7 | 25.7 \pm 0.2 | 3.2 | 35.36 \pm 0.15 | n.s. | |
| | \checkmark | \checkmark | 1.221 (100x) | 358 \pm 0.0 | 33.6 \pm 0.5 | 17.0 \pm 0.6 | 60.8 \pm 0.6 | 43.9 \pm 0.3 | 52.0 \pm 1.7 | 43.8 \pm 0.5 | 18.0 \pm 1.3 | 24.4 \pm 2.3 | 25.6 \pm 0.2 | 3.6 | 35.41 \pm 0.06 | n.s. | |
| | \checkmark | \checkmark | 1.175 (100x) | 348 \pm 0.3 | 32.8 \pm 0.2 | 16.5 \pm 0.2 | 60.8 \pm 0.3 | 43.8 \pm 0.3 | 54.2 \pm 0.8 | 43.4 \pm 0.2 | 17.6 \pm 0.8 | 21.6 \pm 2.0 | 25.7 \pm 0.1 | 4.1 | 35.42 \pm 0.11 | n.s. | |
| | \checkmark | \checkmark | 1.180 (100x) | 344 \pm 0.1 | 31.2 \pm 1.1 | 16.8 \pm 0.7 | 61.6 \pm 0.0 | 44.2 \pm 0.4 | 53.4 \pm 0.8 | 44.2 \pm 0.1 | 18.0 \pm 0.3 | 24.8 \pm 3.3 | 25.6 \pm 0.2 | 3.2 | 35.42 \pm 0.19 | n.s. | |
| | \checkmark | \checkmark | 1.176 (100x) | 348 \pm 0.0 | 34.0 \pm 0.4 | 16.6 \pm 0.0 | 62.4 \pm 0.3 | 44.2 \pm 0.1 | 51.6 \pm 0.1 | 42.6 \pm 0.4 | 18.0 \pm 0.8 | 26.6 \pm 2.3 | 25.4 \pm 0.2 | 3.2 | 35.61 \pm 0.12 | ref. | |
| GPT2-m | x | — | 35.3 | 40.4 | 18.6 | 66.3 | 40.1 | 64.8 | 50.2 | 19.9 | 29.1 | 25.3 | 5.4 | 37.70 | p<0.001 | | |
| | \checkmark | \checkmark | 2.541 (100x) | 331 \pm 0.6 | 42.2 \pm 0.3 | 19.1 \pm 0.0 | 67.7 \pm 0.4 | 48.5 \pm 1.3 | 55.5 \pm 1.0 | 54.2 \pm 2.7 | 18.9 \pm 0.7 | 28.2 \pm 2.3 | 25.1 \pm 0.1 | 4.5 | 38.5 \pm 0.18 | p<0.001 | |
| | \checkmark | \checkmark | 2.580 (100x) | 333 \pm 0.1 | 44.0 \pm 0.3 | 18.6 \pm 0.6 | 66.7 \pm 0.1 | 44.5 \pm 0.6 | 52.1 \pm 1.6 | 51.7 \pm 1.6 | 19.5 \pm 1.8 | 25.5 \pm 0.1 | 25.8 \pm 0.3 | 4.8 | 38.16 \pm 0.25 | p<0.001 | |
| | \checkmark | \checkmark | 2.780 (100x) | 334 \pm 0.0 | 40.4 \pm 1.3 | 18.9 \pm 0.8 | 66.6 \pm 0.1 | 42.2 \pm 0.1 | 50.8 \pm 1.4 | 51.8 \pm 2.5 | 19.8 \pm 1.4 | 26.5 \pm 0.1 | 26.1 \pm 0.0 | 4.4 | 38.15 \pm 0.29 | p<0.001 | |
| | \checkmark | \checkmark | 2.675 (100x) | 338 \pm 0.8 | 44.2 \pm 0.8 | 18.5 \pm 0.6 | 66.2 \pm 0.1 | 42.4 \pm 0.1 | 54.4 \pm 1.1 | 54.0 \pm 2.0 | 18.6 \pm 0.6 | 28.9 \pm 2.1 | 25.3 \pm 0.4 | 4.7 | 38.61 \pm 0.03 | p<0.001 | |
| | \checkmark | \checkmark | 2.662 (100x) | 330 \pm 0.4 | 45.2 \pm 0.4 | 18.6 \pm 0.6 | 66.4 \pm 0.4 | 41.4 \pm 0.2 | 55.1 \pm 0.6 | 53.6 \pm 2.3 | 18.8 \pm 0.4 | 29.0 \pm 1.0 | 25.0 \pm 0.1 | 4.2 | 38.76 \pm 0.20 | p<0.001 | |
| | \checkmark | \checkmark | 2.673 (100x) | 338 \pm 0.1 | 45.2 \pm 0.8 | 19.2 \pm 0.2 | 67.5 \pm 0.1 | 43.4 \pm 0.1 | 56.4 \pm 1.1 | 53.8 \pm 1.1 | 28.1 \pm 0.4 | 28.6 \pm 0.1 | 25.7 \pm 0.1 | 4.2 | 39.25 \pm 0.15 | ref. | |
| GPT2-x1 | x | — | 33.6 | 47.8 | 19.6 | 71.8 | 38.9 | 58.4 | 54.0 | 22.4 | 26.2 | 25.6 | — | 39.81 | — | | |
| | \checkmark | \checkmark | 30.4 | 49.4 | 22.2 | 71.8 | 38.0 | 57.2 | 58.0 | 25.6 | 27.0 | 25.2 | — | 40.89 | — | | |
| Qwen3-0.6B | x | — | 35.0 | 43.0 | 19.5 | 66.5 | 40.7 | 60.5 | 61.5 | 32.5 | 26.5 | 49.4 | 5.7 | 44.11 | p<0.001 | | |
| | \checkmark | \checkmark | 4.676 (100x) | 325 \pm 0.3 | 52.0 \pm 0.3 | 21.5 \pm 0.5 | 67.5 \pm 1.4 | 44.3 \pm 0.9 | 61.8 \pm 0.8 | 68.0 \pm 2.8 | 33.0 \pm 1.3 | 29.5 \pm 0.9 | 30.0 \pm 0.1 | 4.1 | 45.95 \pm 0.37 | p<0.001 | |
| | \checkmark | \checkmark | 4.711 (100x) | 325 \pm 0.0 | 48.7 \pm 0.7 | 21.5 \pm 0.7 | 67.2 \pm 1.1 | 44.5 \pm 0.7 | 62.7 \pm 1.1 | 68.5 \pm 0.7 | 33.0 \pm 1.3 | 29.5 \pm 0.9 | 30.0 \pm 0.1 | 4.1 | 46.17 \pm 0.28 | p<0.001 | |
| | \checkmark | \checkmark | 4.829 (100x) | 350 \pm 0.7 | 49.0 \pm 4.2 | 20.0 \pm 2.1 | 67.2 \pm 0.4 | 48.8 \pm 0.9 | 58.8 \pm 1.8 | 69.0 \pm 1.4 | 34.2 \pm 1.1 | 36.8 \pm 0.4 | 49.3 \pm 0.0 | 4.0 | 46.79 \pm 0.27 | p<0.05 | |
| | \checkmark | \checkmark | 4.959 (100x) | 350 \pm 0.0 | 50.5 \pm 0.4 | 19.5 \pm 1.8 | 67.5 \pm 1.1 | 47.0 \pm 2.3 | 59.5 \pm 1.8 | 68.5 \pm 3.2 | 30.0 \pm 1.8 | 34.0 \pm 0.4 | 49.8 \pm 0.2 | 3.5 | 46.82 \pm 0.25 | p<0.001 | |
| | \checkmark | \checkmark | 4.868 (100x) | 335 \pm 0.0 | 46.0 \pm 2.1 | 19.0 \pm 0.7 | 66.0 \pm 0.4 | 47.4 \pm 0.2 | 59.1 \pm 1.4 | 69.0 \pm 3.2 | 35.0 \pm 1.8 | 40.0 \pm 0.4 | 46.9 \pm 0.1 | 4.8 | 46.25 \pm 0.04 | p<0.001 | |
| | \checkmark | \checkmark | 4.936 (100x) | 368 \pm 1.1 | 50.0 \pm 0.7 | 22.0 \pm 0.7 | 65.5 \pm 0.4 | 47.1 \pm 0.1 | 56.5 \pm 2.1 | 71.5 \pm 3.2 | 33.0 \pm 1.1 | 36.0 \pm 0.7 | 48.9 \pm 0.2 | 4.2 | 46.89 \pm 0.20 | p<0.001 | |
| | \checkmark | \checkmark | 4.878 (100x) | 358 \pm 0.0 | 48.5 \pm 0.8 | 22.5 \pm 0.6 | 65.0 \pm 1.6 | 49.8 \pm 1.1 | 58.1 \pm 1.6 | 72.5 \pm 3.3 | 34.5 \pm 0.8 | 37.5 \pm 1.3 | 49.2 \pm 0.2 | 3.2 | 47.47 \pm 0.16 | ref. | |
| Qwen3-1.7B | x | — | 39.5 | 53.0 | 29.0 | 71.5 | 47.6 | 60.5 | 62.5 | 32.0 | 50.0 | 45.5 | 6.11 | 51.1 | p<0.001 | | |
| | \checkmark | \checkmark | 9.148 (100x) | 400 \pm 0.9 | 60.7 \pm 0.3 | 28.5 \pm 0.9 | 74.3 \pm 0.6 | 49.1 \pm 0.2 | 61.3 \pm 0.4 | 71.3 \pm 1.4 | 47.7 \pm 1.2 | 49.7 \pm 1.0 | 63.1 \pm 0.0 | 4.3 | 54.57 \pm 0.17 | p<0.001 | |
| | \checkmark | \checkmark | 9.345 (100x) | 380 \pm 0.0 | 60.5 \pm 2.4 | 28.0 \pm 0.0 | 70.5 \pm 0.7 | 50.2 \pm 0.9 | 60.8 \pm 3.2 | 73.2 \pm 0.4 | 48.5 \pm 0.0 | 48.5 \pm 2.1 | 62.5 \pm 0.6 | 5.0 | 53.78 \pm 0.07 | p<0.001 | |
| | \checkmark | \checkmark | 9.288 (100x) | 382 \pm 1.0 | 61.0 \pm 1.4 | 28.0 \pm 1.4 | 72.5 \pm 0.0 | 49.9 \pm 0.2 | 69.5 \pm 0.4 | 74.2 \pm 0.4 | 47.2 \pm 1.1 | 49.0 \pm 1.6 | 62.4 \pm 0.1 | 4.9 | 53.98 \pm 0.07 | p<0.001 | |
| | \checkmark | \checkmark | 9.538 (100x) | 395 \pm 1.5 | 60.7 \pm 1.0 | 28.3 \pm 1.2 | 74.8 \pm 1.0 | 49.7 \pm 0.5 | 61.5 \pm 0.9 | 72.7 \pm 1.5 | 49.7 \pm 1.0 | 49.7 \pm 1.0 | 63.5 \pm 0.2 | 3.1 | 55.01 \pm 0.13 | p<0.001 | |
| | \checkmark | \checkmark | 9.460 (100x) | 385 \pm 0.0 | 47.8 \pm 0.6 | 23.0 \pm 0.0 | 70.8 \pm 1.2 | 42.1 \pm 0.2 | 57.5 \pm 0.9 | 68.3 \pm 1.9 | 41.3 \pm 1.8 | 38.5 \pm 2.6 | | | | | |