

---

# Scale Dependent Data Duplication

---

Anonymous Authors<sup>1</sup>

## Abstract

Data duplication during pretraining can degrade generalization and lead to memorization, motivating aggressive deduplication pipelines. However, at web scale, it is unclear what constitutes a “duplicate”: beyond surface-form matches, semantically equivalent documents (e.g. translations) may induce redundant training signals once models become sufficiently capable. Practically, this means that semantic duplicates operate increasingly like exact duplicates during training. We present evidence that duplication is scale-dependent in two ways. First, as model capability increases, cross-entropy loss gradients for semantically equivalent documents become more aligned. Smaller models, by contrast, produce gradients that reflect surface similarity (e.g., shared tokens) rather than semantic similarity. Second, we embedded all 192 million FineWeb-Edu-Dedup documents using EmbeddingGemma-300m. For moderate corpus sizes, the cosine similarity between nearest-neighbors follows an isotropic power law baseline. However, as corpus size grows to hundreds of billions of tokens, the nearest-neighbor similarities deviate sharply, indicating accelerated semantic collisions. Finally, controlled pretraining on data sampled with replacement from pools of finite unique documents shows that limited uniqueness yields mild degradation for small models, but rapidly increasing loss penalties for larger models, breaking naive scaling extrapolation. We derive explicit scaling laws that allow practitioners to estimate deviation from expected scaling due to limited semantic uniqueness of the pretraining corpus. Our results identify and resolve an unstudied source of scale-dependence, allowing for more accurate prediction at scale.

---

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

## 1. Introduction

Modern language models scale by increasing parameters, compute, and training tokens. For example, Llama 1 (Touvron et al., 2023) trained on  $\sim 1\text{T}$  tokens, while the Llama4 herd (Adcock et al., 2026) trained on up to  $40\text{T}$  tokens. At these scales, even small fractions of duplicated data can materially reduce the number of distinct training examples and harm downstream performance, emphasizing the importance of deduplication (Carlini et al., 2021; Hernandez et al., 2022; Lee et al., 2022; Comanici et al., 2025).

Deduplication is often framed as a dataset property: to deduplicate, simply remove exact duplicates and near-duplicates using simhashing techniques (Broder, 1997; Manku et al., 2007; Khan et al., 2025; Lee et al., 2022). Yet, what practically counts as a “duplicate” depends on the model as well: two documents that appear distinct may, to a sufficiently capable model, provide redundant training signal, and thus degrade training just as exact duplicates would. This work identifies a previously unknown source of scale dependence: as models become more capable, *semantic duplicates* induce the same gradients during training. In tandem, capable models are trained on larger corpora, in which the number of semantic collisions rapidly increases. Together, these effects create a recipe for model degradation.

### Contributions:

1. We quantify the **emergence of semantic sensitivity** during training by measuring cosine similarity between per-document cross-entropy gradients across a suite of models and semantic-preserving transformations. We find that in more capable models, semantic duplicates induce similar gradients during training.
2. We study **semantic collisions** by embedding 192M documents from FineWeb-Edu-Dedup (Penedo et al., 2024) documents and analyzing nearest-neighbor (NN) statistics across dataset scales from  $10^4$  to  $10^8$  documents. We discover that power laws governing scaling for moderate corpus sizes break down for large corpora. This collapse of scaling laws occurs earlier for synthetic corpora, revealing lower semantic diversity.
3. We examine the consequences for predictability by training scaling ladders on streams sampled with replacement from finite pools of  $K$  unique documents,

showing that limited uniqueness breaks naive scaling extrapolation. We derive more complete scaling laws that explicitly quantify the effects of limited uniqueness, restoring predictability. Furthermore, we show how to estimate an effective  $K$  directly from mean nearest-neighbor cosine similarity.

We defer related work to Appendix A.

## 2. Emergence of Semantics

As model capabilities increase, semantically equivalent documents induce similar training signals, as measured by the gradient of the per-document cross-entropy loss. Consequently, if two documents are semantic duplicates (e.g., translations), then a sufficiently capable model will update its parameters in similar directions when trained on both documents. Practically, this means that semantic duplicates operate increasingly like exact duplicates during training.

### 2.1. Experimental Setup

We sample  $N = 1000$  texts  $\{x_i\}_{i=1}^N$  from FineWeb-EduDedup (Penedo et al., 2024). To reduce variance due to length, each text is truncated to at most  $T = 2000$  tokens using the tokenizer of the model under evaluation.

We compute the per-document full-parameter gradient

$$g(x; \theta) = \nabla_{\theta} \ell(x_i; \theta), \quad (1)$$

where  $\ell$  is the mean next-token cross-entropy:

$$\ell(x; \theta) = \frac{1}{|x|} \sum_{u=1}^{|x|} \text{CE}(f_{\theta}(x)_u, x_{u+1}). \quad (2)$$

To establish a null baseline, we sample unrelated English documents  $(x_i, x_j)$ ,  $i \neq j$  and compute cosine similarity

$$\text{sim}(x_i, x_j) = \frac{\langle g(x_i; \theta), g(x_j; \theta) \rangle}{\|g(x_i; \theta)\|_2 \|g(x_j; \theta)\|_2}. \quad (3)$$

We repeat this across many random pairings to estimate the baseline mean  $\mu^-$  and standard deviation  $\sigma^-$ .

**Transformations.** We construct a set of transformations  $\mathcal{T} = \{\tau_1, \dots, \tau_L\}$  intended to preserve semantic content while perturbing surface form:

- Swap Characters: with probability 0.05, randomly replace each ascii character with another.
- Drop Words: Randomly delete each word with probability 0.05.
- Capitalize Humps: Capitalize every other character.
- Translate to Chinese/French/German.

For translations, we use Google’s Translate API (Google).

For each document  $x_i$  and transformation  $\tau$ , we compute

$$s_i^+(\tau) := \text{sim}_{\theta}(x_i, \tau(x_i)). \quad (4)$$

**Separability Metrics (Z-scores and AUC).** To summarize separation between positives and negatives, we define:

$$z(\tau) := \frac{\mu^+(\tau) - \mu^-}{\sigma^-}, \quad (5)$$

$$\mu^+(\tau) := \frac{1}{N} \sum_{i=1}^N s_i^+(\tau), \quad (6)$$

$$\mu^- := \frac{1}{|\mathcal{S}^-|} \sum_{(i,j) \in \mathcal{S}^-} \text{sim}_{\theta}(x_i, x_j), \quad (7)$$

$$(\sigma^-)^2 := \text{Var}_{(i,j) \in \mathcal{S}^-} [\text{sim}_{\theta}(x_i, x_j)]. \quad (8)$$

We also report AUC for distinguishing transformed gradients  $\{g_{\theta}(\tau(x_i))\}$  (positives) from unrelated gradients  $\{g(x_j; \theta)\}$  (negatives) using the score  $\text{sim}(x_i, \cdot)$ .

### 2.2. Results

Figure 1 reports mean gradient cosine similarities for both unrelated document pairs (negative baseline) and semantic-preserving transformations (positives). For smaller/weaker models, positive similarities for several transformations are comparable to or below the negative baseline, indicating that gradient direction is dominated by superficial features (e.g., language identity or capitalization). As model capability increases, transformed counterparts become consistently more aligned than unrelated pairs.

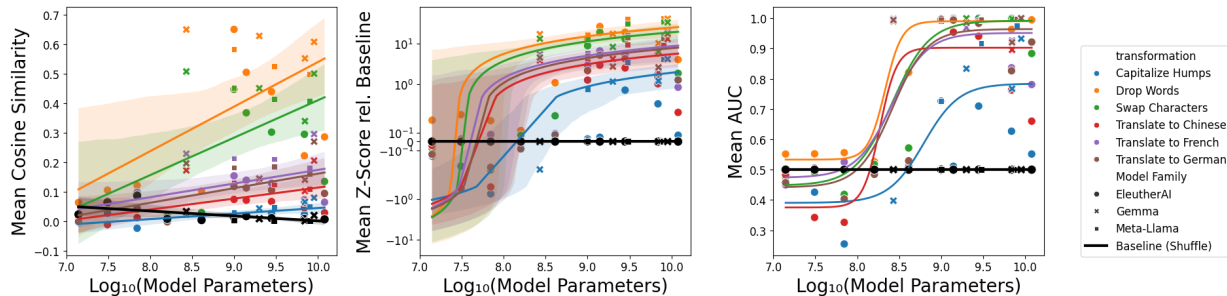
To quantify separability, we compute  $z(\tau)$  in Eq. (5) and AUC for the binary task described above. Figure 2 further shows that AUC increases with training progress for a fixed family and is achieved earlier by larger models.

**Interpretation.** Our findings suggest that semantic and exact duplicates have similar training impacts on capable models: if a model encodes meaning robustly, two semantically equivalent documents generate aligned weight updates. This provides a mechanism by which the *same dataset* can have a smaller effective size for more capable models.

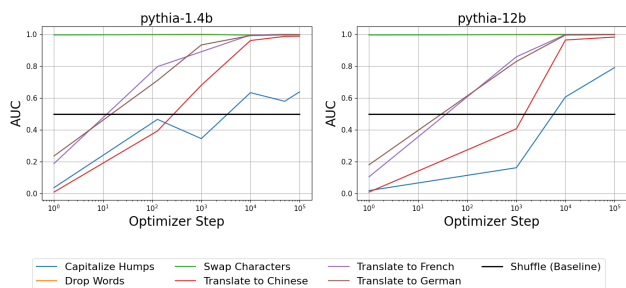
## 3. Semantic Collisions

When training models compute-optimally, corpus size grows in tandem with the number of parameters and model capabilities. In this section, we quantify the number of semantic collisions that occur in a deduplicated corpus of a given magnitude. We find that the rate of near-duplicates follows a predictable scaling law before increasing exponentially.

**Collision metrics.** For a set of unit-normalized embeddings



**Figure 1. Semantic-preserving transformations yield more aligned gradients for larger/stronger models.** We sample  $N=1000$  FineWeb-Edu-Dedup documents and compute per-document gradients of normalized next-token cross-entropy (Eq. 2) for each model. We report mean cosine similarity between (i) unrelated document pairs (negative baseline) and (ii) each document and its transformed counterpart (positives), including translations and light surface perturbations. Smaller/weaker models exhibit gradient similarity dominated by surface cues (language/casing), often failing to separate positives from negatives. As capability increases, positives become consistently more aligned than the negative baseline. Error bars show per-document standard deviation. Per-model-family results are in Figure 8.



**Figure 2. Semantic sensitivity emerges over training and is accelerated by scale.** For a fixed model family, we compute AUC to detect whether a candidate gradient corresponds to a semantic-preserving transformation of the same document versus an unrelated document, with cosine similarity to the original document gradient as the score. Early in training, AUC remains near 0.5 because gradients are dominated by surface-form features (language/casing). With additional optimizer steps, AUC increases, indicating that gradients increasingly reflect semantic content. Larger models reach a given AUC with fewer steps.

$\{v_i\}_{i=1}^N$ , define nearest-neighbor (NN) similarity

$$M_i := \max_{j \neq i} \langle v_i, v_j \rangle \quad \text{and cosine gap} \quad \Delta_i := 1 - M_i.$$

We report (i) estimates of  $\mathbb{E}[M_i]$  as a function of  $N$ , and (ii) tail probabilities  $\mathbb{P}(M_i \geq T)$  for fixed thresholds  $T$ .

### 3.1. Experimental Setup

We embed 190M texts from FineWeb-Edu-Dedup (Penedo et al., 2024) using EmbeddingGemma-300m (Vera et al., 2025). EmbeddingGemma-300m is a Matryoshka Representation Learning (Kusupati et al., 2024) model that produces embeddings of four nested sizes (768, 512, 256, and 128); sub-embeddings are obtained by slicing and re-normalizing.

We sample subsets of embeddings with cardinality ranging from  $10^4$  to  $10^8$  and estimate NN cosine similarities within each pool using FAISS (Douze et al., 2024).

### 3.2. Results

Figures 3 and 4 show that nearest-neighbor collision statistics initially match a power law, but deviate sharply at larger dataset sizes. Collisions occur more quickly in smaller embedding spaces, as expected. Beyond a scale threshold, the mean cosine gap decreases *faster than any fitted power law calibrated on smaller  $N$* . In log-linear coordinates, the decrease in NN cosine similarity is approximately linear over document corpus sizes less than 1M, before decaying much more quickly for corpora with over 10M documents. This presents a potentially compound threat to language models trained at scale: larger models that are more capable of identifying semantic duplicates are trained on more data, which contains more semantic duplicates than log-linear scaling laws would predict. Thus, models for which semantic duplicates are recognizable also experience far more of these duplicates, which could lead to loss of predictable scaling.

**A Note on Synthetic Data:** Recently, synthetic data has become a popular supplement for real data during pretraining and continued pretraining (Mishra et al., 2022; Chen et al., 2024a; Yang et al., 2024; Kang et al., 2025; Qin et al., 2025), though questions remain about whether it has sufficient diversity to provide a future alternative for real data. We repeat the experiment described in Section 3.1 for the fully-synthetic, 44M-document Recycling-the-Web pretraining corpus (Nguyen et al., 2025). We find that divergence from power law scaling (Figure 3) appears an order of magnitude earlier for synthetic pretraining data (Figure 5).

### 4. Impact on Training

We now probe practical implications. If semantic collisions reduce effective uniqueness, scaling-ladder extrapolation can fail. Because we cannot train models at the scale where *semantic* duplicates are recognized in our controlled setting, we model semantic collisions via *exact* document repeats

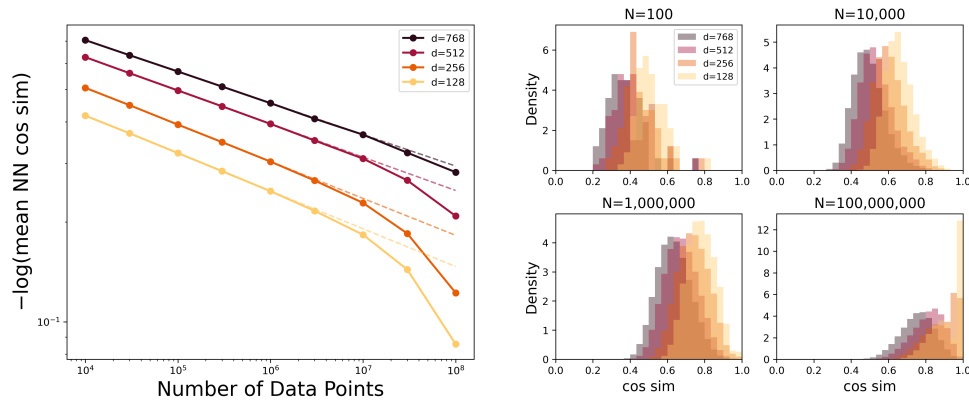


Figure 3. NN cosine similarity scaling deviates sharply at large corpus sizes. We embed 190M FineWeb-Edu-Dedup documents with EmbeddingGemma-300m and sample subsets of size ranging from  $10^4$ - $10^8$  without replacement. For each  $N$ , we estimate the mean nearest-neighbor cosine similarity using FAISS. Dashed lines show best-fit power laws over the small- $N$  regime where the uniform/vMF null predicts  $\mathbb{E}[\Delta_i] \propto N^{-2/d}$ . Beyond a scale threshold, the empirical curve steepens (smaller gaps than predicted), indicating substantially more near neighbors than expected under isotropic baselines.

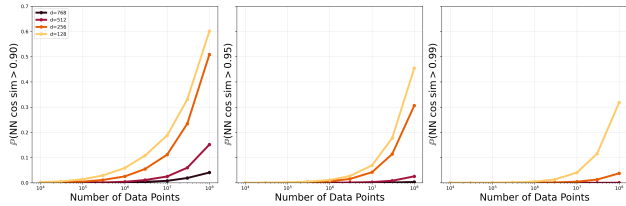


Figure 4. Tail collision rates accelerate with dataset size. For fixed thresholds  $T$ , we estimate the fraction of points with nearest-neighbor similarity  $M_i \geq T$ . These increase exponentially, as predicted under an isotropic baseline.

(sampling with replacement), which provides a pessimistic, worst-case proxy for repeated training signals.

#### 4.1. Experimental Setup

We sample pools of unique data of size  $K$  ranging from  $10^5$  through  $10^8$  unique documents sampled from FineWeb-Edu-Dedup. We construct training streams by sampling with replacement from each pool, inducing exact repeats. As a reference, we also train on streams constructed to minimize repeats (“approximately infinite unique data”) by sampling without replacement from Fineweb-Edu-Dedup.

We train scaling ladders of decoder-only, Chinchilla-optimal (Hoffmann et al., 2022) transformers based on the Qwen architecture ranging from 34M–344M parameters (Qwen et al., 2025; Yang et al., 2025). We match runs by compute (FLOPs) and report train and validation cross-entropy.

#### 4.2. Results and Discussion

Figure 6 shows that limiting  $K$  produces a scale-dependent degradation pattern. For smaller models, train and validation losses are consistent with standard scaling extrapolations, even when  $K$  is small; this can mislead scaling-ladder

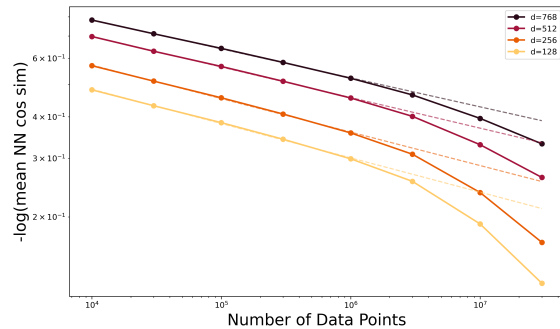


Figure 5. Nearest-neighbor cosine similarity scaling laws collapse an order of magnitude earlier for synthetic datasets: We embed the fully-synthetic pretraining dataset Recycling-the-Web (Nguyen et al., 2025) and find that the scaling law discovered in Figure 3 occurs an order of magnitude earlier for synthetic data, suggesting that the diversity of synthetic pretraining datasets should be improved.

planning. For larger FLOP budgets, finite- $K$  streams yield increasing loss penalties, breaking naive interpolation from smaller ladders trained under the same  $K$  constraint.

Although eval losses do not scale predictably with FLOP budgets under unique data constraints, *fractional loss increase* relative to the approximately-infinite baseline remains predictable:

$$\text{FracInc}(K) := \frac{L(K) - L(\infty)}{L(\infty)}. \quad (9)$$

This poses a challenge for those who predict scaling behavior, since we lack an infinite-unique-data baseline. Section 5.1 develops theory to resolve this problem and restore predictable scaling in the presence of semantic duplicates.

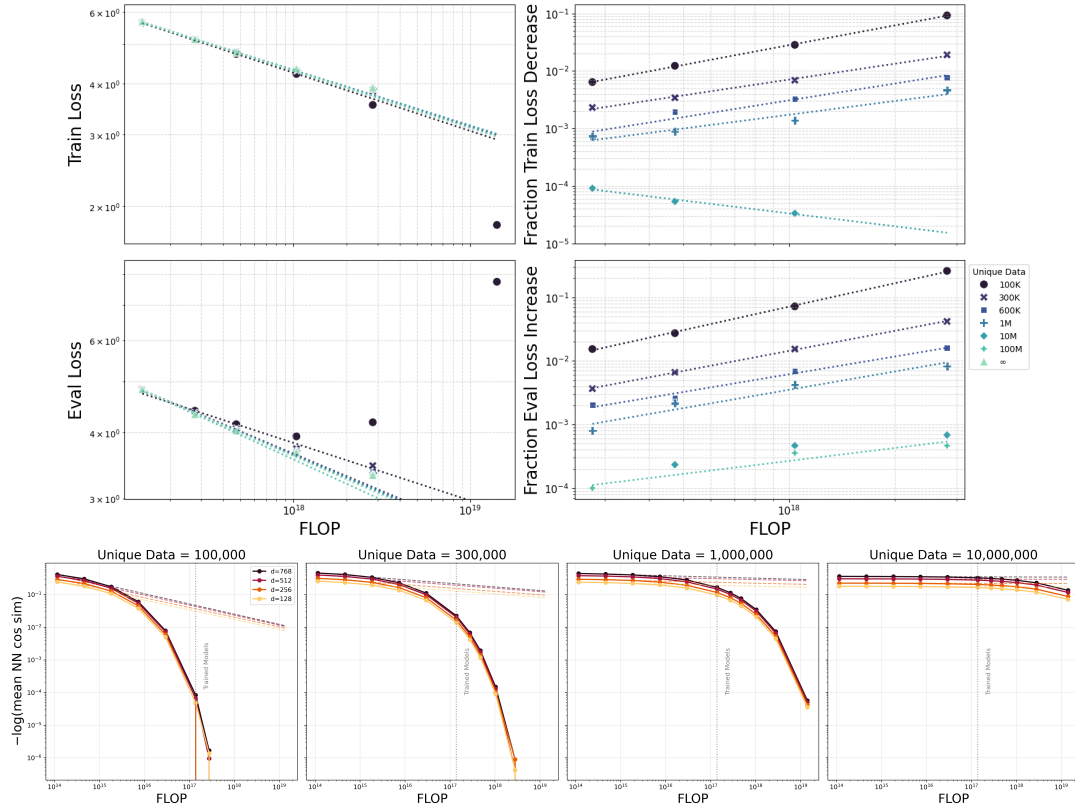


Figure 6. **Finite unique data pools induce scale-dependent degradation and break naive scaling extrapolation.** We train model ladders at matched compute while sampling training documents with replacement from pools of size  $K$  (exact repeats allowed). We compare against an approximately-infinite baseline with negligible repeats. Left: train and validation loss versus compute/scale for each  $K$ . Right: fractional loss change relative to the baseline (Eq. 9). Small models scale normally under small  $K$ , while larger models exhibit rapidly increasing penalties, implying that scaling ladders can underestimate main-run loss when effective uniqueness is limited.

## 5. Theory: Scale-Dependent Effective Duplicates and Restored Scaling

### 5.1. Semantics as Hierarchical Latents and Semantic Duplicates

We model “same meaning, different surface form” via latent semantics and transformations. Let  $z$  denote a *semantic latent* (meaning), and let  $\tau$  denote a *surface transformation* (language, paraphrase, formatting, casing, etc.). A document  $x$  is generated by

$$z \sim p(z), \quad \tau \sim p(\tau | z), \quad x = \mathcal{G}(z, \tau). \quad (10)$$

Two documents  $x$  and  $x'$  are *semantic duplicates* if they share the same  $z$  but differ in  $\tau$ . This abstraction covers translations:  $x = \mathcal{G}(z, \tau_{EN})$  and  $x' = \mathcal{G}(z, \tau_{ZH})$ .

To capture compositional structure, we allow  $z$  itself to be hierarchical:

$$z^{(0)} \rightarrow z^{(1)} \rightarrow \dots \rightarrow z^{(L)} \rightarrow x, \quad (11)$$

where  $z^{(0)}$  is coarse semantics (topic/world knowledge) and  $z^{(L)}$  is closest to surface form. In this view, “duplicates” are

not a binary dataset property: two documents can share an ancestor latent at some depth but not others. A model that only learns shallow latents may treat translations as distinct, while a model that learns deeper invariances collapses them to the same effective representation.

The hierarchy in Eq. (11) is an abstract model of compositional structure, where coarser latents  $z^{(0)}$  capture broad topics/semantics while deeper latents capture increasingly fine-grained meaning and surface realization. This perspective is closely related to recent theoretical models of compositional data such as the Random Hierarchy Model (RHM), which generates examples by composing features along a tree (analogous to a grammar derivation) and predicts scale-dependent learnability of deeper levels (Cagnetta et al., 2024b). In our setting, increasing capability corresponds to learning deeper invariances in the latent hierarchy, which enlarges the set of surface variants that collide into the same effective semantic latent, increasing redundancy.

Using notation from Section 2, we formalize “duplication” in terms of training signal rather than surface form. Let  $f_\theta$  be a language model trained by next-token prediction.

**Definition 5.1** (Effective duplicates). Fix  $\varepsilon \in (0, 1)$ . We call  $x$  and  $x'$   $\varepsilon$ -effective duplicates at  $\theta$  if

$$\text{sim}_{\theta}(x, x') \geq 1 - \varepsilon. \quad (12)$$

This definition is explicitly model-dependent: as capability/scale increases, the relation (12) can merge previously distinct examples (e.g. translations).

To connect semantics to gradients, we use a minimal decomposition. Let  $z = z(x)$  denote the semantic latent for  $x$ . We write the per-document gradient as

$$g(x; \theta) = \underbrace{\mu(\theta)}_{\text{global}} + \underbrace{\delta_z(\theta)}_{\text{semantic}} + \underbrace{\xi_x(\theta)}_{\text{surface/idiosyncratic}}, \quad (13)$$

where  $\mathbb{E}[\delta_z] = 0$  and  $\mathbb{E}[\xi_x | z] = 0$ . Intuitively,  $\delta_z$  captures the update direction shared by all surface forms of the same meaning, while  $\xi_x$  captures surface-specific variations.

A convenient summary of semantic sensitivity at scale  $s$  (parameters/compute/training time) is the fraction of gradient energy explained by the semantic component:

$$\rho(s) := \frac{\mathbb{E}\|\delta_z(\theta(s))\|_2^2}{\mathbb{E}\|g(x; \theta(s)) - \mu(\theta(s))\|_2^2} \in [0, 1]. \quad (14)$$

Under mild assumptions that the surface/idiosyncratic term  $\xi_x$  is approximately isotropic and independent across different surface forms of the same latent  $z$ ,  $\rho(s)$  controls expected gradient cosine similarity. Concretely, for semantic duplicates  $x = \mathcal{G}(z, \tau)$  and  $x' = \mathcal{G}(z, \tau')$  with the same  $z$ , the numerator satisfies  $\mathbb{E}\langle g(x) - \mu, g(x') - \mu \rangle \approx \mathbb{E}\|\delta_z\|_2^2$ , while the denominator is  $\mathbb{E}\|g(x) - \mu\|_2^2 \approx \mathbb{E}\|\delta_z\|_2^2 + \mathbb{E}\|\xi_x\|_2^2$ , yielding the approximation

$$\begin{aligned} \mathbb{E}[\text{sim}_{\theta(s)}(x, x') | z] &\approx \rho(s), \\ \mathbb{E}[\text{sim}_{\theta(s)}(x, \tilde{x})] &\approx 0 \text{ for unrelated } \tilde{x}. \end{aligned} \quad (15)$$

Our gradient experiments (Section 2) provide direct empirical evidence that  $\rho(s)$  increases with both training progress and model capability: transformations that preserve  $z$  (e.g. translations) become increasingly aligned in gradient space.

## 5.2. Replication, Redundancy, and Effective Uniqueness

Consider training on a stream constructed by sampling *with replacement* from an underlying distribution over semantic latents  $z \in \mathcal{Z}$  with mixture weights  $\{w_z\}$ . In our controlled experiments (Section 4), this corresponds to uniform sampling from a pool of  $K$  unique documents (so  $w_z = 1/K$ ), but the latent view also covers non-uniform frequencies.

A key quantity is the (Simpson) latent collision probability (Simpson, 1949)

$$p_{\text{lat}} := \mathbb{P}(z = z') = \sum_z w_z^2,$$

and the associated effective latent count

$$K_{\text{eff}} := \frac{1}{p_{\text{lat}}} = \frac{1}{\sum_z w_z^2}. \quad (16)$$

(When  $w_z \equiv 1/K$ , we have  $K_{\text{eff}} = K$ .)

Let  $x_1, \dots, x_n$  be  $n$  iid draws from this mixture and define the averaged centered gradient  $\bar{g}_n := \frac{1}{n} \sum_{t=1}^n (g(x_t; \theta) - \mu(\theta))$ . Assume the following simplified correlation structure consistent with Eq. (13)

$$C(x, x') \approx \begin{cases} \sigma^2 & z(x) = z(x') \text{ and } x = x', \\ \rho(s) \sigma^2 & z(x) = z(x') \text{ and } x \neq x', \\ 0 & z(x) \neq z(x'), \end{cases} \quad (17)$$

where  $C(x, x') \equiv E\langle g(x; \theta) - \mu, g(x'; \theta) - \mu \rangle$ .

**Proposition 5.2** (Saturation of independent training signal). Under (17) and uniform sampling over  $K$  classes,

$$\begin{aligned} \mathbb{E}\|\bar{g}_n\|_2^2 &\approx \frac{\sigma^2}{n} (1 + \rho(s)(n-1)p_{\text{lat}}) \\ &= \frac{\sigma^2}{n} \left(1 + \rho(s) \frac{n-1}{K_{\text{eff}}}\right). \end{aligned} \quad (18)$$

Equivalently, the averaged gradient behaves like an iid average with effective sample size

$$\begin{aligned} n_{\text{eff}}(n, K_{\text{eff}}; s) &:= \frac{n}{1 + \rho(s)(n-1)/K_{\text{eff}}} \\ &\approx \min\left\{n, \frac{K_{\text{eff}}}{\rho(s)}\right\}. \end{aligned} \quad (19)$$

**Interpretation.** When  $n \ll K/\rho(s)$ , redundancy is negligible and signal scales like  $1/n$ . When  $n \gg K/\rho(s)$ , semantic redundancy dominates and the number of effectively independent update directions saturates at  $K/\rho(s)$ . Because  $\rho(s)$  increases with capability (Section 2), the same finite- $K$  stream becomes *more redundant* for larger/stronger models, i.e. effective uniqueness  $K/\rho(s)$  shrinks with scale.

## 5.3. From Effective Reuse to a Restored Scaling Law

Let  $C$  denote training compute. Let  $L(C, K_{\text{eff}})$  be eval loss when sampling with replacement from an effective semantic pool size  $K_{\text{eff}}$ , and let  $L_{\infty}(C)$  be the baseline with effectively infinite uniqueness (negligible repeats). Define the normalized degradation

$$\Delta(C, K) := \frac{L(C, K) - L_{\infty}(C)}{L_{\infty}(C)}. \quad (20)$$

The redundancy picture suggests that the relevant control variable is an *effective reuse ratio*

$$r_{\text{eff}}(C, K_{\text{eff}}) := \frac{\rho(C)n(C)}{K_{\text{eff}}}, \quad (21)$$

where  $n(C)$  is the number of documents trained on at compute  $C$ , and  $\rho(C)$  captures semantic alignment (Section 5.1).

**Assumption (Power Law Penalty in Effective Reuse).** Over the regime where scaling laws are measured, we posit

$$\Delta(C, K) \approx \lambda r_{\text{eff}}(C, K)^\eta. \quad (22)$$

This is a parsimonious way to encode that (i) no penalty occurs when reuse is negligible and (ii) penalty grows smoothly with semantic redundancy.

**Compute Dependence and the Plane Law.** Over a limited compute range, we approximate both  $n(C)$  and  $\rho(C)$  by power laws

$$n(C) \propto C^u, \quad \rho(C) \propto C^v. \quad (23)$$

Then (22) yields

$$\Delta(C, K_{\text{eff}}) \approx a C^\beta K_{\text{eff}}^{-\gamma}, \quad \beta = \eta(u+v), \quad \gamma = \eta, \quad (24)$$

where  $a > 0$  absorbs constants. Equation (24) is the *minimal global scaling correction* compatible with: (i) reuse increasing with compute ( $u > 0$ ), and (ii) semantic sensitivity increasing with compute ( $v \geq 0$ ). A special *ratio-only* law  $\Delta \propto (\sqrt{C}/K)^\eta$  corresponds to  $u = 1/2$  and  $v = 0$ , which can be too restrictive when  $\rho(C)$  grows with scale.

**Restored Predictivity.** Combining (20) and (24) gives the restored loss prediction:

$$L_{\text{pred}}(C, K_{\text{eff}}) = L_\infty(C) \left(1 + a C^\beta K_{\text{eff}}^{-\gamma}\right). \quad (25)$$

In our experiments (Section 4),  $L_\infty(C)$  is measured directly from the ‘‘approximately infinite unique data’’ runs at the same compute, so restoring predictivity requires fitting only  $(a, \beta, \gamma)$ . In App. E, we provide a collision-aware scaling correction can be derived by combining a Hutter-style learning curve (Hutter, 2021) with an effective-sample-size reduction induced by duplicate/semantic-collision gradients.

**Empirical Validation and Minimality.** On our controlled scaling ladders, the 3-parameter plane law (24) accurately predicts *all* eval losses across  $(C, K)$ , including the breakdown regime, with small average relative error, whereas the 2-parameter ratio-only constraint can substantially underpredict the catastrophic  $K = 10^5$  main run. This supports the interpretation that semantic sensitivity  $\rho(C)$  contributes nontrivially to the compute exponent  $\beta$ .

#### 5.4. Estimating an Effective Semantic Pool Size from Mean Nearest-Neighbor Cosine

In real pretraining, the ‘‘number of unique semantic items’’  $K$  is not directly observable. However, our restored scaling law only requires an *effective uniqueness*—the rate at which training samples collide under the model’s semantic resolution. Here we show how to estimate an effective

$K_{\text{eff}}$  using only a *mean nearest-neighbor cosine* statistic computed from embeddings of the *sampled training stream* (which includes repeats).

**Setup: Cosine is Measured on a Fixed Embedding Subsample of the Stream.** For each training run we take a subsample of  $N_{\text{meas}}$  training documents from the run’s data stream (including repeats), embed each document with a fixed embedding model, and unit-normalize to obtain vectors  $v_t \in \mathbb{S}^{d-1}$ . We then compute the nearest-neighbor cosine for each embedded sample

$$M_t := \max_{s \neq t} \langle v_t, v_s \rangle, \quad \overline{M}_{N_{\text{meas}}} := \frac{1}{N_{\text{meas}}} \sum_{t=1}^{N_{\text{meas}}} M_t.$$

All quantities below refer to this fixed measurement size  $N_{\text{meas}}$ . (In our controlled ladder,  $N_{\text{meas}}$  is constant across runs; if  $N_{\text{meas}}$  is not logged, it can be inferred from the small- $K$  regime where exact repeats are frequent.) Crucially,  $\overline{M}_N$  is computed on the *stream* of size  $N$ , which depends on  $C$  (through the number of examples processed), not on the unknown pool size  $K$ .

#### Step 1: Background NN Similarity without Collisions.

Let  $m_0(N)$  denote the expected mean NN cosine when the nearest neighbor is *not* a semantic collision (i.e., no same-latent partner appears among the  $N - 1$  other samples). In practice we estimate  $m_0(N)$  from a high-uniqueness reference stream (largest- $K$  pool or without-replacement stream), where exact repeats are negligible, using the same embedding pipeline.

**Step 2: A Two-Component Model for  $\overline{M}_N$ .** We model each  $M_t$  as either: (i) a background neighbor with mean  $m_0(N)$ , or (ii) a collision neighbor (same latent) with typical similarity  $m_+ \in (m_0(N), 1]$ . Let  $q_N$  be the probability that a given sample has at least one collision neighbor among the other  $N - 1$  samples. Then

$$\mathbb{E}[\overline{M}_N] \approx (1 - q_N) m_0(N) + q_N m_+. \quad (26)$$

Solving gives the estimator

$$\hat{q}_N := \text{clip}\left(\frac{\overline{M}_N - m_0(N)}{m_+ - m_0(N)}, 0, 1\right). \quad (27)$$

In our controlled experiment where collisions correspond to *exact repeats*, we take  $m_+ = 1$  (up to numerical precision). For semantic (non-exact) collisions,  $m_+ < 1$  can be calibrated using known semantic-duplicate pairs (e.g. translations).

#### Step 3: Invert $\hat{q}_N$ into an Effective Latent Count $K_{\text{eff}}$ .

Let  $z$  be the semantic latent with mixture weights  $\{w_z\}$  and collision probability  $p_{\text{lat}} = \mathbb{P}(z = z') = \sum_z w_z^2$ . Define the effective number of latents (Simpson effective size)

$$K_{\text{eff}} := \frac{1}{p_{\text{lat}}} = \frac{1}{\sum_z w_z^2}. \quad (28)$$

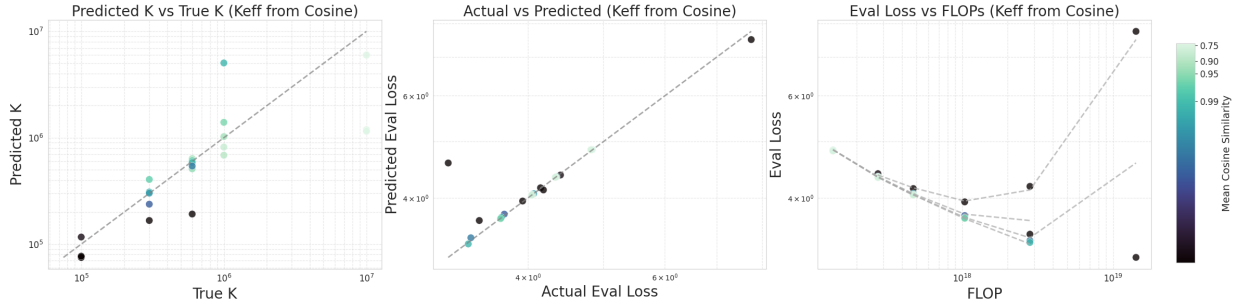


Figure 7. **Predictable scaling can be restored by accounting for limited semantic diversity:** We use dataset size and mean cosine similarity to estimate  $K$  via Equation (29) (left). We then plug our estimate of  $\hat{K}_{\text{eff}}$  into Equation (25) to estimate the loss (Center). This produces scaling curves that align closely with the empirical eval losses (right).

For a latent mixture with weights  $\{w_z\}$ , the probability that a given draw has at least one same-latent partner among the other  $N_{\text{meas}} - 1$  draws is

$$q_{N_{\text{meas}}} = 1 - \sum_z w_z (1 - w_z)^{N_{\text{meas}} - 1} \approx 1 - \exp\left(-\frac{N_{\text{meas}} - 1}{K_{\text{eff}}}\right),$$

where the approximation holds when the mixture has no heavy modes (all  $w_z \ll 1$ ); in the uniform- $K$  case it is exact up to the standard  $\log(1 - x) \approx -x$  approximation, see App. D. Inverting yields

$$\hat{K}_{\text{eff}} := \frac{N_{\text{meas}} - 1}{-\log(1 - \hat{q}_{N_{\text{meas}}})}. \quad (29)$$

**Step 4: A  $K$ -free Restored Scaling Law.** Our restored degradation model is

$$\Delta(C, K) := \frac{L(C, K) - L_{\infty}(C)}{L_{\infty}(C)} \approx a C^{\beta} K^{-\gamma}.$$

Replacing  $K$  by  $\hat{K}_{\text{eff}}$  gives a correction depending only on observable stream geometry  $\Delta(C) \approx a C^{\beta} \hat{K}_{\text{eff}}^{-\gamma}$  as

$$L_{\text{pred}}(C) = L_{\infty}(C)(1 + \Delta(C)). \quad (30)$$

**Validation on the Controlled Ladder.** On the common evaluation set of runs in our controlled  $K$ -pool experiment, the plane law using the true pool size  $K$  achieves mean absolute relative error  $\approx 0.77\%$  (median  $\approx 0.28\%$ ). Replacing  $K$  with  $\hat{K}_{\text{eff}}$  estimated from mean NN cosine via Eqs. Equation (26)–(29) achieves  $\approx 0.90\%$  (median  $\approx 0.24\%$ ). Thus, even with access only to a mean cosine statistic,  $\hat{K}_{\text{eff}}$  recovers most of the predictivity of the true- $K$  scaling correction. See Figure 7.

*Remark 5.3* (Identifiability from mean NN cosine). The mapping  $\overline{M}_{N_{\text{meas}}} \mapsto \hat{q}_{N_{\text{meas}}}$  requires specifying both a background term  $m_0(N_{\text{meas}})$  and a collision-similarity level  $m_+$ . With only the mean NN cosine available,  $m_+$  cannot

be identified without external calibration; in the controlled with-replacement experiment, exact repeats imply  $m_+ \approx 1$ , while for semantic (non-exact) collisions one can calibrate  $m_+$  using known semantic-duplicate pairs (e.g. translations) embedded by the same model.

## 6. Discussion and Future Directions

We discover an insidious source of scale-dependence that can impact the training of large language models, but not smaller language models: as model capabilities increase, training signals from semantically equivalent documents align. Thus, semantically equivalent documents in the corpora may act similarly to exact duplicates, harming model quality. Moreover, as training data scale, the number of semantic collisions increases far more quickly than one would expect based on trends gleaned from small corpora. We model this effect on small language models and propose scaling laws that account for semantic diversity in the dataset, restoring predictable scaling.

Our experiments have profound implications for the future of language models. Until now, industry convention has been to bet trillions of dollars on the success of the bitter lesson: scale, scale, scale, and super-intelligence will follow (Sutton, 2019). The only obstruction on this path has been the limited number of training data in web-scale corpora. Frontier labs have tried to sidestep this obstacle by synthesizing massive corpora comprised of LLM-generated text. Our findings tell a cautionary tale about this approach: even if one can scale the raw number of tokens to asymptotically high regimes, semantic diversity may be just as important as data volume. As we show in Figure 5, synthetic data scales poorly with respect to semantic diversity. Our experiments emphasize the importance of seeding semantic diversity in synthetic data. There is only one other path: if the sum total of extant semantically distinct human thoughts is insufficient to train modern LMs, then labs must invest in more data-efficient training and architectures. We discuss limitations and future work in Appendix B.

440 **Impact Statement**

441 This paper presents work whose goal is to advance the field  
442 of machine learning. There are many potential societal  
443 consequences of our work, none of which we feel must be  
444 specifically highlighted here.  
445

446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494

References

Adcock, A., Srivastava, A., Dubey, A., Jauhri, A., Pande, A., Pandey, A., Sharma, A., Kadian, A., Kumawat, A., Kelsey, A., Stelle, A., Cheema, A., Kabiljo, A., Katz, A., Gangidi, A., Tayade, A., Victoria, A., Alastuey, A. S., Conrath, A., Mohiuddin, A., Sharif, A., Siddiqui, A., Goldstand, A., Li, A., Boyd, A., Daliri, A. K., Iqbal, A., Menon, A., Mathews, A., Mathur, A., Agarwal, A., Schelten, A., Shine, A., Muñoz, A. C., Guliaev, A., Radovic, A., Song, A., Vaughan, A., Simeonov, A., Rezende, A., Rezende, A., Baevski, A., Roubaud, A., Ma, A., Lee, A., Pereira, A., Ahmed, A., Shankar, A., Kallet, A., Budhiraja, A., Khandekar, A., Benhalloum, A., Gershman, A., Nagpal, A., Zohar, A., Sharaf, A., Desai, A., Razdaibiedina, A., Agape, A., Kurghinyan, A., Perunicic, A., Madotto, A., Darabanov, A., Alvarado, A., Brown, A., Cohen, A., Fang, A., Freeman, A., Gallagher, A., Gu, A., Jo, A. P., Ryan, A., Steffen, A., Wei, A., Rusakov, A., Golovei, A., Shang, A., Fan, A., Fan, A., Flewelling, A., Pathak, A., Goyal, A., Ramchandani, A., Pai, A., Singh, A., Garg, A., Xing, A., Cai, A., Grosul, A., Prochowska, A., Sun, A., Dong, A., Franco, A., Hu, A., Chawla, A., Hartshorn, A., Sheng, A., Thomas, A., Goyal, A., De, A., Bodiwala, A., Bodiwala, A., Yang, A., Saraf, A., Samudra, A., Mun, A., Rahnema, A., Mitra, A., Sravankumar, A., Gupta, A., Haghighi, A., Stolerman, A., Chowdhury, A., Choudhury, A., Korenev, A., Guo, A., Hinsvark, A., Mallya, A., Neelakantan, A., Talebzadeh, A., Shah, A., Shetty, A. J., Bharambe, A., Islam, A., Zhang, A., Gregerson, A., Lewis, A., Ibrahim, A., Minhas, A., Dahan, A., Dabah, A. R., Tang, B., Ulman, B., Sadeghi, B., Jedrzejewski, B., Skarabahaty, B., Zhu, B., Li, B., Bharier, B., Leonhardi, B., Muller, B., Plessala, B., Huang, B., Loyd, B., Paranjape, B., Sheth, B., Bonner, B., Holland, B., Wang, B., Liu, B., Tang, B., Liu, B., Wu, B., Li, B., Yu, B., Chen, B.-C., Araya, B., Vidolov, B., Chen, B., Peng, B., Ni, B., Davis, B., Wasti, B., Adams, B., Taylor, B., Wu, B., Swidler, B., Chiang, B., Clerkin, B., Fuller, B., Cutter, B., Novais, B., Gmyrek, B., Easton, B., Campos, C., Case, C., Fu, C. C., Burton, C., Diaz, C., Cole, C., Liu, C., Fougere, C., Peng, C., Peng, C., Zhao, C., Wang, C., Kim, C., Shaib, C., Zhou, C., Caucheteux, C., Nguyen, C., Sitawarin, C., Nayak, C., Asher, C., Fan, C., Zhu, C., Cheng, C., Zhang, C., Zhu, C., Ruan, C., Yu, C., Hua, C., Whitehouse, C., Holloway, C., Chu, C.-H., Chuang, C.-Y., Karande, C., Nagpal, C., Bakalar, C., Bi, C., Cai, C., Marra, C., McConnell, C., Thi, C., Tindal, C., Waterson, C., Deverall, C., Fuegen, C., Keller, C., Cheng, C., Jou, C., Smith, C., Wang, C., Feichtenhofer, C., Touret, C., Luc, C., Sauper, C., Zhuge, C., Sung, C.-Y., Tang, C., Wu, C., Siegel, C., Heale, C., Wilbourn, C., White, C., Xia, C., Wong, C., Rat, C., Ferrer, C. C., Habis, C., Nikolaidis, C., Lohachov, D.,

Ju, D., Flanagan, D., Allonsius, D., Civin, D., Johnson, D., Bolya, D., Francisco, D., Fried, D., Hawthorne, D., Haziza, D., Ho, D., Kreymer, D., Li, D., Machlab, D., McKinnon, D., Obenshain, D., Rodriguez, D., Song, D., Tse, D., Pintz, D., Livshits, D., Rodrigo, D. J., Huynh, D., Askarov, D., Brandfonbrener, D., Esiobu, D., Kant, D., Levin, D., Renardy, D., Soofian, D., Stevens, D., Xu, D., Zhang, D., Shah, D., David, D., Douglas, D., Boyda, D., Raj, D., Hazarika, D., Mekala, D., Choudhary, D., Mahajan, D., Jin, D., Coll-Vinent, D. S., Foss, D., Garcia-Olano, D., Perino, D., Hupkes, D., Su, D., Madathil, D., Govindasamy, D., Yeduguru, D., Vengertsev, D., He, D., Li, D., Wang, D., Li, D., Le, D., Hin, D., Holland, D., Nguyen, D., Nguyen, D., Dowling, E., Litt, E., Lakomkin, E., AlBadawy, E., Ardestani, E. K., Eckstein, E., Dabir, E., Montgomery, E., Lobanova, E., Abramoviz, E., Hedeman, E., Li, E., Hilbert, E., Tan, E. X., Yun, E., Stener, E., Stoimenov, E., Garreau, E., Dinan, E., Hahn, E., Wood, E., Li, E., Ademuwagun, E., Seker, E., Alamillo, E., Gan, E., Han, E., Huang, E., Smith, E. M., Le, E.-T., Chang, E., Helenowski, E., Elnikety, E., Arcaute, E., Myers, E., Nho, E., Poliukhovych, E., Dunbar, E., Litvinenko, E., Altıntaş, E., Hochman, E., Shtrauch, E., Mastenbroek, F., Zeb, F., Ahmad, F., Farahbakhshian, F., Kou, F., Sun, F., Chen, F., Chung, F., Tian, F., Xu, F., Radenovic, F., Kokkinos, F., Barbieri, F., Caggioni, F., Esparza, F., Guzmán, F., Kanayet, F., Seide, F., Zhang, F., Lewis, F., Huang, F., Wang, F., Synnaeve, G., Jacques-Silva, G., Schwarz, G., Ghardhora, G., Elfer, G., Dickson, G., Chaurasia, G., Sewani, G., Shingi, G., Zuo, G., Jeong, G., Puthanpurackal, G., Swee, G., Bertran, G. M.-T., Keren, G., Ling, G., Stasa, G., Saha, G., Safran, G., French, G., Rajendran, G., Thattai, G., Cineas, G., Nail, G., Fletcher, G., Mialon, G., Adams, G., Sizov, G., Pang, G., Elsahar, H., Tran, H. D., Nguyen, H., Wu, H., Inan, H., Eghbalzadeh, H., Fang, H., Zou, H., Doyle, H., Korevaar, H., Wang, H., Werbel, H., Zha, H., Morsy, H., Ma, H., Zhang, H., Sun, H., Wang, H., Shah, H., Habeeb, H., Rudolph, H., Gupta, H., Poddar, H., Parikh, H., Zhang, H., Wang, H., Li, H., Sharma, H., Nguyen, H. P., Zhang, H., Qiu, H., Lv, H., Xu, H., Zhan, H., Hamooni, H., Huang, H., Xu, H., Laurençon, H., Touvron, H., Dinh, H., Goldman, H., Mehanna, H., Nguyen, H., Tsuo, H., Graves, I., Yu, I., Damlaj, I., Cohen, I., Tufanov, I., Goldenstein, I., Leontiadis, I., Zarov, I., Ahmed, I., Djiofack, I., Spulber, I., Veliche, I.-E., Ramos, I., Misra, I., Gal, I., Evtimov, I., Evtimov, I., Obraztsov, I., Wu, J., Vertino, J. R., Koo, J., Lee, J., Jung, J., Weissman, J., Beldock, J., Crnkovich, J., Grinage, J., Zeng, J. H., Kohli, J., Tian, J., Cahill, J., Gefert, J., Seidel, J., Seidel, J., Tracey, J., Cho, J. H., Wei, J., Kahn, J., Howell, J., Vu, J. L., Park, J., Yan, J., Yip, J., Li, J., Mahadeokar, J., Goluguri, J. B. R., Mehar, J., Gaya, J.-B., Shah, J., Hanson, J., Marcus, J., Walsh, J., Yang, J., van der Linde, J., Fan, J., Chan, J., Zhen, J., Lee, J.,

550 Fu, J., Reizenstein, J., Teboul, J., He, J., Zhong, J., Hou,  
551 J., Yang, J., Ding, J., Hu, J., Zhu, J., Guo, J., Wang, J.,  
552 Ouyang, J., Chi, J., Huang, J., Zhao, J., Yang, J., Zhou, J.,  
553 Zhao, J., Liu, J., Wang, J., You, J., Yu, J., Schwiep, J., Wu,  
554 J., Huang, J., Li, J., Koh, J. Y., Zhang, J., Chen, J., Yang,  
555 J., Shen, J., Hwang, J., Guo, J., Khatiawada, J., Bitton, J.,  
556 Li, J., Quanaim, J., Beales, J., Schuijt, J., Chang, J., Quan,  
557 J., Chan, J., Shepard, J., Harris, J., Rubin, J., Janzen, J.,  
558 Kaldor, J., Silva, J. L., Leitao, J., Greer, J., Moon, J.,  
559 Rocca, J., Tighe, J., Fromm, J., Deng, J., Fernandes, J.,  
560 Saxe, J., Zheng, J., Pino, J., Prigent, J., Chen, J., Tian,  
561 J., Qi, J., Wang, J., Jia, J., Baker, K., Londenberg, K.,  
562 Wang, K., Peng, K., Peng, K., Yang, K., Alwala, K. V.,  
563 Yu, K. H., Narang, K., Chadha, K., Sikka, K., Zhang,  
564 K., Schuberts, K., Mandyam, K., Sankararaman, K. A.,  
565 Padthe, K., Prasad, K., Sivakumar, K., Upasani, K., Plaw-  
566 iak, K., Saenko, K., Žmolíková, K., Stadler, K., Matosich,  
567 K., Doulgass, K., Hassani, K., Ji, K., Li, K., Heafield, K.,  
568 Yu, K., Li, K., Ma, K. C.-Y., Hannan, K., Man, K., Chen,  
569 K., El-Arini, K., Hutsulyak, K., Nash, K., Jagadeesh, K.,  
570 Bartelt, K., Topaloglou-Mundy, K., Chatziioannou, K.,  
571 Karanasos, K., Vougioukas, K., Tsiampouris, K., Hamill,  
572 K., Choi, K., Iyer, K., Malik, K., Chiu, K., Huang, K.,  
573 Bhalla, K., Chawla, K., Li, K., Lakhotia, K., Monk, K.,  
574 Garg, L., Chourey, L., Hamre, L., Gustafson, L., Deason,  
575 L., Rouesnel, L., van der Maaten, L., A, L., Chen, L.,  
576 Jang, L., Silva, L., Sari, L., Hetherington, L., Zhang, L.,  
577 Zhao, L., Chen, L., Li, L. C., Yang, L., Zhan, L., Corallo,  
578 L., Tan, L., Yu, L., Liu, L., Mor, L., Lin, L., Li, L., Titus,  
579 L., Jenkins, L., Madaan, L., Fang, L., Yuan, L., Nava, L.,  
580 Pasqualin, L., Switzer, L., Fang, L., Sun, L., Tadic, L.,  
581 Blecher, L., Landzaat, L., Zhang, L., Rao, M., Khabsa,  
582 M., Miller, M., Kariya, M., Pasupuleti, M., Luthra, M.,  
583 Faruqui, M., Avlani, M., Wang, M., Singh, M., Paluri, M.,  
584 Chakkaravarthy, M., Nair, M., Tiffany, M., Pawlowski,  
585 M., Wu, M., Lomeli, M., Consuegra, M., Boiteux, M.,  
586 Galanis, M. A., Chen, M., Gleize, M., Fazel-Zarandi, M.,  
587 Hasson, M., Oldham, M., Rita, M., Dordal, M., Setzler,  
588 M., Staats, M., Staats, M., Wilde, M., Clark, M., Grange,  
589 M., Lennie, M., Schmohl, M., Raphael, M., Naumov,  
590 M., Samoylov, M., Lecanu, M., Pavlova, M., Jawaid, M.  
591 T. B., Keneally, M., Kambadur, M., Zhang, M., Liu, M.,  
592 Lin, M., Wang, M., Abraham, M., Liu, M., Au-Yeung,  
593 M., Feldergraf, M., Man, M., Matheny, M., Suo, M.,  
594 Tontchev, M., Meyer, M., Ma, M., Patel, M., Kale, M. S.,  
595 Vyatskov, M., Alexander, M., Andersland, M., Clark, M.,  
596 Lewis, M., Li, M., Macey, M., Macey, M., Seltzer, M.,  
597 Fernandez, M. J., Antonov, M., Plekhanov, M., Zhou,  
598 M., Si, M., Qiao, M., Ma, M., Zhang, M., Liang, M.,  
599 Hermoso, M. J., Suzgun, M., Skarica, M., Singh, M. K.,  
600 Kabbani, M., Rastegari, M., Sarantakos, M., Sim, M.,  
601 Gangapuram, M., Moshe, M., Doulaty, M., Metanat, M.,  
602 Chen, M., Kumar, M., Bansal, M., Ramarao, M., Li, N.,  
603 Azaria, N., Malik, N., Goyal, N., Balderas, N. V., Wang,  
604 N., Kanda, N., Gimelshein, N., Neverova, N., Aclander,  
N., Sithiviraporn, N., Kumar, N. M., Newton, N., Bahl,  
N., Ghorbani, N., Patel, N., lee Golan, N., Longenbaugh,  
N., Egebo, N., Johri, N., Mehta, N., Naik, N., Moritz,  
N., Bashlykov, N., Bogoychev, N., Laptev, N. P., Chat-  
terji, N., Jones, N., Shah, N., Dong, N., Li, N., Li, N.,  
Zhang, N., Yadav, N., Paz, N., Cheng, N., Cheng, N.,  
Adesanya, O., Repin, O., Maksymets, O., Salpekar, O.,  
Harosh, O., Pednekar, O., Çelebi, O., Gafni, O., Edinger,  
O., Hanna, O., Mohammed, O. K., Kalinli, O., Tomasello,  
P., Singh, P., Quevedo, P., Jain, P., Rashidinejad, P., Too-  
ley, P., Parekh, P., Thakkar, P., Taheri, P., Hapuarachchi,  
P., Kesseli, P., Alrassy, P., de Rezende Pinatti, P., Balaji,  
P., Sisodiya, P., Moreira, P. J. F., Rittner, P., Valenzuela,  
P., Sun, P., Zhang, P., Chen, P.-J., Wang, P., Zhang, P.,  
Li, P., Vasic, P., Carras, P., Ney, P., Weng, P., Dumea,  
P., Hayes, P., Woods, P., Andrews, P., Ménard, P., Wu,  
P.-H., Liu, P., Dollar, P., Dzhelepev, P., Zvyagina, P., A,  
P., Agrawal, P., Rajendran, P., Prakash, P., Bhargava, P.,  
Pramono, Shah, P., Dave, P., Jain, P., Dubal, P., Gollakota,  
P., Krishnan, P., Yuvraj, P., Ghosh, P., Koura, P. S., Xu,  
P., Qi, Q., Zhou, Q., Guan, Q., Sun, Q., Liu, Q., He, Q.,  
Zheng, Q., Yang, Q., Guo, Q., You, Q., Carbonneaux, Q.,  
Carbonneaux, Q., Duval, Q., Fettes, Q., Alao, R., Batish,  
R., Guo, R., Rodriguez, R., Bhargava, R., Asuncion, R.,  
Murthy, R., Dutta, R., Jha, R., Kindi, R., Mitra, R., Gana-  
pathy, R., Shah, R., Das, R., Shrivastava, R., Nishtala, R.,  
Shankar, R., Shukhau, R., Calderer, R., Parthasarathy, R.,  
Subramanian, R., Bensadoun, R., Bostan, R., Chaturvedi,  
R., Agrawal, R., Gao, R., Li, R., Kogen, R., Duran, R.  
J. P., Cabral, R. S., Lee, R., Pang, R. Y., Bhalodia, R.,  
Mansour, R., Singh, R., Godugu, R., Patney, R., Boyle,  
R., Goldfarb, R., Caldwell, R., Kuo, R., Raileanu, R.,  
Battey, R., Sharma, R., Sapra, R., Wang, R., Granata,  
R., Castro, R. D., Paim, R., Maheshwari, R., Varma, R.,  
Girdhar, R., Patel, R., Sumbaly, R., Sheaffer, R., Silva,  
R., Buchillon, R. R., Hou, R., Xie, R., Mavlyutov, R.,  
Semenov, R., Dinov, R., Bao, R., Fox, R., Kilpatrick,  
R., Kwan, R., Lim, R., Smith, R., Narayan, S., Qiao, S.,  
Mehta, S., Siby, S., Jain, S., Hosseini, S., Gur-Ari, S.,  
Chennabasappa, S., Geyik, S., Bondu, S. J., Nekkhalapudi,  
S. M. C., Hasan, S., Okabayashi, S., Rambhatla, S., Sawh-  
ney, S., Dunster, S., Zhao, S., Keon, S., Azadi, S., Sapra,  
S., Dooley, S., Datta, S., Parab, S., Xie, S. M., Singh, S.,  
Chen, S., Behn, S., Khodeir, S., Shirazyan, S., Dhillon, S.,  
Pumma, S., Sidorov, S., Adaime, S., Khanna, S., Wani,  
S., Brenton, S., Bell, S., Kelly, S., Koger, S., Nunley,  
S., Perry, S., Caicedo, S., Dahlgren, S., Ruder, S., Ya-  
mamoto, S., Mehretu, S., Ravi, S. S., Lyu, S., Chellapan,  
S., Mellos, S., Edunov, S., Royt, S., Cohen, S., Peng, S.,  
Adams, S., Nie, S., Ramaswamy, S., Narang, S., Pisupati,  
S., Gandham, S., Lim, S., Lindsay, S., Artrip, S., Sheynin,  
S., Yan, S., Feng, S., Shen, S., Zheng, S., Lin, S., Bi, S.,  
Zha, S. C., Wan, S., Qian, S., Cai, S., Shao, S., Shahidi,

- 605 S., Li, S., Bernholtz, S., Wang, S., Patil, S. G., Verma,  
 606 S., P. S. S., Chen, S., Yaida, S., Debnath, S., Siravara, S.,  
 607 Bhosale, S., Ma, S., Zhang, S., Tang, S., Zhang, S., Zhou,  
 608 S., Che, S., Srinivisan, S., Bhattacharya, S., Patki, S.,  
 609 Chen, S., Chen, S., Vandenhende, S., Merello, S., Wang,  
 610 S., Barzily, S., Yi, S., Lin, S., Bong, S., Yin, S., Agarwal,  
 611 S., Agarwal, S., Lieve, S., Sajuyigbe, S., Jiang, S., Li, S.,  
 612 Kim, S., Khosla, S., Maiti, S., Whitman, S., Popuri, S.,  
 613 Tallam, S., Vaidyanathan, S., Vaidyanathan, S., Sootla,  
 614 S., Collot, S., Ding, S., Chen, S., Cai, S., Gururangan, S.,  
 615 Govindaprasad, S., Young, S., Dewakar, S., Gonugondla,  
 616 S. K., Bhandari, S., Gumudavelli, S., Gumudavelli, S.,  
 617 Gupta, S., Deng, S., Cho, S., Ganapathy, S., Dhal, S.,  
 618 Fedynak, S., Contrera, S., Kim, S., Rebuffi, S., Chahande,  
 619 T., Herman, T., Li, T., Xu, T., Fowler, T., Sheasha, T.,  
 620 Anand, T., Kalluri, T., Singh, T., Shavrina, T., Li, T., Rao,  
 621 T., Patil, T., Li, T., Bui, T., Quach, T., Alharbash, T., Vo,  
 622 T. V., Kooburat, T., Koehler, T., Georgiou, T., Scialom, T.,  
 623 Ye, T., Li, T., Zhang, T., Li, T., Blankevoort, T., Willi, T.,  
 624 Chou, T., Leung, T., Lee, T., Mihaylov, T., Heatwole, T.,  
 625 Xiao, T., Cao, T., Lee, T., Le, T., Rice, T., Chan, T. K. S.,  
 626 Tran, T., Tiplea, T., Baumgartner, T., Savagaonkar, U.,  
 627 Karn, U., Araiza, U. M., Farooq, U., Cohen, U., Sharif,  
 628 U., Murarka, U., Phung, V., Joginpalli, V., Saravagi, V.,  
 629 Sharma, V., Viswamurthy, V., Goswami, V., Seth, V.,  
 630 Ramesh, V., Ramesh, V., Gupta, V., Montanez, V., Natarajan,  
 631 V., Sarma, V., Ramanathan, V., Kerkez, V., Rao, V.,  
 632 Gonguet, V., Mauge, V., Do, V., Vogeti, V., Chaudhary,  
 633 V., Sankaran, V., Albiero, V., Miglani, V., Pai, V., Cojanu,  
 634 V., Shubin, V., Mihailescu, V. T., Petrovic, V., Ivanov, V.,  
 635 Vorotilov, V., Bhutada, V., Ng, W. I., Cheng, W., Sun,  
 636 W., Tu, W., Wei, W., Zhou, W., Hsu, W.-N., Chu, W.,  
 637 Yuan, W., Wang, W., Zhao, W., Jiang, W., Fu, W., Jiang,  
 638 W., Meers, W., Constable, W., Wang, W., Wong, W. R.,  
 639 Martinet, X., Lin, X. V., Yan, X., Yin, X., Li, X., Rui,  
 640 X., Yang, X., Tang, X., Wang, X., Wang, X., Wang, X.,  
 641 Dai, X., Peng, X., Li, X., Meng, X., Zhang, X., Xia, X.,  
 642 Jin, X., xinbo Gao, Xie, X., Zhou, X., Ma, X., Ju, X.,  
 643 Zhao, X., Liu, X., Jia, X., Zhang, X., Cao, X., Wang,  
 644 X., Wu, X., Xu, X., Ma, X., Wang, X., Cui, Y., Chen,  
 645 Y., Li, Y., Shu, Y., Xia, Y., Chen, Y., Zhou, Y., Mehta,  
 646 Y., Patel, Y., Tekena, Y., Gaur, Y., Babaei, Y., Zhou, Y.,  
 647 Hu, Y., Qi, Y., Lee, Y., Wen, Y., Liu, Y.-C., Wu, Y. B.,  
 648 Pan, Y., Yang, Y., Lin, Y.-H., Wang, Y., Wu, Y., Yang,  
 649 Y., Huang, Y., Aharon, Y. B., Yang, Y., You, Y., Xu, Y.,  
 650 Zhang, Y., Yuan, Y., Liu, Y., Ma, Y., Yang, Y., Lu, Y.,  
 651 Komornik, Y., Lin, Y., Goyhman, Y., Mamo, Y. M., Nam,  
 652 Y., Wang, Y., Lu, Y., Zhao, Y., Hsieh, Y.-H., Lo, Y.-J.,  
 653 Tian, Y., Zhang, Y., Xiong, Y., Yao, Y., Hao, Y., Zhang,  
 654 Y., Li, Y., Cao, Y., Yu, Y., Zhao, Y., Guo, Y., Wang, Y.,  
 655 Huang, Y., Lu, Y., Shi, Y., Wang, Y., He, Y., Wang, Y.,  
 656 Qian, Y., Wang, Y., Tang, Y., Mao, Y., Li, Y., Dai, Y.,  
 657 Hulovatyy, Y., Hu, Y., Sun, Y., Rait, Z., Wentz, Z., Coudert,  
 658 Z. D., Collins, Z., Hankir, Z., He, Z., Ahmed, Z.,  
 Ahmed, Z., RosnBrick, Z., Shu, Z., Rohalska, Z., Wen,  
 Z., Liu, Z., Liu, Z., Qiao, Z., Xu, Z., Zhou, Z., Chen,  
 Z., Tang, Z., Wu, Z., Ouyang, Z., Lei, Z., Hong, Z., Xiu,  
 Z., Zhao, Z., Meng, Z., Jin, Z., Zeng, Z., Liu, Z., Meng,  
 Z., Qiao, Z., Zheng, Z., Qi, Z., Luo, Z., Birkhead, Z. F.,  
 Sun, Z., and Achdut, Z. The llama 4 herd: Architecture,  
 training, evaluation, and deployment notes, 2026. URL  
<https://arxiv.org/abs/2601.11659>.
- Aljaafari, N., Carvalho, D. S., and Freitas, A. Trace for tracking  
 the emergence of semantic representations in trans-  
 formers, 2025. URL <https://arxiv.org/abs/2505.17998>.
- Bhattamishra, S., Patel, A., and Goyal, N. On the ability and  
 limitations of transformers to recognize formal languages.  
*arXiv preprint arXiv:2009.11264*, 2020.
- Broder, A. On the resemblance and containment of docu-  
 ments. In *Proceedings. Compression and Complexity of*  
*SEQUENCES 1997 (Cat. No.97TB100171)*, pp. 21–29,  
 1997. doi: 10.1109/SEQUEN.1997.666900.
- Cagnetta, F. and Wyart, M. Towards a theory of how the  
 structure of language is acquired by deep neural  
 networks. In Globerson, A., Mackey, L., Belgrave,  
 D., Fan, A., Paquet, U., Tomczak, J., and Zhang, C.  
 (eds.), *Advances in Neural Information Processing*  
*Systems*, volume 37, pp. 83119–83163. Curran As-  
 sociates, Inc., 2024. doi: 10.52202/079017-2645.  
 URL [https://proceedings.neurips.cc/paper\\_files/paper/2024/file/9740dalc07c7b451af14e11523f95271-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/9740dalc07c7b451af14e11523f95271-Paper-Conference.pdf).
- Cagnetta, F., Cornacchia, F., and Wyart, M. Towards a  
 theory of how the structure of language is acquired by  
 deep neural networks. *arXiv preprint arXiv:2406.00048*,  
 2024a.
- Cagnetta, F., Petrini, L., Tomasini, U. M., Favero, A., and  
 Wyart, M. How deep neural networks learn compositional  
 data: The random hierarchy model. *Physical Review X*, 14  
 (3), July 2024b. ISSN 2160-3308. doi: 10.1103/physrevx.  
 14.031001. URL <http://dx.doi.org/10.1103/PhysRevX.14.031001>.
- Cagnetta, F., Kang, H., and Wyart, M. Learning curves  
 theory for hierarchically compositional data with power-  
 law distributed features, 2025. URL <https://arxiv.org/abs/2505.07067>.
- Carlini, N., Tramer, F., Wallace, E., Jagielski, M., Herbert-  
 Voss, A., Lee, K., Roberts, A., Brown, T., Song, D.,  
 Erlingsson, U., Oprea, A., and Raffel, C. Extracting  
 training data from large language models, 2021. URL  
<https://arxiv.org/abs/2012.07805>.

- 660 Chen, H., Waheed, A., Li, X., Wang, Y., Wang, J., Raj, B.,  
661 and Abdin, M. I. On the diversity of synthetic data and its  
662 impact on training large language models, 2024a. URL  
663 <https://arxiv.org/abs/2410.15226>.
- 664 Chen, H., Yang, X., Zhu, J., and Wang, W. Quantifying  
665 semantic emergence in language models, 2024b. URL  
666 <https://arxiv.org/abs/2405.12617>.
- 668 Comanici, G., Bieber, E., Schaekermann, M., Pasupat, I.,  
669 Sachdeva, N., Dhillon, I., Blistein, M., Ram, O., Zhang,  
670 D., Rosen, E., Marris, L., Petulla, S., Gaffney, C., Aha-  
671 roni, A., Lintz, N., Pais, T. C., Jacobsson, H., Szpektor,  
672 I., Jiang, N.-J., Haridasan, K., Omran, A., Saunshi, N.,  
673 Bahri, D., Mishra, G., Chu, E., Boyd, T., Hekman, B.,  
674 Parisi, A., Zhang, C., Kawintiranon, K., Bedrax-Weiss,  
675 T., Wang, O., Xu, Y., Purkiss, O., Mendlovic, U., Deu-  
676 tel, I., Nguyen, N., Langley, A., Korn, F., Rossazza, L.,  
677 Ramé, A., Waghmare, S., Miller, H., Byrd, N., Sheshan,  
678 A., Hadsell, R., Bhardwaj, S., Janus, P., Rissa, T., Horgan,  
679 D., Abdagic, A., Belenki, L., Allingham, J., Singh, A.,  
680 Guidroz, T., Srinivasan, S., Schmit, H., Chiafullo, K.,  
681 Elisseeff, A., Jha, N., Kolhar, P., Berrada, L., Ding, F., Si,  
682 X., Mallick, S. B., Och, F., Erell, S., Ni, E., Latkar, T.,  
683 Yang, S., Sirkovic, P., Feng, Z., Leland, R., Hornung, R.,  
684 Wu, G., Blundell, C., Alvari, H., Huang, P.-S., Yip, C.,  
685 Deur, S., Liu, L., Surita, G., Duque, P., Damen, D., Jia, J.,  
686 Guez, A., Mircea, M., Sinha, A., Magni, A., Stradomski,  
687 P., Marian, T., Galić, V., Chen, W., Husain, H., Sing-  
688 hal, A., Grewe, D., Aubet, F.-X., Song, S., Blanco, L.,  
689 Rechis, L., Ho, L., Munoz, R., Zheng, K., Hamrick, J.,  
690 Mather, K., Taitelbaum, H., Rutherford, E., Lei, Y., Chen,  
691 K., Shukla, A., Moreira, E., Doi, E., Isik, B., Shabat,  
692 N., Rogozińska, D., Kolipaka, K., Chang, J., Vušak, E.,  
693 Venkatachary, S., Noghabi, S., Bharti, T., Jun, Y., Zaks,  
694 A., Green, S., Challagundla, J., Wong, W., Mohammad,  
695 M., Hirsch, D., Cheng, Y., Naim, I., Proleev, L., Vincent,  
696 D., Singh, A., Krikun, M., Krishnan, D., Ghahramani, Z.,  
697 Atias, A., Aggarwal, R., Kirov, C., Vytiniotis, D., Koh,  
698 C., Chronopoulou, A., Dogra, P., Ion, V.-D., Tyen, G.,  
699 Lee, J., Weissenberger, F., Strohmman, T., Balakrishna, A.,  
700 Rae, J., Velic, M., de Liedekerke, R., Elyada, O., Yuan,  
701 W., Liu, C., Shani, L., Kishchenko, S., Alessio, B., Li, Y.,  
702 Song, R., Kwei, S., Jankowski, O., Pappu, A., Namiki, Y.,  
703 Ma, Y., Tripuraneni, N., Cherry, C., Ikonomidis, M., Ling,  
704 Y.-C., Ji, C., Westberg, B., Wright, A., Yu, D., Parkin-  
705 son, D., Ramaswamy, S., Connor, J., Yeganeh, S. H.,  
706 Grover, S., Kenwright, G., Litchev, L., Apps, C., Tomala,  
707 A., Halim, F., Castro-Ros, A., Li, Z., Boral, A., Sho, P.,  
708 Yarom, M., Malmi, E., Klinghoffer, D., Lin, R., Ansell,  
709 A., S. P. K., Zhao, S., Zuo, S., Santoro, A., Cheng, H.-T.,  
710 Demmessie, S., Liu, Y., Brichtova, N., Culp, A., Braun,  
711 N., Graur, D., Ng, W., Mehta, N., Phillips, A., Sundberg,  
712 P., Godbole, V., Liu, F., Katariya, Y., Rim, D., Seyed-  
713 hosseini, M., Ammirati, S., Valfridsson, J., Malihi, M.,  
714 Knight, T., Toor, A., Lampe, T., Ittycheriah, A., Chiang,  
L., Yeung, C., Fréchette, A., Rao, J., Wang, H., Srivas-  
tava, H., Zhang, R., Rhodes, R., Brand, A., Weesner, D.,  
Figotin, I., Gimeno, F., Fellingner, R., Marcenac, P., Leal,  
J., Marcus, E., Cotruta, V., Cabrera, R., Luo, S., Gar-  
rette, D., Axelrod, V., Baltateanu, S., Barker, D., Chen,  
D., Toma, H., Ingram, B., Riesa, J., Kulkarni, C., Zhang,  
Y., Liu, H., Wang, C., Polacek, M., Wu, W., Hui, K.,  
Reyes, A. N., Su, Y., Barnes, M., Malhi, I., Siddiqui, A.,  
Feng, Q., Damaschin, M., Pighin, D., Steiner, A., Yang,  
S., Boppana, R. S., Ivanov, S., Kandoor, A., Shah, A.,  
Mujika, A., Huang, D., Choquette-Choo, C. A., Patel,  
M., Yu, T., Creswell, T., Jerry, Liu, Barros, C., Razeghi,  
Y., Roy, A., Culliton, P., Xiong, B., Pan, J., Strohmman,  
T., Powell, T., Seal, B., DeCarlo, D., Shyam, P., Katir-  
cioglu, K., Wang, X., Hardin, C., Odisho, I., Broder, J.,  
Chang, O., Nair, A., Shtefan, A., O’Brien, M., Agarwal,  
M., Potluri, S., Goyal, S., Jhinal, A., Thakur, S., Stuken,  
Y., Lyon, J., Toutanova, K., Feng, F., Wu, A., Horn, B.,  
Wang, A., Cullum, A., Taubman, G., Shrivastava, D.,  
Shi, C., Tomlinson, H., Patel, R., Tu, T., Oflazer, A. M.,  
Pongetti, F., Yang, M., Taïga, A. A., Perot, V., Pierse,  
N. W., Han, F., Drori, Y., Iturrate, I., Chakrabarti, A.,  
Yeung, L., Dopson, D., ting Chen, Y., Kulshreshtha, A.,  
Guo, T., Pham, P., Schuster, T., Chen, J., Polozov, A.,  
Xing, J., Zhou, H., Kacham, P., Kukliansky, D., Miech,  
A., Yaroshenko, S., Chi, E., Douglas, S., Fei, H., Blondel,  
M., Myla, P., Madmoni, L., Wu, X., Keysers, D., Kjemis,  
K., Albuquerque, I., Yu, L., D’sa, J., Plantan, M., Ionescu,  
V., Elias, J. S., Gupta, A., Vuyyuru, M. R., Alcober, F.,  
Zhou, T., Ji, K., Hartmann, F., Puttagunta, S., Song, H.,  
Amid, E., Stefanoiu, A., Lee, A., Pucciarelli, P., Wang, E.,  
Raul, A., Petrov, S., Tian, I., Anklin, V., Nti, N., Gomes,  
V., Schumacher, M., Vesom, G., Panagopoulos, A., Bous-  
malis, K., Andor, D., Jacob, J., Zhang, Y., Rosgen, B.,  
Kecman, M., Tung, M., Belias, A., Goodman, N., Cov-  
ington, P., Wieder, B., Saxena, N., Davoodi, E., Huang,  
M., Maddineni, S., Roulet, V., Campbell-Ajala, F., Sessa,  
P. G., Xintian, Wu, Lai, G., Collins, P., Haig, A., Sakenas,  
V., Xu, X., Giustina, M., Shafey, L. E., Charoenpanit, P.,  
Garg, S., Ainslie, J., Severson, B., Arenas, M. G., Pathak,  
S., Rajayogam, S., Feng, J., Bakker, M., Li, S., Wich-  
ers, N., Rogers, J., Geng, X., Li, Y., Jagerman, R., Jia,  
C., Olmert, N., Sharon, D., Mauger, M., Mariserla, S.,  
Ma, H., Mohabey, M., Kim, K., Andreev, A., Pollom, S.,  
Love, J., Jain, V., Agrawal, P., Schroecker, Y., Fortin, A.,  
Warmuth, M., Liu, J., Leach, A., Blok, I., Girirajan, G. P.,  
Aharoni, R., Uria, B., Sozanschi, A., Goldberg, D., Ionita,  
L., Ribeiro, M. T., Zlocha, M., Birodkar, V., Lachgar, S.,  
Yuan, L., Choudhury, H., Ginsberg, M., Zheng, F., Dibb,  
G., Graves, E., Lokhande, S., Rasskin, G., Muraru, G.-  
C., Quick, C., Tata, S., Sermanet, P., Chawla, A., Karo,  
I., Wang, Y., Zhang, S., Keller, O., Dragan, A., Su, G.,  
Chou, I., Liu, X., Tao, Y., Prabhakara, S., Wilson, M.,

715 Liu, R., Wang, S., Evans, G., Du, D., Castaño, A., Prasad,  
716 G., Mahdy, M. E., Gerlach, S., Reid, M., Kahn, J., Zait,  
717 A., Pillai, T. S., Ulrich, T., Wang, G., Wassenberg, J.,  
718 Farkash, E., Yalasang, K., Wang, C., Bauza, M., Bucher,  
719 S., Liu, T., Yan, J., Leung, G., Sindhvani, V., Barnes, P.,  
720 Singh, A., Jurin, I., Chang, J., Bhumihar, N. K., Eiger, S.,  
721 Citovsky, G., Withbroe, B., Li, Z., Xue, S., Santo, N. D.,  
722 Stoyanov, G., Raimond, Y., Zheng, S., Gao, Y., Listik,  
723 V., Kwasiborski, S., Saputro, R., Ozturk, A., Mallya, G.,  
724 Majmundar, K., West, R., Caron, P., Wei, J., Castrejon,  
725 L., Vikram, S., Ramachandran, D., Dhawan, N., Park,  
726 J., Smoot, S., van den Driessche, G., Blau, Y., Malik,  
727 C., Liang, W., Hirsch, R., dos Santos, C. N., Weinstein,  
728 E., van den Oord, A., Lall, S., FitzGerald, N., Jiang, Z.,  
729 Yang, X., Webster, D., Elqursh, A., Pope, A., Rotival,  
730 G., Raposo, D., Zhu, W., Dean, J., Alabed, S., Tran, D.,  
731 Gupta, A., Gleicher, Z., Austin, J., Rosseel, E., Umekar,  
732 M., Das, D., Sun, Y., Chen, K., Misiunas, K., Zhou, X.,  
733 Di, Y., Loo, A., Newlan, J., Li, B., Ramasesh, V., Xu,  
734 Y., Chen, A., Gandhe, S., Soricut, R., Gupta, N., Hu,  
735 S., El-Sayed, S., Garcia, X., Brusilovsky, I., Chen, P.-C.,  
736 Bolt, A., Huang, L., Gurney, A., Zhang, Z., Pritzel, A.,  
737 Wilkiewicz, J., Seybold, B., Shamanna, B. K., Fischer,  
738 F., Dean, J., Gill, K., McIlroy, R., Bhowmick, A., Sel-  
739 lier, J., Yang, A., Cheng, D., Magay, V., Tan, J., Varma,  
740 D., Walder, C., Kocisky, T., Nakashima, R., Natsev, P.,  
741 Kwong, M., Gog, I., Zhang, C., Dieleman, S., Jimma,  
742 T., Ryabtsev, A., Brahma, S., Steiner, D., Du, D., Žužul,  
743 A., Žanić, M., Raghavachari, M., Gierke, W., Zheng, Z.,  
744 Petrova, D., Dauphin, Y., Liu, Y., Kessler, I., Hand, S.,  
745 Duvarney, C., Kim, S., Lee, H., Hussenot, L., Hui, J.,  
746 Smith, J., Jain, D., Xia, J., Tomar, G. S., Amiri, K., Phan,  
747 D., Fuchs, F., Weyand, T., Tomasev, N., Cordell, A., Liu,  
748 X., Mallinson, J., Joshi, P., Crawford, A., Suggala, A.,  
749 Chien, S., Fernando, N., Sanchez-Vargas, M., Williams,  
750 D., Crone, P., Luo, X., Karpov, I., Shan, J., Thurk, T.,  
751 Strudel, R., Voigtlaender, P., Patil, P., Dozat, T., Kho-  
752 daei, A., Singla, S., Ambroszczyk, P., Wu, Q., Chang, Y.,  
753 Roark, B., Hegde, C., Ding, T., Filos, A., Wu, Z., Pinto,  
754 A. S., Liu, S., Khanna, S., Pandey, A., Mcloughlin, S.,  
755 Li, Q., Haves, S., Zhou, A., Buchatskaya, E., Leal, I.,  
756 de Boursac, P., Akazawa, N., Anderson, N., Chen, T.,  
757 Somandepalli, K., Liang, C., Goenka, S., Winkler, S.,  
758 Grushetsky, A., Ding, Y., Smith, J., Ye, F., Pont-Tuset, J.,  
759 Li, E., Li, R., Golany, T., Wegner, D., Jiang, T., Barak,  
760 O., Shangguan, Y., Vértés, E., Wong, R., Bornschein,  
761 J., Tudor, A., Bevilacqua, M., Schaul, T., Rawat, A. S.,  
762 Zhao, Y., Axiotis, K., Meng, L., McLean, C., Lai, J.,  
763 Beattie, J., Kushman, N., Liu, Y., Kutzman, B., Lang,  
764 F., Ye, J., Netrapalli, P., Mishra, P., Khan, M., Goel, M.,  
765 Willoughby, R., Tian, D., Zhuang, H., Chen, J., Tsai,  
766 Z., Kementsietsidis, T., Khare, A., Keeling, J., Xu, K.,  
767 Waters, N., Altché, F., Popat, A., Mittal, B., Saxton, D.,  
768 Badawy, D. E., Mathieu, M., Zheng, Z., Zhou, H., Ranka,  
769 N., Shin, R., Duan, Q., Salimans, T., Mihailescu, I., Sha-  
ham, U., Chang, M.-W., Assael, Y., Dikkala, N., Izzard,  
M., Cohen-Addad, V., Graves, C., Feinberg, V., Chung,  
G., Strouse, D., Karmon, D., Sharifzadeh, S., Ashwood,  
Z., Pham, K., Blanton, J., Vasiloff, A., Barber, J., Geller,  
M., Zhou, A., Zubach, F., Huang, T.-K., Zhang, L., Gupta,  
H., Young, M., Proskurnia, J., Votel, R., Gabeur, V., Bar-  
cik, G., Tripathi, A., Yu, H., Yan, G., Changpinyo, B.,  
Pavetić, F., Coyle, A., Fujii, Y., Mendez, J. G., Zhou,  
T., Rajamani, H., Hechtman, B., Cao, E., Juan, D.-C.,  
Tan, Y.-X., Dalibard, V., Du, Y., Clay, N., Yao, K., Jia,  
W., Vijaykumar, D., Zhou, Y., Bai, X., Hung, W.-C.,  
Pecht, S., Todorov, G., Khadke, N., Gupta, P., Lahoti,  
P., Autef, A., Duddu, K., Lee-Thorp, J., Bykovsky, A.,  
Misiunas, T., Flennerhag, S., Thangaraj, S., McGiffin, J.,  
Nado, Z., Kunesch, M., Noever, A., Hertz, A., Liang, M.,  
Stone, V., Palmer, E., Daruki, S., Pramanik, A., Pöder, S.,  
Kyker, A., Khan, M., Sluzhaev, E., Ritter, M., Ruderman,  
A., Zhou, W., Nagpal, C., Vodrahalli, K., Necula, G.,  
Barham, P., Pavlick, E., Hartford, J., Shafran, I., Zhao, L.,  
Mikuła, M., Eccles, T., Shimokawa, H., Garg, K., Vilnis,  
L., Chen, H., Shumailov, I., Lee, K.-H., Abdelhamed, A.,  
Xie, M., Cohen, V., Hlavnova, E., Malkin, D., Sitawarin,  
C., Lottes, J., Coquiot, P., Yu, T., Kumar, S., Zhang, J.,  
Mahendru, A., Ahmed, Z., Martens, J., Chen, T., Boag,  
A., Peng, D., Devin, C., Klimovskiy, A., Phuong, M.,  
Vainstein, D., Xie, J., Ramabhadran, B., Howard, N., Yu,  
X., Goswami, G., Cui, J., Shleifer, S., Pinto, M., Yeh,  
C.-K., Yang, M.-H., Javanmardi, S., Ethier, D., Lee, C.,  
Orbay, J., Kotecha, S., Bromberg, C., Shaw, P., Thorn-  
ton, J., Rosenthal, A. G., Gu, S., Thomas, M., Gemp,  
I., Ayyar, A., Ushio, A., Selvan, A., Wee, J., Liu, C.,  
Majzoubi, M., Yu, W., Abernethy, J., Liechty, T., Pan,  
R., Nguyen, H., Qiong, Hu, Perrin, S., Arora, A., Pitler,  
E., Wang, W., Shivakumar, K., Prost, F., Limonchik, B.,  
Wang, J., Gao, Y., Cour, T., Buch, S., Gui, H., Ivanova,  
M., Neubeck, P., Chan, K., Kim, L., Chen, H., Goyal, N.,  
Chung, D.-W., Liu, L., Su, Y., Petrushkina, A., Shen, J.,  
Joulin, A., Xu, Y., Lin, S. X., Kulizhskaya, Y., Chelba,  
C., Vasudevan, S., Collins, E., Bashlovkina, V., Lu, T.,  
Fritz, D., Park, J., Zhou, Y., Su, C., Tanburn, R., Sushkov,  
M., Rasquinha, M., Li, J., Prendki, J., Li, Y., LV, P.,  
Sharma, S., Fitoussi, H., Huang, H., Dai, A., Dao, P.,  
Burrows, M., Prior, H., Qin, D., Pundak, G., Sjoesund,  
L. L., Khurshudov, A., Zhu, Z., Webson, A., Kemp, E.,  
Tan, T., Agrawal, S., Sargsyan, S., Cheng, L., Stephan, J.,  
Kwiatkowski, T., Reid, D., Byravan, A., Michaely, A. H.,  
Heess, N., Zhou, L., Goenka, S., Carpenter, V., Levskaya,  
A., Wang, B., Roberts, R., Leblond, R., Chikkerur, S.,  
Ginzburg, S., Chang, M., Riachi, R., Chuqiao, Xu, Bor-  
sos, Z., Pliskin, M., Pawar, J., Lustman, M., Kirkwood,  
H., Anand, A., Chaudhary, A., Kalb, N., Milan, K., Au-  
genstein, S., Goldie, A., Prince, L., Raman, K., Sun, Y.,  
Xia, V., Cohen, A., Huo, Z., Camp, J., Ellis, S., Zilka,

770 L., Torres, D. V., Patel, L., Arora, S., Chan, B., Adler, J.,  
771 Ayoub, K., Liang, J., Jamil, F., Jiang, J., Baumgartner, S.,  
772 Sun, H., Karov, Y., Akulov, Y., Zheng, H., Cai, I., Fantacci,  
773 C., Rubin, J., Acha, A. R., Wang, M., D'Souza, N.,  
774 Sathyanarayana, R., Dai, S., Rowe, S., Simanovsky, A.,  
775 Goldman, O., Kuang, Y., Pan, X., Rosenberg, A., Rojas-  
776 Esponda, T., Dutta, P., Zeng, A., Jurenka, I., Farquhar, G.,  
777 Bansal, Y., Iqbal, S., Roelofs, B., Joung, G.-Y., Beak, P.,  
778 Ryu, C., Poplin, R., Wu, Y., Alayrac, J.-B., Buthpitiya,  
779 S., Ronneberger, O., Habtegebriel, C., Li, W., Cavallaro,  
780 P., Wei, A., Bensch, G., Denk, T., Ganapathy, H., Stan-  
781 way, J., Joshi, P., Bertolini, F., Lo, J., Ma, O., Charles,  
782 Z., Sampemane, G., Sahni, H., Chen, X., Askham, H.,  
783 Gaddy, D., Young, P., Tan, J., Eyal, M., Bražinskis, A.,  
784 Zhong, L., Wu, Z., Epstein, M., Bailey, K., Hard, A., Lee,  
785 K., Goldshtein, S., Ruiz, A., Badawi, M., Lochbrunner,  
786 M., Kearns, J., Brown, A., Pardo, F., Weber, T., Yang, H.,  
787 Jiang, P.-P., Akin, B., Fu, Z., Wainwright, M., Zou, C.,  
788 Gaba, M., Manzagol, P.-A., Kan, W., Song, Y., Zainullina,  
789 K., Lin, R., Ko, J., Deshmukh, S., Jindal, A., Svensson, J.,  
790 Tyam, D., Zhao, H., Kaeser-Chen, C., Baird, S., Moradi,  
791 P., Hall, J., Guo, Q., Tsang, V., Liang, B., Pereira, F.,  
792 Ganesh, S., Korotkov, I., Adamek, J., Thiagarajan, S.,  
793 Tran, V., Chen, C., Tar, C., Jain, S., Dasgupta, I., Bilal, T.,  
794 Reitter, D., Zhao, K., Vezzani, G., Gehman, Y., Mehta,  
795 P., Beltrone, L., Dotiwalla, X., Guadarrama, S., Abbas,  
796 Z., Karp, S., Georgiev, P., Ferng, C.-S., Brockschmidt,  
797 M., Peng, L., Hirnschall, C., Verma, V., Bi, Y., Xiao, Y.,  
798 Dabush, A., Xu, K., Wallis, P., Parker, R., Wang, Q., Xu,  
799 Y., Safarli, I., Tewari, D., Zhang, Y., Kim, S., Gesmundo,  
800 A., Thomas, M., Levi, S., Chowdhury, A., Rao, K., Garst,  
801 P., Conway-Rahman, S., Ran, H., McKinney, K., Xiao, Z.,  
802 Yu, W., Agrawal, R., Stjerngren, A., Ionescu, C., Chen, J.,  
803 Sharma, V., Chiu, J., Liu, F., Franko, K., Sanford, C., Cai,  
804 X., Michel, P., Ganapathy, S., Labanowski, J., Garrett, Z.,  
805 Vargas, B., Sun, S., Gale, B., Buschmann, T., Desjardins,  
806 G., Ghelani, N., Jain, P., Verma, M., Asawaroengchai,  
807 C., Eisenschlos, J., Harlalka, J., Kazawa, H., Metzler, D.,  
808 Howland, J., Jian, Y., Ades, J., Shah, V., Gangwani, T.,  
809 Lee, S., Ring, R., Hernandez, S. M., Reich, D., Sinha,  
810 A., Sathe, A., Kovac, J., Gill, A., Kannan, A., D'olimpio,  
811 A., Sevenich, M., Whang, J., Kim, B., Sim, K. C., Chen,  
812 J., Zhang, J., Lall, S., Matias, Y., Jia, B., Friesen, A.,  
813 Nasso, S., Thapliyal, A., Perozzi, B., Yu, T., Shekhawat,  
814 A., Huda, S., Grabowski, P., Wang, E., Sreevatsa, A.,  
815 Dib, H., Hassen, M., Schuh, P., Milutinovic, V., Welty,  
816 C., Quinn, M., Shah, A., Wang, B., Barth-Maron, G.,  
817 Frye, J., Axelsson, N., Zhu, T., Ma, Y., Giannoumis, I.,  
818 Sedghi, H., Ye, C., Luan, Y., Aydin, K., Chandra, B.,  
819 Sampathkumar, V., Huang, R., Lavrenko, V., Eleryan,  
820 A., Hong, Z., Hansen, S., Carthy, S. M., Samanta, B.,  
821 Čevič, D., Wang, X., Li, F., Voznesensky, M., Hoffman,  
822 M., Terzis, A., Schwag, V., Fidel, G., He, L., Cai, M., He,  
823 Y., Feng, A., Nikoltchev, M., Phatale, S., Chase, J., Law-  
824 ton, R., Zhang, M., Ouyang, T., Tragut, M., Manshadi,  
M. H., Narayanan, A., Shen, J., Gao, X., Bolukbasi, T.,  
Roy, N., Li, X., Golovin, D., Panait, L., Qin, Z., Han,  
G., Anthony, T., Kudugunta, S., Patraucean, V., Ray, A.,  
Chen, X., Yang, X., Bhatia, T., Talluri, P., Morris, A.,  
Ražnatović, A., Brownfield, B., An, J., Peng, S., Kane,  
P., Zheng, C., Duduta, N., Kessinger, J., Noraky, J., Liu,  
S., Rong, K., Veličković, P., Rush, K., Goldin, A., Wei,  
F., Garlapati, S. M. R., Pantofaru, C., Kwon, O., Ni, J.,  
Noland, E., Trapani, J. D., Beaufays, F., Roy, A. G., Chow,  
Y., Turker, A., Cideron, G., Mei, L., Clark, J., Dou, Q.,  
Bošnjak, M., Leith, R., Du, Y., Yazdanbakhsh, A., Nasr,  
M., Kwak, C., Sheth, S. S., Kaskasoli, A., Anand, A.,  
Lakshminarayanan, B., Jerome, S., Bieber, D., Chu, C.-  
T., Senges, A., Shen, T., Sridhar, M., Ndebele, N., Beyret,  
B., Mohamed, S., Chen, M., Freitag, M., Guo, J., Liu,  
L., Roit, P., Chen, H., Yan, S., Stone, T., Co-Reyes, J.,  
Cole, J., Scellato, S., Azizi, S., Hashemi, H., Jin, A., Iyer,  
A., Valentine, M., György, A., Ahuja, A., Diaz, D. H.,  
Lee, C.-Y., Clement, N., Kong, W., Garmon, D., Watts, I.,  
Bhatia, K., Gupta, K., Miecznikowski, M., Vallet, H., Taly,  
A., Loper, E., Joshi, S., Atwood, J., Chick, J., Collier, M.,  
Iliopoulos, F., Trostle, R., Gunel, B., Leal-Cavazos, R.,  
Hrafnkelsson, A. M., Guzman, M., Ju, X., Forbes, A.,  
Emond, J., Chauhan, K., Caine, B., Xiao, L., Zeng, W.,  
Moufarek, A., Murphy, D., Meng, M., Gupta, N., Riedel,  
F., Das, A., Lawal, E., Narayan, S., Sosea, T., Swirhun, J.,  
Friso, L., Neyshabur, B., Lu, J., Girgin, S., Wunder, M.,  
Yvinec, E., Pyne, A., Carbune, V., Rijhwani, S., Guo, Y.,  
Doshi, T., Briukhov, A., Bain, M., Hitron, A., Wang, X.,  
Gupta, A., Chen, K., Du, C., Zhang, W., Shah, D., Akula,  
A., Dylla, M., Kachra, A., Kuo, W., Zou, T., Wang, L., Xu,  
L., Zhu, J., Snyder, J., Menon, S., Firat, O., Mordatch, I.,  
Yuan, Y., Ponomareva, N., Blevins, R., Moore, L., Wang,  
W., Chen, P., Scholz, M., Dwornik, A., Lin, J., Li, S., An-  
tognini, D., I. T., Song, X., Miller, M., Kalra, U., Raveret,  
A., Akerlund, O., Wu, F., Nystrom, A., Godbole, N., Liu,  
T., DeBalsi, H., Zhao, J., Liu, B., Caciularu, A., Lax, L.,  
Khandelwal, U., Langston, V., Bailey, E., Lattanzi, S.,  
Wang, Y., Kovelamudi, N., Mondal, S., Guruganesh, G.,  
Hua, N., Roval, O., Wesołowski, P., Ingale, R., Halcrow,  
J., Sohn, T., Angermueller, C., Raad, B., Stickgold, E.,  
Lu, E., Kosik, A., Xie, J., Lillicrap, T., Huang, A., Zhang,  
L. L., Paulus, D., Farabet, C., Wertheim, A., Wang, B.,  
Joshi, R., Ling Ko, C., Wu, Y., Agrawal, S., Lin, L., Sheng,  
X., Sung, P., Breland-King, T., Butterfield, C., Gawde,  
S., Singh, S., Zhang, Q., Apte, R., Shetty, S., Hutter, A.,  
Li, T., Salesky, E., Lebron, F., Kanerva, J., Paganini, M.,  
Nguyen, A., Vallu, R., Peter, J.-T., Velury, S., Kao, D.,  
Hoover, J., Bortsova, A., Bishop, C., Jakobovits, S., Agos-  
tini, A., Agarwal, A., Liu, C., Kwong, C., Tavakkol, S.,  
Bica, I., Greve, A., GP, A., Marcus, J., Hou, L., Duerig,  
T., Moroshko, R., Lacey, D., Davis, A., Amelot, J., Wang,  
G., Kim, F., Strinopoulos, T., Wan, H., Lan, C. L., Krish-

- 825 nan, S., Tang, H., Humphreys, P., Bai, J., Shtacher, I. H.,  
 826 Machado, D., Pang, C., Burke, K., Liu, D., Aravamud-  
 827 han, R., Song, Y., Hirst, E., Singh, A., Jou, B., Bai, L.,  
 828 Piccinno, F., Fu, C. K., Alazard, R., Meiri, B., Winter, D.,  
 829 Chen, C., Zhang, M., Heitkaemper, J., Lambert, J., Lee,  
 830 J., Frömmgen, A., Rogulenko, S., Nair, P., Niemczyk, P.,  
 831 Bulyenov, A., Xu, B., Shemtov, H., Zadimoghaddam, M.,  
 832 Toropov, S., Wirth, M., Dai, H., Gollapudi, S., Zheng, D.,  
 833 Kurakin, A., Lee, C., Bullard, K., Serrano, N., Balazevic,  
 834 I., Li, Y., Schalkwyk, J., Murphy, M., Zhang, M., Se-  
 835 queira, K., Datta, R., Agrawal, N., Sutton, C., Attaluri, N.,  
 836 Chiang, M., Farhan, W., Thornton, G., Lin, K., Choma,  
 837 T., Nguyen, H., Dasgupta, K., Robinson, D., Comşa, I.,  
 838 Riley, M., Pillai, A., Mustafa, B., Golan, B., Zandieh, A.,  
 839 Lespiau, J.-B., Porter, B., Ross, D., Rajayogam, S., Agar-  
 840 wal, M., Venugopalan, S., Shahriari, B., Yan, Q., Xu, H.,  
 841 Tobin, T., Dubov, P., Shi, H., Recasens, A., Kovsharov,  
 842 A., Borgeaud, S., Dery, L., Vasanth, S., Gribovskaya, E.,  
 843 Qiu, L., Mahdieh, M., Skut, W., Nielsen, E., Zheng, C.,  
 844 Yu, A., Bostock, C. G., Gupta, S., Archer, A., Rawles,  
 845 C., Davies, E., Svyatkovskiy, A., Tsai, T., Halpern, Y.,  
 846 Reisswig, C., Wydrowski, B., Chang, B., Puigcerver, J.,  
 847 Taege, M. H., Li, J., Schnider, E., Li, X., Dena, D., Xu,  
 848 Y., Telang, U., Shi, T., Zen, H., Kastner, K., Ko, Y., Sub-  
 849 ramaniam, N., Kumar, A., Blois, P., Dai, Z., Wieting,  
 850 J., Lu, Y., Zeldes, Y., Xie, T., Hauth, A., Tifrea, A., Li,  
 851 Y., El-Husseini, S., Abolafia, D., Zhou, H., Ding, W.,  
 852 Ghalebikesabi, S., Guía, C., Maksai, A., Ágoston Weisz,  
 853 Arik, S., Sukhanov, N., Świetlik, A., Jia, X., Yu, L., Wang,  
 854 W., Brand, M., Bloxwich, D., Kirmani, S., Chen, Z., Go,  
 855 A., Sprechmann, P., Kannen, N., Carin, A., Sandhu, P.,  
 856 Edkins, I., Nooteboom, L., Gupta, J., Maggiore, L., Azizi,  
 857 J., Pritch, Y., Yin, P., Gupta, M., Tarlow, D., Smith, D.,  
 858 Ivanov, D., Babaeizadeh, M., Goel, A., Kambala, S., Chu,  
 859 G., Kastelic, M., Liu, M., Soltau, H., Stone, A., Agrawal,  
 860 S., Kim, M., Soparkar, K., Tadepalli, S., Bunyan, O.,  
 861 Soh, R., Kannan, A., Kim, D., Chen, B. J., Halumi, A.,  
 862 Roy, S., Wang, Y., Sercinoglu, O., Gibson, G., Bhatnagar,  
 863 S., Sano, M., von Dincklage, D., Ren, Q., Mitrevski, B.,  
 864 Olšák, M., She, J., Doersch, C., Jilei, Wang, Liu, B., Tan,  
 865 Q., Yakar, T., Warkentin, T., Ramirez, A., Lebsack, C.,  
 866 Dillon, J., Mathews, R., Cogley, T., Wu, Z., Chen, Z.,  
 867 Simon, J., Nath, S., Sainath, T., Bendebury, A., Julian, R.,  
 868 Mankalale, B., Čurko, D., Zacchello, P., Brown, A. R.,  
 869 Sodhia, K., Howard, H., Caelles, S., Gupta, A., Evans, G.,  
 870 Bulanova, A., Katzen, L., Goldenberg, R., Tsitsulin, A.,  
 871 Stanton, J., Schillings, B., Kovalev, V., Fry, C., Shah, R.,  
 872 Lin, K., Upadhyay, S., Li, C., Radpour, S., Maggioni, M.,  
 873 Xiong, J., Haas, L., Brennan, J., Kamath, A., Savinov, N.,  
 874 Nagrani, A., Yacovone, T., Kappedal, R., Andriopoulos,  
 875 K., Lao, L., Li, Y., Rozhdestvenskiy, G., Hashimoto, K.,  
 876 Audibert, A., Austin, S., Rodriguez, D., Ruoss, A., Honke,  
 877 G., Karkhanis, D., Xiong, X., Wei, Q., Huang, J., Leng,  
 878 Z., Premachandran, V., Bileschi, S., Evangelopoulos, G.,  
 879 Mensink, T., Pavagadhi, J., Teplyashin, D., Chang, P.,  
 Xue, L., Tanzer, G., Goldman, S., Patel, K., Li, S., Wies-  
 ner, J., Zheng, I., Stewart-Binks, I., Han, J., Li, Z., Luo,  
 L., Lenc, K., Lučić, M., Xue, F., Mullins, R., Guseynov,  
 A., Chang, C.-C., Galatzer-Levy, I., Zhang, A., Bingham,  
 G., Hu, G., Hartman, A., Ma, Y., Griffith, J., Irpan, A.,  
 Radebaugh, C., Yue, S., Fan, L., Ungureanu, V., Sorokin,  
 C., Teufel, H., Li, P., Anil, R., Paparas, D., Wang, T.,  
 Lin, C.-C., Peng, H., Shum, M., Petrovic, G., Brady, D.,  
 Nguyen, R., Macherey, K., Li, Z., Singh, H., Yenugula,  
 M., Iinuma, M., Chen, X., Kopparapu, K., Stern, A.,  
 Dave, S., Thekkath, C., Perot, F., Kumar, A., Li, F., Xiao,  
 Y., Bilotti, M., Bateni, M. H., Noble, I., Lee, L., Vázquez-  
 Reina, A., Salazar, J., Yang, X., Wang, B., Gruzewska,  
 E., Rao, A., Raghuram, S., Xu, Z., Ben-David, E., Mei,  
 J., Dalmia, S., Zhang, Z., Liu, Y., Bansal, G., Pankov, H.,  
 Schwarcz, S., Burns, A., Chan, C., Sanghai, S., Liang,  
 R., Liang, E., He, A., Stuart, A., Narayanan, A., Zhu,  
 Y., Frank, C., Fatemi, B., Sabne, A., Lang, O., Bhat-  
 tacharya, I., Settle, S., Wang, M., McMahan, B., Tac-  
 chetti, A., Soares, L. B., Hadian, M., Cabi, S., Chung, T.,  
 Putikhin, N., Li, G., Chen, J., Tarango, A., Michalewski,  
 H., Kazemi, M., Masoom, H., Sheftel, H., Shivanna, R.,  
 Vadali, A., Comanescu, R., Reid, D., Moore, J., Neelakan-  
 tan, A., Sander, M., Herzig, J., Rosenberg, A., Dehghani,  
 M., Choi, J., Fink, M., Hayes, R., Ge, E., Weng, S., Ho,  
 C.-H., Karro, J., Krishna, K., Thiet, L. N., Skerry-Ryan,  
 A., Eppens, D., Andreetto, M., Sarma, N., Bonacina,  
 S., Ayan, B. K., Nawhal, M., Shan, Z., Dusenberry, M.,  
 Thakoor, S., Gubbi, S., Nguyen, D. D., Tsarfaty, R., Al-  
 banie, S., Mitrović, J., Gandhi, M., Chen, B.-J., Epasto,  
 A., Stephanov, G., Jin, Y., Gehman, S., Amini, A., Weber,  
 J., Behbahani, F., Xu, S., Allamanis, M., Chen, X., Ott,  
 M., Sha, C., Jastrzebski, M., Qi, H., Greene, D., Wu, X.,  
 Toki, A., Vlasic, D., Shapiro, J., Kotikalapudi, R., Shen,  
 Z., Saeki, T., Xie, S., Cassirer, A., Bharadwaj, S., Kiyono,  
 T., Bhojanapalli, S., Rosenfeld, E., Ritter, S., Mao, J.,  
 Oliveira, J. G., Egyed, Z., Bandemer, B., Parisotto, E., Ki-  
 noshita, K., Pluto, J., Maniatis, P., Li, S., Guo, Y., Ghiasi,  
 G., Tarbouriech, J., Chatterjee, S., Jin, J., Katrina, Xu,  
 Palomaki, J., Arnold, S., Sewak, M., Piccinini, F., Sharma,  
 M., Albrecht, B., Purser-haskell, S., Vaswani, A., Chen,  
 C., Wisniewski, M., Cao, Q., Aslanides, J., Phu, N. M.,  
 Sieb, M., Agubuzu, L., Zheng, A., Sohn, D., Selvi, M.,  
 Andreassen, A., Subudhi, K., Eruvbetine, P., Woodman,  
 O., Mery, T., Krause, S., Ren, X., Ma, X., Luo, J., Chen,  
 D., Fan, W., Griffiths, H., Schuler, C., Li, A., Zhang, S.,  
 Sarr, J.-M., Luo, S., Patana, R., Watson, M., Naboulsi,  
 D., Collins, M., Sidhwani, S., Hoogeboom, E., Silver,  
 S., Caveness, E., Zhao, X., Rodriguez, M., Deines, M.,  
 Bai, L., Griffin, P., Tagliasacchi, M., Xue, E., Babbula,  
 S. R., Pang, B., Ding, N., Shen, G., Peake, E., Crocker, R.,  
 Raghvendra, S. S., Swisher, D., Han, W., Singh, R., Wu,  
 L., Pchelin, V., Munkhdalai, T., Alon, D., Bacon, G., Rob-

- 880 les, E., Bulian, J., Johnson, M., Powell, G., Ferreira, F. T.,  
 881 Li, Y., Benzing, F., Velimirović, M., Soyer, H., Kong,  
 882 W., Tony, Nguyễn, Yang, Z., Liu, J., van Amersfoort, J.,  
 883 Gillick, D., Sun, B., Rauschmayr, N., Zhang, K., Zhan, S.,  
 884 Zhou, T., Frolov, A., Yang, C., Vnukov, D., Rouillard, L.,  
 885 Li, H., Mandhane, A., Fallen, N., Venkataraman, R., Hu,  
 886 C. H., Brennan, J., Lee, J., Chang, J., Sundermeyer, M.,  
 887 Pan, Z., Ke, R., Tong, S., Fabrikant, A., Bono, W., Gu, J.,  
 888 Foley, R., Mao, Y., Delakis, M., Bhaswar, D., Frostig, R.,  
 889 Li, N., Zipori, A., Hope, C., Kozlova, O., Mishra, S., Djo-  
 890 longa, J., Schiff, C., Merey, M. A., Briakou, E., Morgan,  
 891 P., Wan, A., Hassidim, A., Skerry-Ryan, R., Sengupta,  
 892 K., Jasarevic, M., Kallakuri, P., Kunkle, P., Brennan, H.,  
 893 Lieber, T., Mansoor, H., Walker, J., Zhang, B., Xie, A.,  
 894 Žužić, G., Chukwuka, A., Druinsky, A., Cho, D., Yao, R.,  
 895 Naeem, F., Butt, S., Kim, E., Jia, Z., Jordan, M., Lelkes,  
 896 A., Kurzeja, M., Wang, S., Zhao, J., Over, A., Chak-  
 897 ladar, A., Prasetya, M., Jha, N., Ganapathy, S., Cong, Y.,  
 898 Shroff, P., Saroufim, C., Miryoosefi, S., Hammad, M.,  
 899 Nasir, T., Xi, W., Gao, Y., Maeng, Y., Hora, B., Cheng,  
 900 C.-Y., Haghani, P., Lewenberg, Y., Lu, C., Matysiak, M.,  
 901 Raisinghani, N., Wang, H., Baugher, L., Sukthankar, R.,  
 902 Giang, M., Schultz, J., Fiedel, N., Chen, M., Lee, C.-C.,  
 903 Dey, T., Zheng, H., Paul, S., Smith, C., Ly, A., Wang,  
 904 Y., Bansal, R., Perz, B., Ricco, S., Blank, S., Keshava,  
 905 V., Sharma, D., Chow, M., Lad, K., Jalan, K., Osindero,  
 906 S., Swanson, C., Scott, J., Ilić, A., Li, X., Jonnalagadda,  
 907 S. R., Soudagar, A. S., Xiong, Y., Batsaikhan, B.-O., Jar-  
 908 rett, D., Kumar, N., Shah, M., Lawlor, M., Waters, A.,  
 909 Graham, M., May, R., Ramos, S., Lefdal, S., Cankara, Z.,  
 910 Cano, N., O'Donoghue, B., Borovik, J., Liu, F., Grimstad,  
 911 J., Alnahlawi, M., Tsihlas, K., Hudson, T., Grigorev, N.,  
 912 Jia, Y., Huang, T., Igwe, T. P., Lebedev, S., Tang, X., Kri-  
 913 vokon, I., Garcia, F., Tan, M., Jia, E., Stys, P., Vashishth,  
 914 S., Liang, Y., Venkatraman, B., Gu, C., Kementsietsidis,  
 915 A., Zhu, C., Jung, J., Bai, Y., Hosseini, M. J., Ahmed, F.,  
 916 Gupta, A., Yuan, X., Ashraf, S., Nigam, S., Vasudevan,  
 917 G., Awasthi, P., Gilady, A. M., Mariet, Z., Eskander, R.,  
 918 Li, H., Hu, H., Garrido, G., Schlattner, P., Zhang, G., Sax-  
 919 ena, R., Dević, P., Muralidharan, K., Murthy, A., Zhou,  
 920 Y., Choi, M., Wongpanich, A., Wang, Z., Shah, P., Xu,  
 921 Y., Huang, Y., Spencer, S., Chen, A., Cohan, J., Wang,  
 922 J., Tompson, J., Wu, J., Haroun, R., Li, H., Huergo, B.,  
 923 Yang, F., Yin, T., Wendt, J., Bendersky, M., Chaabouni,  
 924 R., Snaider, J., Ferret, J., Jindal, A., Thompson, T., Xue,  
 925 A., Bishop, W., Phal, S. M., Sharma, A., Sung, Y., Rad-  
 926 hakrishnan, P., Shomrat, M., Ingle, R., Vij, R., Gilmer, J.,  
 927 Istin, M. D., Sobell, S., Lu, Y., Nottage, E., Sadigh, D.,  
 928 Willcock, J., Zhang, T., Xu, S., Brown, S., Lee, K., Wang,  
 929 G., Zhu, Y., Tay, Y., Kim, C., Gutierrez, A., Sharma, A.,  
 930 Xian, Y., Seo, S., Cui, C., Pochernina, E., Baetu, C., Jas-  
 931 trzebski, K., Ly, M., Elhawaty, M., Suh, D., Sezener, E.,  
 932 Wang, P., Yuen, N., Tucker, G., Cai, J., Yang, Z., Wang,  
 933 C., Muzio, A., Qian, H., Yoo, J., Lockhart, D., McKee,  
 934 K. R., Guo, M., Mehrotra, M., Mendonça, A., Mehta,  
 S. V., Ben, S., Tekur, C., Mu, J., Zhu, M., Krakovna,  
 V., Lee, H., Maschinot, A., Cevey, S., Choe, H., Bai,  
 A., Srinivasan, H., Gasaway, D., Young, N., Siegler, P.,  
 Holtmann-Rice, D., Piratla, V., Baumli, K., Yogev, R.,  
 Hofer, A., van Hasselt, H., Grant, S., Chervonyi, Y., Sil-  
 ver, D., Hogue, A., Agarwal, A., Wang, K., Singh, P.,  
 Flynn, F., Lipschultz, J., David, R., Bellot, L., Yang, Y.-  
 Y., Le, L., Graziano, F., Olszewska, K., Hui, K., Maurya,  
 A., Parotsidis, N., Chen, W., Oguntebi, T., Kelley, J., Bad-  
 depudi, A., Mauerer, J., Shaw, G., Siegman, A., Yang,  
 L., Shetty, S., Roy, S., Song, Y., Stokowiec, W., Bur-  
 nell, R., Savant, O., Busa-Fekete, R., Miao, J., Ghosh, S.,  
 MacDermed, L., Lippe, P., Dektiarev, M., Behrman, Z.,  
 Mentzer, F., Nguyen, K., Wei, M., Verma, S., Knutsen, C.,  
 Dasari, S., Yan, Z., Mitrichev, P., Wang, X., Shejwalkar,  
 V., Austin, J., Sunkara, S., Potti, N., Virin, Y., Wright, C.,  
 Liu, G., Riva, O., Pot, E., Kochanski, G., Le, Q., Balasub-  
 ramaniam, G., Dhar, A., Liao, Y., Bloniarz, A., Shukla,  
 D., Cole, E., Lee, J., Zhang, S., Kafle, S., Vashishtha, S.,  
 Mahmoudieh, P., Chen, G., Hoffmann, R., Srinivasan, P.,  
 Lago, A. D., Shalom, Y. B., Wang, Z., Elabd, M., Sharma,  
 A., Oh, J., Kothawade, S., Le, M., Monteiro, M., Yang,  
 S., Alarakyia, K., Geirhos, R., Mincu, D., Garnes, H.,  
 Kobayashi, H., Mariooryad, S., Krasowiak, K., Zhixin,  
 Lai, Mourad, S., Wang, M., Bu, F., Aharoni, O., Chen,  
 G., Goyal, A., Zubov, V., Bapna, A., Dabir, E., Kothari,  
 N., Lamerigts, K., Cao, N. D., Shar, J., Yew, C., Kulkarni,  
 N., Mahaarachchi, D., Joshi, M., Zhu, Z., Lichtarge, J.,  
 Zhou, Y., Muckenhirn, H., Selo, V., Vinyals, O., Chen,  
 P., Brohan, A., Mehta, V., Cogan, S., Wang, R., Geri, T.,  
 Ko, W.-J., Chen, W., Viola, F., Shivam, K., Wang, L.,  
 Elish, M. C., Popa, R. A., Pereira, S., Liu, J., Koster, R.,  
 Kim, D., Zhang, G., Ebrahimi, S., Talukdar, P., Zheng,  
 Y., Poklukar, P., Mikhalap, A., Johnson, D., Vijayaku-  
 mar, A., Omernick, M., Dibb, M., Dubey, A., Hu, Q.,  
 Suman, A., Aggarwal, V., Kornakov, I., Xia, F., Lowe,  
 W., Kolganov, A., Xiao, T., Nikolaev, V., Hemingray, S.,  
 Li, B., Iljazi, J., Rybiński, M., Sandhu, B., Lu, P., Luong,  
 T., Jenatton, R., Govindaraj, V., Hui, Li, Dulac-Arnold,  
 G., Park, W., Wang, H., Modi, A., Pouget-Abadie, J.,  
 Grellier, K., Gupta, R., Berry, R., Ramachandran, P., Xie,  
 J., McCafferty, L., Wang, J., Gupta, K., Lim, H., Bratanič,  
 B., Brock, A., Akolzin, I., Sproch, J., Karliner, D., Kim,  
 D., Goedeckemeyer, A., Shazeer, N., Schmid, C., Calan-  
 driello, D., Bhatia, P., Choromanski, K., Montgomery, C.,  
 Dua, D., Ramalho, A., King, H., Gao, Y., Nguyen, L.,  
 Lindner, D., Pitta, D., Johnson, O., Salama, K., Ardila,  
 D., Han, M., Farnese, E., Odoom, S., Wang, Z., Ding,  
 X., Rink, N., Smith, R., Lehri, H. T., Cohen, E., Vats, N.,  
 He, T., Gopavarapu, P., Paszke, A., Patel, M., Gansbeke,  
 W. V., Loher, L., Castro, L., Voitovich, M., von Glehn, T.,  
 George, N., Niklaus, S., Eaton-Rosen, Z., Rakićević, N.,  
 Jue, E., Perel, S., Zhang, C., Bahat, Y., Pouget, A., Xing,

- 935 Z., Huot, F., Shenoy, A., Bos, T., Coriou, V., Richter,  
 936 B., Noy, N., Wang, Y., Ontanon, S., Qin, S., Makarchuk,  
 937 G., Hassabis, D., Li, Z., Sharma, M., Venkatesan, K.,  
 938 Kemaev, I., Daniel, R., Huang, S., Shah, S., Ponce, O.,  
 939 Warren, Chen, Faruqui, M., Wu, J., Andačić, S., Payrits,  
 940 S., McDuff, D., Hume, T., Cao, Y., Tessler, M., Wang,  
 941 Q., Wang, Y., Rendulic, I., Agustsson, E., Johnson, M.,  
 942 Lando, T., Howard, A., Padmanabhan, S. G. S., Daswani,  
 943 M., Banino, A., Kilgore, M., Heek, J., Ji, Z., Caceres,  
 944 A., Li, C., Kassner, N., Vlaskin, A., Liu, Z., Grills, A.,  
 945 Hou, Y., Sukkerd, R., Cheon, G., Shetty, N., Markeeva,  
 946 L., Stanczyk, P., Iyer, T., Gong, Y., Gao, S., Gopalakr-  
 947 ishnan, K., Blyth, T., Reynolds, M., Bhoopchand, A.,  
 948 Bilenko, M., Gharibian, D., Zayats, V., Faust, A., Singh,  
 949 A., Ma, M., Jiao, H., Vijayanarasimhan, S., Aroyo, L.,  
 950 Yadav, V., Chakera, S., Kakarla, A., Meshram, V., Gregor,  
 951 K., Botea, G., Senter, E., Jia, D., Kovacs, G., Sharma,  
 952 N., Baur, S., Kang, K., He, Y., Zhuo, L., Kostelac, M.,  
 953 Laish, I., Peng, S., O'Bryan, L., Kasenberg, D., Rao,  
 954 G. R., Leurent, E., Zhang, B., Stevens, S., Salazar, A.,  
 955 Zhang, Y., Lobov, I., Walker, J., Porter, A., Redshaw, M.,  
 956 Ke, H., Rao, A., Lee, A., Lam, H., Moffitt, M., Kim, J.,  
 957 Qiao, S., Koo, T., Dadashi, R., Song, X., Sundararajan,  
 958 M., Xu, P., Kawamoto, C., Zhong, Y., Barbu, C., Reddy,  
 959 A., Verzetti, M., Li, L., Papamakarios, G., Klimczak-  
 960 Plucińska, H., Cassin, M., Kavukcuoglu, K., Swavely, R.,  
 961 Vaucher, A., Zhao, J., Hemsley, R., Tschannen, M., Ge,  
 962 H., Menghani, G., Yu, Y., Ha, N., He, W., Wu, X., Song,  
 963 M., Sterneck, R., Zinke, S., Calian, D. A., Marsden, A.,  
 964 Ruiz, A. C., Hessel, M., Gueta, A., Lee, B., Farris, B.,  
 965 Gupta, M., Li, Y., Saleh, M., Misra, V., Xiao, K., Men-  
 966 dolicchio, P., Buttimore, G., Krayvanova, V., Nayakanti,  
 967 N., Wiethoff, M., Pande, Y., Mirhoseini, A., Lao, N.,  
 968 Liu, J., Hua, Y., Chen, A., Malkov, Y., Kalashnikov, D.,  
 969 Gupta, S., Audhkhasi, K., Zhai, Y., Kopalle, S., Jain, P.,  
 970 Ofek, E., Meyer, C., Baatarsukh, K., Strejček, H., Qian,  
 971 J., Freedman, J., Figueira, R., Sokolik, M., Bachem, O.,  
 972 Lin, R., Kharrat, D., Hidey, C., Xu, P., Duan, D., Li, Y.,  
 973 Ersoy, M., Everett, R., Cen, K., Santamaria-Fernandez,  
 974 R., Taubenfeld, A., Mackinnon, I., Deng, L., Zablot-  
 975 skaia, P., Viswanadha, S., Goel, S., Yates, D., Deng, Y.,  
 976 Choy, P., Chen, M., Sinha, A., Mossin, A., Wang, Y.,  
 977 Szlam, A., Hao, S., Rubenstein, P. K., Toksoz-Exley, M.,  
 978 Aperghis, M., Zhong, Y., Ahn, J., Isard, M., Lacombe,  
 979 O., Luisier, F., Anastasiou, C., Kalley, Y., Prabhu, U.,  
 980 Dunleavy, E., Bijwadia, S., Mao-Jones, J., Chen, K., Pa-  
 981 sumarathi, R., Wood, E., Dostmohamed, A., Hurley, N.,  
 982 Simsa, J., Parrish, A., Pajarskas, M., Harvey, M., Skopek,  
 983 O., Kochinski, Y., Rey, J., Rieser, V., Zhou, D., Lee, S. J.,  
 984 Acharya, T., Li, G., Jiang, J., Zhang, X., Gipson, B.,  
 985 Mahintorabi, E., Gelmi, M., Khajehouri, N., Yeh, A.,  
 986 Lee, K., Matthey, L., Baker, L., Pham, T., Fu, H., Pak,  
 987 A., Gupta, P., Vasconcelos, C., Sadovsky, A., Walker, B.,  
 988 Hsiao, S., Zochbauer, P., Marzoca, A., Velan, N., Zeng,  
 989 J., Baechler, G., Driess, D., Jain, D., Huang, Y., Tao, L.,  
 Maggs, J., Levine, N., Schneider, J., Gemzer, E., Petit, S.,  
 Han, S., Fisher, Z., Zelle, D., Biles, C., Ie, E., Fadeeva,  
 A., Liu, C., Franco, J. V., Collister, A., Zhang, H., Wang,  
 R., Zhao, R., Kieliger, L., Shuster, K., Zhu, R., Gong,  
 B., Chan, L., Sun, R., Basu, S., Zimmermann, R., Hayes,  
 J., Bapna, A., Snoek, J., Yang, W., Datta, P., Abdallah,  
 J. A., Kilgour, K., Li, L., Mah, S., Jun, Y., Rivière, M.,  
 Karmarkar, A., Spalink, T., Huang, T., Gonzalez, L., Tran,  
 D.-H., Nowak, A., Palowitch, J., Chadwick, M., Talius,  
 E., Mehta, H., Sellam, T., Fränken, P., Nicosia, M., He,  
 K., Kini, A., Amos, D., Basu, S., Jobe, H., Shaw, E.,  
 Xu, Q., Evans, C., Ikeda, D., Yan, C., Jin, L., Wang, L.,  
 Yadav, S., Labzovsky, I., Sampath, R., Ma, A., Schu-  
 mann, C., Siddhant, A., Shah, R., Youssef, J., Agarwal,  
 R., Dabney, N., Tonioni, A., Ambar, M., Li, J., Guyon,  
 I., Li, B., Soergel, D., Fang, B., Karadzhov, G., Udrescu,  
 C., Trinh, T., Raunak, V., Noury, S., Guo, D., Gupta, S.,  
 Finkelstein, M., Petek, D., Liang, L., Billock, G., Sun,  
 P., Wood, D., Song, Y., Yu, X., Matejovicova, T., Cohen,  
 R., Andra, K., D'Ambrosio, D., Deng, Z., Nallatamby,  
 V., Songhori, E., Dangovski, R., Lampinen, A., Botadra,  
 P., Hillier, A., Cao, J., Baddi, N., Kuncoro, A., Yoshino,  
 T., Bhagatwala, A., Ranzato, M., Schaeffer, R., Liu, T.,  
 Ye, S., Sarvana, O., Nham, J., Kuang, C., Gao, I., Baek,  
 J., Mittal, S., Wahid, A., Gergely, A., Ni, B., Feldman, J.,  
 Muir, C., Lamblin, P., Macherey, W., Dyer, E., Kilpatrick,  
 L., Campos, V., Bhutani, M., Fort, S., Ahmad, Y., Sev-  
 eryl, A., Chatziprimou, K., Ferludin, O., Dimarco, M.,  
 Kusunpati, A., Heyward, J., Bahir, D., Villela, K., Milli-  
 can, K., Marcus, D., Bahargam, S., Unlu, C., Roth, N.,  
 Wei, Z., Gopal, S., Ghoshal, D., Lee, E., Lin, S., Lees, J.,  
 Lee, D., Hosseini, A., Fan, C., Neel, S., Wu, M., Altun,  
 Y., Cai, H., Piqueras, E., Woodward, J., Bissacco, A.,  
 Haykal, S., Bordbar, M., Sundaram, P., Hodgkinson, S.,  
 Toyama, D., Polovets, G., Myers, A., Sinha, A., Levin-  
 boim, T., Krishnakumar, K., Chhaparia, R., Sholokhova,  
 T., Gundavarapu, N. B., Jawahar, G., Qureshi, H., Hu, J.,  
 Momchev, N., Rahtz, M., Wu, R., S. A. P., Dhamdhare,  
 K., Guo, M., Gupta, U., Eslami, A., Schain, M., Blokzijl,  
 M., Welling, D., Orr, D., Bolelli, L., Perez-Nieves, N.,  
 Sirotenko, M., Prasad, A., Kar, A., Pigem, B. D. B., Terzi,  
 T., Weisz, G., Ghosh, D., Mavalankar, A., Madeka, D.,  
 Daugaard, K., Adam, H., Shah, V., Berman, D., Tran,  
 M., Baker, S., Andrejczuk, E., Chole, G., Raboshchuk,  
 G., Mirzazadeh, M., Kagohara, T., Wu, S., Schallhart,  
 C., Orlando, B., Wang, C., Rustemi, A., Xiong, H., Liu,  
 H., Vezer, A., Ramsden, N., yiin Chang, S., Mudgal, S.,  
 Li, Y., Vieillard, N., Hoshen, Y., Ahmad, F., Slone, A.,  
 Hua, A., Potikha, N., Rossini, M., Stritar, J., Prakash, S.,  
 Wang, Z., Dong, X., Nazari, A., Nehoran, E., Tekelioglu,  
 K., Li, Y., Badola, K., Funkhouser, T., Li, Y., Yerram,  
 V., Ganeshan, R., Formoso, D., Langner, K., Shi, T., Li,  
 H., Yamamori, Y., Panda, A., Saade, A., Scarpatti, A. S.,

- 990 Breaux, C., Carey, C., Zhou, Z., Hsieh, C.-J., Bridgers, S.,  
 991 Butryna, A., Gupta, N., Tulsyan, V., Woo, S., Eltyshev,  
 992 E., Grathwohl, W., Parks, C., Benjamin, S., Panigrahy, R.,  
 993 Dodhia, S., Freitas, D. D., Sauer, C., Song, W., Alet, F.,  
 994 Tolins, J., Paduraru, C., Zhou, X., Albert, B., Zhang, Z.,  
 995 Shu, L., Bansal, M., Nguyen, S., Globerson, A., Xiao, O.,  
 996 Manyika, J., Hennigan, T., Rong, R., Matak, J., Bakalov,  
 997 A., Sharma, A., Sinopalnikov, D., Pierson, A., Roller, S.,  
 998 Brown, G., Gao, M., Fukuzawa, T., Ghafouri, A., Vas-  
 999 sigh, K., Barr, I., Wang, Z., Korsun, A., Jayaram, R.,  
 1000 Ren, L., Zaman, T., Khan, S., Lunts, Y., Deutsch, D.,  
 1001 Uthus, D., Katz, N., Samsikova, M., Khalifa, A., Sethi,  
 1002 N., Sun, J., Tang, L., Alon, U., Luo, X., Yu, D., Nayyar,  
 1003 A., Petrini, B., Truong, W., Hellendoorn, V., Chinaev, N.,  
 1004 Alberti, C., Wang, W., Hu, J., Mirrokni, V., Balashankar,  
 1005 A., Aharon, A., Mehta, A., Iscen, A., Kready, J., Man-  
 1006 ning, L., Mohananeey, A., Chen, Y., Tripathi, A., Wu, A.,  
 1007 Petrovski, I., Hwang, D., Baeuml, M., Chandrakaladha-  
 1008 ran, S., Liu, Y., Coaguila, R., Chen, M., Ma, S., Tafti,  
 1009 P., Tatineni, S., Spitz, T., Ye, J., Vicol, P., Rosca, M.,  
 1010 Puigdomènech, A., Yahav, Z., Ghemawat, S., Lin, H.,  
 1011 Kirk, P., Nabulsi, Z., Brin, S., Bohnet, B., Caluwaerts,  
 1012 K., Veerubhotla, A. S., Zheng, D., Dai, Z., Petrov, P., Xu,  
 1013 Y., Mehran, R., Xu, Z., Zintgraf, L., Choi, J., Hombaiyah,  
 1014 S. A., Thoppilan, R., Reddi, S., Lew, L., Li, L., Webster,  
 1015 K., Sawhney, K., Lamprou, L., Shakeri, S., Lunayach, M.,  
 1016 Chen, J., Bagri, S., Salcianu, A., Chen, Y., Donchev, Y.,  
 1017 Magister, C., Nørly, S., Rodrigues, V., Izo, T., Noga, H.,  
 1018 Zou, J., Köppe, T., Zhou, W., Lee, K., Long, X., Eisen-  
 1019 bud, D., Chen, A., Schenck, C., To, C. M., Zhong, P.,  
 1020 Taropa, E., Truong, M., Levy, O., Martins, D., Zhang,  
 1021 Z., Semturs, C., Zhang, K., Yakubovich, A., Moreno, P.,  
 1022 McConnaughey, L., Lu, D., Redmond, S., Weerts, L.,  
 1023 Bitton, Y., Refice, T., Lacasse, N., Conmy, A., Tallec,  
 1024 C., Odell, J., Forbes-Pollard, H., Socala, A., Hoech, J.,  
 1025 Kohli, P., Walton, A., Wang, R., Sazanovich, M., Zhu,  
 1026 K., Kapishnikov, A., Galt, R., Denton, M., Murdoch, B.,  
 1027 Sikora, C., Mohamed, K., Wei, W., First, U., McConnell,  
 1028 T., Cobo, L. C., Qin, J., Avrahami, T., Balle, D., Watan-  
 1029 abe, Y., Louis, A., Kraft, A., Ariaifar, S., Gu, Y., Rives,  
 1030 E., Yoon, C., Rusu, A., Cobon-Kerr, J., Hahn, C., Luo, J.,  
 1031 Yuvein, Zhu, Ahuja, N., Benenson, R., Kaufman, R. L.,  
 1032 Yu, H., Hightower, L., Zhang, J., Ni, D., Hendricks, L. A.,  
 1033 Wang, G., Yona, G., Jain, L., Barrio, P., Bhupatiraju, S.,  
 1034 Velusamy, S., Dafoe, A., Riedel, S., Thomas, T., Yuan, Z.,  
 1035 Bellaïche, M., Panthaplackel, S., Kloboves, K., Jauhari,  
 1036 S., Akbulut, C., Davchev, T., Gladchenko, E., Madras,  
 1037 D., Chuklin, A., Hill, T., Yuan, Q., Madhavan, M., Leon-  
 1038 hard, L., Scandinaro, D., Chen, Q., Niu, N., Douillard,  
 1039 A., Damoc, B., Onoe, Y., Pedregosa, F., Bertsch, F., Le-  
 1040 ichner, C., Pagadora, J., Malmaud, J., Ponda, S., Twigg,  
 1041 A., Duzhyi, O., Shen, J., Wang, M., Garg, R., Chen, J.,  
 1042 Evcı, U., Lee, J., Liu, L., Kojima, K., Yamaguchi, M., Ra-  
 1043 jendran, A., Piergiovanni, A., Rajendran, V. K., Fornoni,  
 1044 M., Ibagon, G., Ragan, H., Khan, S. M., Blitzer, J., Bun-  
 ner, A., Sun, G., Kosakai, T., Lundberg, S., Elue, N.,  
 Guu, K., Park, S., Park, J., Narayanaswamy, A., Wu, C.,  
 Mudigonda, J., Cohn, T., Mu, H., Kumar, R., Graesser,  
 L., Zhang, Y., Killam, R., Zhuang, V., Giménez, M., Jishi,  
 W. A., Ley-Wild, R., Zhai, A., Osawa, K., Cedillo, D.,  
 Liu, J., Upadhyay, M., Sieniek, M., Sharma, R., Paine,  
 T., Angelova, A., Addepalli, S., Parada, C., Majumder,  
 K., Lamp, A., Kumar, S., Deng, X., Myaskovsky, A.,  
 Sabolić, T., Dudek, J., York, S., de Chaumont Quitry, F.,  
 Nie, J., Cattle, D., Gunjan, A., Piot, B., Khawaja, W.,  
 Bang, S., Wang, S., Khodadadeh, S., R. R., Rawlani, P.,  
 Powell, R., Lee, K., Griesser, J., Oh, G., Magalhaes, C.,  
 Li, Y., Tokumine, S., Vogel, H. N., Hsu, D., BC, A., Jin-  
 dal, D., Cohen, M., Yang, Z., Yuan, J., de Cesare, D.,  
 Bruguier, T., Xu, J., Roy, M., Jacovi, A., Belov, D., Arya,  
 R., Meadowlark, P., Cohen-Ganor, S., Ye, W., Morris-  
 Suzuki, P., Banzal, P., Song, G., Ponnuramu, P., Zhang,  
 F., Scrivener, G., Zaiem, S., Rochman, A. R., Han, K.,  
 Ghazi, B., Lee, K., Drath, S., Suo, D., Girgis, A., Shenoy,  
 P., Nguyen, D., Eck, D., Gupta, S., Yan, L., Carreira, J.,  
 Gulati, A., Sang, R., Mirylenka, D., Cooney, E., Chou,  
 E., Ling, M., Fan, C., Coleman, B., Tubone, G., Kumar,  
 R., Baldridge, J., Hernandez-Campos, F., Lazaridou, A.,  
 Besley, J., Yona, I., Bulut, N., Wellens, Q., Piergiovanni,  
 A., George, J., Green, R., Han, P., Tao, C., Clark, G.,  
 You, C., Abdolmaleki, A., Fu, J., Chen, T., Chaugule, A.,  
 Chandorkar, A., Rahman, A., Thompson, W., Koanan-  
 takool, P., Bernico, M., Ren, J., Vlasov, A., Vassilvitskii,  
 S., Kula, M., Liang, Y., Kim, D., Huang, Y., Ye, C., Lep-  
 ikhin, D., and Helmholz, W. Gemini 2.5: Pushing the  
 frontier with advanced reasoning, multimodality, long  
 context, and next generation agentic capabilities, 2025.  
 URL <https://arxiv.org/abs/2507.06261>.
- Domhan, T., Springenberg, J. T., and Hutter, F. Speeding  
 up automatic hyperparameter optimization of deep neural  
 networks by extrapolation of learning curves. In *Proceed-  
 ings of the 24th International Conference on Artificial In-  
 telligence, IJCAI'15*, pp. 3460–3468. AAAI Press, 2015.  
 ISBN 9781577357384.
- Douze, M., Guzhva, A., Deng, C., Johnson, J., Szilvasy, G.,  
 Mazaré, P.-E., Lomeli, M., Hosseini, L., and Jégou, H.  
 The faiss library. 2024.
- Google. Cloud translation api reference. <https://cloud.google.com/translate/docs/reference/rest>. Accessed: 2026-01-24.
- Hahn, M. The computational power of transformers. *arXiv preprint arXiv:2006.09286*, 2020a.
- Hahn, M. Theoretical limitations of self-attention in neural  
 sequence models. *Transactions of the Association for  
 Computational Linguistics*, 8:156–171, 2020b.

- 1045 Hernandez, D., Brown, T., Conerly, T., DasSarma, N.,  
 1046 Drain, D., El-Showk, S., Elhage, N., Hatfield-Dodds,  
 1047 Z., Henighan, T., Hume, T., Johnston, S., Mann, B.,  
 1048 Olah, C., Olsson, C., Amodei, D., Joseph, N., Ka-  
 1049 plan, J., and McCandlish, S. Scaling laws and inter-  
 1050 pretability of learning from repeated data, 2022. URL <https://arxiv.org/abs/2205.10487>.
- 1051  
 1052 Hestness, J., Narang, S., Ardalani, N., Diamos, G., Jun, H.,  
 1053 Kianinejad, H., Patwary, M. M. A., Yang, Y., and Zhou,  
 1054 Y. Deep learning scaling is predictable, empirically, 2017.  
 1055 URL <https://arxiv.org/abs/1712.00409>.
- 1056  
 1057 Hoffmann, J., Borgeaud, S., Mensch, A., Buchatskaya, E.,  
 1058 Cai, T., Rutherford, E., de Las Casas, D., Hendricks,  
 1059 L. A., Welbl, J., Clark, A., Hennigan, T., Noland, E.,  
 1060 Millican, K., van den Driessche, G., Damoc, B., Guy,  
 1061 A., Osindero, S., Simonyan, K., Elsen, E., Rae, J. W.,  
 1062 Vinyals, O., and Sifre, L. Training compute-optimal large  
 1063 language models, 2022. URL <https://arxiv.org/abs/2203.15556>.
- 1064  
 1065 Hutter, M. Learning curve theory. *arXiv preprint*  
 1066 *arXiv:2102.04074*, 2021.
- 1067  
 1068 Ivgi, M., Carmon, Y., and Berant, J. Scaling laws under the  
 1069 microscope: Predicting transformer performance from  
 1070 small scale experiments, 2022. URL <https://arxiv.org/abs/2202.06387>.
- 1071  
 1072 Jerad, S., Svete, A., Hao, S., Cotterell, R., and Merrill,  
 1073 W. Context-free recognition with transformers. *arXiv*  
 1074 *preprint arXiv:2601.01754*, 2026.
- 1075  
 1076 Jin, C. and Rinard, M. Emergent representations of pro-  
 1077 gram semantics in language models trained on programs.  
 1078 In *Proceedings of the 41st International Conference on*  
 1079 *Machine Learning, ICML’24*. JMLR.org, 2024.
- 1080  
 1081 Kadra, A., Janowski, M., Wistuba, M., and Grabocka,  
 1082 J. Scaling laws for hyperparameter optimization, 2023.  
 1083 URL <https://arxiv.org/abs/2302.00441>.
- 1084  
 1085 Kalra, D. S. and Barkeshli, M. Why warmup the learning  
 1086 rate? underlying mechanisms and improvements, 2024.  
 1087 URL <https://arxiv.org/abs/2406.09405>.
- 1088  
 1089 Kang, F., Ardalani, N., Kuchnik, M., Emad, Y., Elhoushi,  
 1090 M., Sengupta, S., Li, S.-W., Raghavendra, R., Jia, R.,  
 1091 and Wu, C.-J. Demystifying synthetic data in llm pre-  
 1092 training: A systematic study of scaling laws, benefits,  
 1093 and pitfalls, 2025. URL <https://arxiv.org/abs/2510.01631>.
- 1094  
 1095 Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B.,  
 1096 Chess, B., Child, R., Gray, S., Radford, A., Wu, J.,  
 1097 and Amodei, D. Scaling laws for neural language mod-  
 1098 els, 2020. URL <https://arxiv.org/abs/2001.08361>.
- 1099  
 1100 Khan, A., Underwood, R., Siebenschuh, C., Babuji, Y.,  
 1101 Ajith, A., Hippe, K., Gokdemir, O., Brace, A., Chard,  
 1102 K., and Foster, I. Lshbloom: Memory-efficient, extreme-  
 1103 scale document deduplication, 2025. URL <https://arxiv.org/abs/2411.04257>.
- 1104  
 1105 Koh, P. W. and Liang, P. Understanding black-box pre-  
 1106 dictions via influence functions, 2020. URL <https://arxiv.org/abs/1703.04730>.
- 1107  
 1108 Kusupati, A., Bhatt, G., Rege, A., Wallingford, M., Sinha,  
 1109 A., Ramanujan, V., Howard-Snyder, W., Chen, K.,  
 1110 Kakade, S., Jain, P., and Farhadi, A. Matryoshka represen-  
 1111 tation learning, 2024. URL <https://arxiv.org/abs/2205.13147>.
- 1112  
 1113 Lee, K., Ippolito, D., Nystrom, A., Zhang, C., Eck, D.,  
 1114 Callison-Burch, C., and Carlini, N. Deduplicating train-  
 1115 ing data makes language models better, 2022. URL <https://arxiv.org/abs/2107.06499>.
- 1116  
 1117 Manku, G. S., Jain, A., and Das Sarma, A. Detecting near-  
 1118 duplicates for web crawling. In *Proceedings of the 16th*  
 1119 *International Conference on World Wide Web, WWW ’07*,  
 1120 pp. 141–150, New York, NY, USA, 2007. Association for  
 1121 Computing Machinery. ISBN 9781595936547. doi: 10.  
 1122 1145/1242572.1242592. URL <https://doi.org/10.1145/1242572.1242592>.
- 1123  
 1124 McCandlish, S., Kaplan, J., Amodei, D., and Team, O. D.  
 1125 An empirical model of large-batch training, 2018. URL <https://arxiv.org/abs/1812.06162>.
- 1126  
 1127 Mishra, S., Panda, R., Phoo, C. P., Chen, C.-F. R., Karlinsky,  
 1128 L., Saenko, K., Saligrama, V., and Feris, R. S. Task2sim:  
 1129 Towards effective pre-training and transfer from synthetic  
 1130 data. In *Proceedings of the IEEE/CVF Conference on*  
 1131 *Computer Vision and Pattern Recognition (CVPR)*, pp.  
 1132 9194–9204, June 2022.
- 1133  
 1134 Muennighoff, N., Rush, A. M., Barak, B., Scao, T. L., Piktus,  
 1135 A., Tazi, N., Pyysalo, S., Wolf, T., and Raffel, C. Scaling  
 1136 data-constrained language models, 2025. URL <https://arxiv.org/abs/2305.16264>.
- 1137  
 1138 Nguyen, T., Li, Y., Golovneva, O., Zettlemoyer, L., Oh, S.,  
 1139 Schmidt, L., and Li, X. Recycling the web: A method  
 1140 to enhance pre-training data quality and quantity for lan-  
 1141 guage models, 2025. URL <https://arxiv.org/abs/2506.04689>.
- 1142  
 1143 Penedo, G., Kydlíček, H., allal, L. B., Lozhkov, A., Mitchell,  
 1144 M., Raffel, C., Werra, L. V., and Wolf, T. The fineweb  
 1145 datasets: Decanting the web for the finest text data  
 1146 at scale, 2024. URL <https://arxiv.org/abs/2406.17557>.

- 1100 Peng, Z., Wang, Z., and Deng, D. Near-duplicate sequence  
 1101 search at scale for large language model memorization  
 1102 evaluation. *Proc. ACM Manag. Data*, 1(2), June 2023.  
 1103 doi: 10.1145/3589324. URL [https://doi.org/10.](https://doi.org/10.1145/3589324)  
 1104 [1145/3589324](https://doi.org/10.1145/3589324).
- 1105 Porian, T., Wortsman, M., Jitsev, J., Schmidt, L., and Car-  
 1106 mon, Y. Resolving discrepancies in compute-optimal  
 1107 scaling of language models, 2025. URL [https://](https://arxiv.org/abs/2406.19146)  
 1108 [arxiv.org/abs/2406.19146](https://arxiv.org/abs/2406.19146).
- 1109 Pruthi, G., Liu, F., Sundararajan, M., and Kale, S. Es-  
 1110 timating training data influence by tracing gradient  
 1111 descent, 2020. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2002.08484)  
 1112 [2002.08484](https://arxiv.org/abs/2002.08484).
- 1113 Qin, Z., Dong, Q., Zhang, X., Dong, L., Huang, X., Yang,  
 1114 Z., Khademi, M., Zhang, D., Awadalla, H. H., Fung,  
 1115 Y. R., Chen, W., Cheng, M., and Wei, F. Scaling laws of  
 1116 synthetic data for language models, 2025. URL [https://](https://arxiv.org/abs/2503.19551)  
 1117 [arxiv.org/abs/2503.19551](https://arxiv.org/abs/2503.19551).
- 1118 Qwen, ., Yang, A., Yang, B., Zhang, B., Hui, B., Zheng,  
 1119 B., Yu, B., Li, C., Liu, D., Huang, F., Wei, H., Lin, H.,  
 1120 Yang, J., Tu, J., Zhang, J., Yang, J., Yang, J., Zhou, J.,  
 1121 Lin, J., Dang, K., Lu, K., Bao, K., Yang, K., Yu, L.,  
 1122 Li, M., Xue, M., Zhang, P., Zhu, Q., Men, R., Lin, R.,  
 1123 Li, T., Tang, T., Xia, T., Ren, X., Ren, X., Fan, Y., Su,  
 1124 Y., Zhang, Y., Wan, Y., Liu, Y., Cui, Z., Zhang, Z., and  
 1125 Qiu, Z. Qwen2.5 technical report, 2025. URL [https://](https://arxiv.org/abs/2412.15115)  
 1126 [arxiv.org/abs/2412.15115](https://arxiv.org/abs/2412.15115).
- 1127 Rosenfeld, J. S., Rosenfeld, A., Belinkov, Y., and Shavit,  
 1128 N. A constructive prediction of the generalization er-  
 1129 ror across scales, 2019. URL [https://arxiv.org/](https://arxiv.org/abs/1909.12673)  
 1130 [abs/1909.12673](https://arxiv.org/abs/1909.12673).
- 1131 Schaeffer, R., Miranda, B., and Koyejo, S. Are emergent  
 1132 abilities of large language models a mirage?, 2023. URL  
 1133 <https://arxiv.org/abs/2304.15004>.
- 1134 Schoenholz, S. S., Gilmer, J., Ganguli, S., and Sohl-  
 1135 Dickstein, J. Deep information propagation, 2017. URL  
 1136 <https://arxiv.org/abs/1611.01232>.
- 1137 Schulz, L. Y., Mitropolsky, D., and Poggio, T. Unraveling  
 1138 syntax: How language models learn context-free gram-  
 1139 mars. *arXiv preprint arXiv:2510.02524*, 2025.
- 1140 Sclocchi, A., Favero, A., Itzhak Levi, N., and Wyart,  
 1141 M. Probing the latent hierarchical structure of data  
 1142 via diffusion models\*. *Journal of Statistical Mechan-*  
 1143 *ics: Theory and Experiment*, 2025(8):084005, aug 2025.  
 1144 doi: 10.1088/1742-5468/aded6c. URL [https://doi.](https://doi.org/10.1088/1742-5468/aded6c)  
 1145 [org/10.1088/1742-5468/aded6c](https://doi.org/10.1088/1742-5468/aded6c).
- 1146 Simpson, E. H. Measurement of diversity. *Nature*, 163  
 1147 (4148):688–688, 1949. doi: 10.1038/163688a0. URL  
 1148 <https://doi.org/10.1038/163688a0>.
- 1149 Sutton, R. The bitter lesson. [http://](http://www.incompleteideas.net/IncIdeas/BitterLesson.html)  
 1150 [www.incompleteideas.net/IncIdeas/](http://www.incompleteideas.net/IncIdeas/BitterLesson.html)  
 1151 [BitterLesson.html](http://www.incompleteideas.net/IncIdeas/BitterLesson.html), March 2019.
- 1152 Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux,  
 1153 M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro,  
 1154 E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., and  
 Lample, G. Llama: Open and efficient foundation lan-  
 guage models, 2023. URL [https://arxiv.org/](https://arxiv.org/abs/2302.13971)  
[abs/2302.13971](https://arxiv.org/abs/2302.13971).
- Vera, H. S., Dua, S., Zhang, B., Salz, D., Mullins, R., Pa-  
nyam, S. R., Smoot, S., Naim, I., Zou, J., Chen, F., Cer,  
D., Lisak, A., Choi, M., Gonzalez, L., Sanseviero, O.,  
Cameron, G., Ballantyne, I., Black, K., Chen, K., Wang,  
W., Li, Z., Martins, G., Lee, J., Sherwood, M., Ji, J.,  
Wu, R., Zheng, J., Singh, J., Sharma, A., Sreepathihalli,  
D., Jain, A., Elarabawy, A., Co, A., Doumanoglou, A.,  
Samari, B., Hora, B., Potetz, B., Kim, D., Alfonseca,  
E., Moiseev, F., Han, F., Gomez, F. P., Ábrego, G. H.,  
Zhang, H., Hui, H., Han, J., Gill, K., Chen, K., Chen, K.,  
Shanbhogue, M., Boratko, M., Suganthan, P., Duddu, S.  
M. K., Mariserla, S., Ariafar, S., Zhang, S., Zhang, S.,  
Baumgartner, S., Goenka, S., Qiu, S., Dabral, T., Walker,  
T., Rao, V., Khawaja, W., Zhou, W., Ren, X., Xia, Y.,  
Chen, Y., Chen, Y.-T., Dong, Z., Ding, Z., Visin, F., Liu,  
G., Zhang, J., Kenealy, K., Casbon, M., Kumar, R., Mes-  
nard, T., Gleicher, Z., Brick, C., Lacombe, O., Roberts,  
A., Yin, Q., Sung, Y., Hoffmann, R., Warkentin, T., Joulin,  
A., Duerig, T., and Seyedhosseini, M. Embeddinggemma:  
Powerful and lightweight text representations, 2025. URL  
<https://arxiv.org/abs/2509.20354>.
- Wang, H., Ma, S., Dong, L., Huang, S., Zhang, D., and Wei,  
F. Deepnet: Scaling transformers to 1,000 layers, 2022.  
URL <https://arxiv.org/abs/2203.00555>.
- Wang, H., Minervini, P., and Ponti, E. M. Probing the  
emergence of cross-lingual alignment during llm train-  
ing, 2024. URL [https://arxiv.org/abs/2406.](https://arxiv.org/abs/2406.13229)  
[13229](https://arxiv.org/abs/2406.13229).
- Xiong, R., Yang, Y., He, D., Zheng, K., Zheng, S., Xing, C.,  
Zhang, H., Lan, Y., Wang, L., and Liu, T.-Y. On layer  
normalization in the transformer architecture, 2020. URL  
<https://arxiv.org/abs/2002.04745>.
- Yan, T., Wen, H., Li, B., Luo, K., Chen, W., and Lyu,  
K. Larger datasets can be repeated more: A theoretical  
analysis of multi-epoch scaling in linear regression, 2025.  
URL <https://arxiv.org/abs/2511.13421>.
- Yang, A., Li, A., Yang, B., Zhang, B., Hui, B., Zheng,  
B., Yu, B., Gao, C., Huang, C., Lv, C., Zheng, C., Liu,  
D., Zhou, F., Huang, F., Hu, F., Ge, H., Wei, H., Lin,  
H., Tang, J., Yang, J., Tu, J., Zhang, J., Yang, J., Yang,  
J., Zhou, J., Zhou, J., Lin, J., Dang, K., Bao, K., Yang,

1155 K., Yu, L., Deng, L., Li, M., Xue, M., Li, M., Zhang,  
1156 P., Wang, P., Zhu, Q., Men, R., Gao, R., Liu, S., Luo,  
1157 S., Li, T., Tang, T., Yin, W., Ren, X., Wang, X., Zhang,  
1158 X., Ren, X., Fan, Y., Su, Y., Zhang, Y., Zhang, Y., Wan,  
1159 Y., Liu, Y., Wang, Z., Cui, Z., Zhang, Z., Zhou, Z., and  
1160 Qiu, Z. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.  
1161  
1162 Yang, G., Hu, E. J., Babuschkin, I., Sidor, S., Liu, X., Farhi,  
1163 D., Ryder, N., Pachocki, J., Chen, W., and Gao, J. Tensor  
1164 programs v: Tuning large neural networks via zero-shot  
1165 hyperparameter transfer, 2022. URL <https://arxiv.org/abs/2203.03466>.  
1166  
1167 Yang, Z., Band, N., Li, S., Candès, E., and Hashimoto, T.  
1168 Synthetic continued pretraining, 2024. URL <https://arxiv.org/abs/2409.07431>.  
1169  
1170 Zhang, H., Dauphin, Y. N., and Ma, T. Fixup initialization:  
1171 Residual learning without normalization, 2019. URL  
1172 <https://arxiv.org/abs/1901.09321>.  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187  
1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209

## A. Related Work

Predictable scaling has been central to deep learning since early work on learning curves and performance prediction (Domhan et al., 2015; Hestness et al., 2017). Despite the analytic intractability of modern neural networks, empirical scaling laws often predict loss as a function of model size, data size, and compute with high accuracy (Schoenholz et al., 2017; Rosenfeld et al., 2019; Kaplan et al., 2020; Hoffmann et al., 2022). Scaling predictability governs many aspects of training recipes, including parameterization, learning rates, depth-to-width ratios, initialization, warmup, and batch size (Kadra et al., 2023; Yang et al., 2022; Xiong et al., 2020; Wang et al., 2022; Zhang et al., 2019; Kalra & Barkeshli, 2024; McCandlish et al., 2018). A persistent challenge is identifying *scale-dependent* factors that undermine predictable extrapolation (Ivgi et al., 2022; Schaeffer et al., 2023; Porian et al., 2025). Our work highlights a new source of scale dependence linked to semantic duplicates and their frequency at web scale.

We build on work studying repeated or low-uniqueness training data. Hernandez et al. (2022) showed that repeating a small subset of training examples can substantially reduce the effective parameter size predicted by scaling, and subsequent work reports that the effects of repetition can grow with scale. Because near-duplicates are common in web corpora (Peng et al., 2023), practical pipelines deploy hashing and approximate matching to identify and eliminate such “fuzzy duplicates” (Broder, 1997; Manku et al., 2007; Khan et al., 2025). Unlike settings with *explicit* repeats (e.g., injected duplicates or many epochs) (Hernandez et al., 2022; Muennighoff et al., 2025; Yan et al., 2025), we emphasize an *implicit* and *scale-dependent* notion of repetition: as models become semantically sensitive, semantically equivalent documents may function as duplicates, and the prevalence of semantic collisions grows with corpus scale.

Our gradient-based measurements relate to work that treats gradients as training signals and influence proxies (Pruthi et al., 2020; Koh & Liang, 2020). Our observation that semantic structure becomes more salient over training aligns with prior work on the emergence and probing of semantic representations during pretraining (Jin & Rinard, 2024; Chen et al., 2024b; Aljaafari et al., 2025; Wang et al., 2024). Unlike prior work that treats semantic emergence as a purely beneficial phenomenon, we connect it to a potential failure mode: semantic duplicates can create redundant training signals that disproportionately affect capable LMs.

A complementary line of work studies how neural networks learn hierarchical and compositional structure. Recent theory introduces stylized latent-data models such as the Random Hierarchy Model (RHM), in which examples are generated by composing features along a tree (analogous to a grammar derivation), yielding sharp predictions about which levels of the hierarchy are learnable at a given scale (Cagnetta et al., 2024b; 2025; Cagnetta & Wyart, 2024; Sclocchi et al., 2025). In language modeling, formal-language and grammar-based probes have been used to analyze whether attention-based architectures can represent and generalize hierarchical dependencies, including theoretical limitations of self-attention (Hahn, 2020b) and empirical studies of Transformer recognition of formal languages (Bhattachamishra et al., 2020). Most recently, Schulz et al. (2025) directly characterizes how language models learn context-free grammars over training. Our work connects to these perspectives by highlighting a distinct consequence of learning deeper invariances: as models become semantically/compositionally sensitive, semantically equivalent documents increasingly behave as effective duplicates, amplifying the impact of semantic collisions at corpus scale.

## B. Limitations and Future Work

This work has several limitations.

First, because we are unable to train multi-billion-parameter models on our compute budget, we simulated the impact of semantic duplicates using exact duplicates. We drew appropriate comparisons by basing estimates on mean cosine similarity of semantic embeddings, but behavior may differ for large models trained on semantic rather than exact duplicates. Future work could validate our findings on large models by training on semantically deduplicated datasets. The formation of such datasets requires future research attention and resources.

As another limitation, the semantic embeddings that we used were from EmbeddingGemma-300m. Although this is a state-of-the-art embedding model used by many frontier labs today for data exploration, it still does not produce perfectly isotropic or representative embeddings. Training better embedding models that utilize the latent space more efficiently could improve confidence in results.

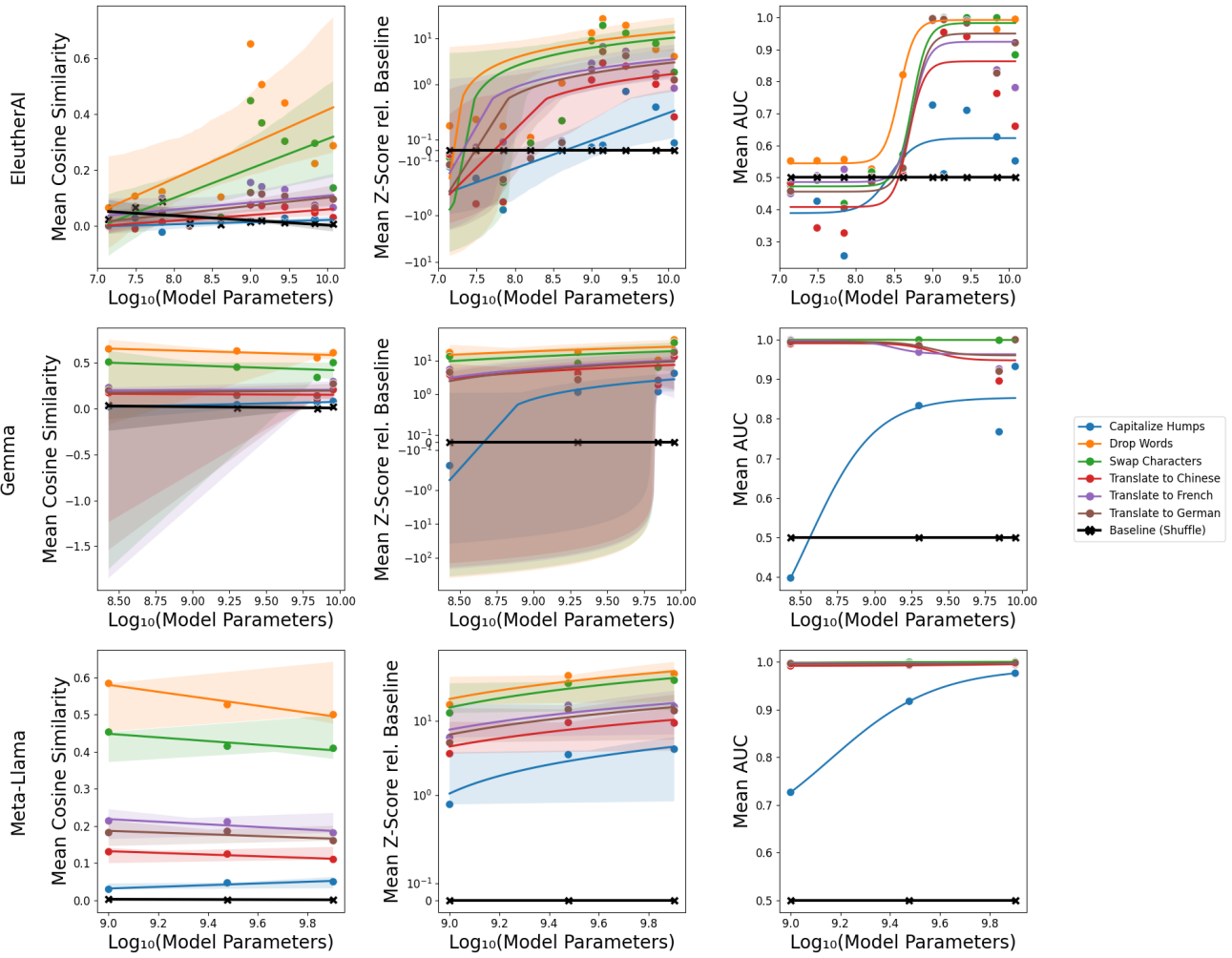


Figure 8. Semantic-preserving transformations yield more aligned gradients for larger/stronger models: We display the same data as in Figure 1.

## C. Algorithm for Gradient Comparison

---

**Algorithm 1** Experiment 1: Gradient similarity for semantic duplicates with negative baseline distribution

---

**Require:** Base texts  $\{x_i\}_{i=1}^N$ ; transformation set  $\mathcal{T}$ ; models  $\{f^{(k)}\}_{k=1}^K$  with checkpoints  $\{\theta_{k,s}\}$ ; loss  $\ell$ ; token budget  $T$ ; number of negative pairing rounds  $R$

**Ensure:** For each  $(k, s, \tau)$ : positives  $\{s_i^+\}$ , baseline negatives  $\mathcal{S}^-$ , AUC, Z-score summary

- 1: Truncate each  $x_i$  to at most  $T$  tokens (model tokenizer)
- 2: **for**  $k \leftarrow 1$  **to**  $K$  **do**
- 3:     **for all** checkpoints  $s$  of model  $k$  **do**
- 4:         Load parameters  $\theta \leftarrow \theta_{k,s}$
- 5:         **for**  $i \leftarrow 1$  **to**  $N$  **do**
- 6:              $g_i \leftarrow \nabla_{\theta} \ell(x_i; \theta)$
- 7:         **end for**
- 8:          $\mathcal{S}^- \leftarrow []$
- 9:         **for**  $r \leftarrow 1$  **to**  $R$  **do**
- 10:             Sample pairing map  $j_r(\cdot)$  such that  $j_r(i) \neq i$  for all  $i$
- 11:             **for**  $i \leftarrow 1$  **to**  $N$  **do**
- 12:                  $\mathcal{S}^- \text{.append}(\cos(g_i, g_{j_r(i)}))$
- 13:             **end for**
- 14:         **end for**
- 15:         Compute  $\mu^-$  and  $\sigma^-$  from  $\mathcal{S}^-$
- 16:         **for all**  $\tau \in \mathcal{T}$  **do**
- 17:             **for**  $i \leftarrow 1$  **to**  $N$  **do**
- 18:                  $x_i^{(\tau)} \leftarrow \tau(x_i)$
- 19:                  $g_i^{(\tau)} \leftarrow \nabla_{\theta} \ell(x_i^{(\tau)}; \theta)$
- 20:                  $s_i^+ \leftarrow \cos(g_i, g_i^{(\tau)})$
- 21:                  $z_i \leftarrow (s_i^+ - \mu^-) / \sigma^-$
- 22:             **end for**
- 23:             Compute AUC using positives  $\{s_i^+\}_{i=1}^N$  vs negatives  $\mathcal{S}^-$
- 24:             Summarize Z-scores (e.g., mean/median over  $i$ )
- 25:         **end for**
- 26:     **end for**
- 27: **end for**

▷ Compute gradients for all base texts once

▷ Build a negative baseline distribution from many random pairings

▷ Evaluate each transformation on all texts

---

## D. Deriving The Partner Probability and the Effective Latent Count Approximation

This appendix derives Eq. (29) step by step, starting from an exact expression for a general latent mixture and then giving a controlled approximation that yields the compact form  $1 - \exp(-(N_{\text{meas}} - 1)/K_{\text{eff}})$ .

**Setup (Latent-Mixture Model).** Let  $Z$  be a discrete semantic latent taking values in an index set  $\mathcal{Z}$  with mixture weights  $\{w_z\}_{z \in \mathcal{Z}}$ , i.e.  $\mathbb{P}(Z = z) = w_z$  and  $\sum_z w_z = 1$ . Let  $Z_1, \dots, Z_N \stackrel{iid}{\sim} \{w_z\}$  denote the latents of  $N$  independent draws (e.g.  $N = N_{\text{meas}}$  samples from the training stream).

For a given draw  $i$ , define the event that it has at least one same-latent partner among the other  $N - 1$  draws:

$$A_i := \{\exists j \neq i : Z_j = Z_i\}.$$

By exchangeability,  $\mathbb{P}(A_i)$  does not depend on  $i$ , so we analyze  $A_1$ .

**D.1. Exact Expression for  $q_N$** 

Define

$$q_N := \mathbb{P}(A_1) = \mathbb{P}(\exists j \neq 1 : Z_j = Z_1).$$

Condition on  $Z_1 = z$ . Then each of the remaining  $N - 1$  draws matches  $z$  with probability  $w_z$ , independently, so the number of matches among draws  $2, \dots, N$  is Binomial( $N - 1, w_z$ ). Hence,

$$\mathbb{P}(\text{no partner for draw 1} \mid Z_1 = z) = \mathbb{P}(Z_2 \neq z, \dots, Z_N \neq z \mid Z_1 = z) = (1 - w_z)^{N-1}. \quad (31)$$

Averaging over  $Z_1$  gives the exact identity

$$\mathbb{P}(\text{no partner for draw 1}) = \sum_{z \in \mathcal{Z}} \mathbb{P}(Z_1 = z) (1 - w_z)^{N-1} = \sum_z w_z (1 - w_z)^{N-1}. \quad (32)$$

Therefore,

$$q_N = 1 - \sum_z w_z (1 - w_z)^{N-1}. \quad (33)$$

This is the first line of Eq. (29) and is *exact* for any discrete mixture.

**D.2. From Mixture Weights to  $K_{\text{eff}}$** 

A key quantity is the probability that *two independent draws* share the same latent:

$$p_{\text{lat}} := \mathbb{P}(Z = Z') = \sum_z \mathbb{P}(Z = z) \mathbb{P}(Z' = z) = \sum_z w_z^2. \quad (34)$$

This is the Simpson collision probability. It induces the *Simpson effective number of latents*

$$K_{\text{eff}} := \frac{1}{p_{\text{lat}}} = \frac{1}{\sum_z w_z^2}. \quad (35)$$

In the uniform- $K$  case ( $w_z = 1/K$  for  $z = 1, \dots, K$ ), we have  $p_{\text{lat}} = 1/K$  and thus  $K_{\text{eff}} = K$ .

**D.3. Approximation: Rare-Collision / No-Heavy-Modes Regime**

We now explain the approximation

$$q_N \approx 1 - \exp\left(- (N - 1) \sum_z w_z^2\right) = 1 - \exp\left(- \frac{N - 1}{K_{\text{eff}}}\right).$$

**Step 1: Poissonizing the Binomial for Small  $w_z$ .** For small  $w_z$ , the binomial Binomial( $N - 1, w_z$ ) is well-approximated by Poisson( $\lambda_z$ ) with rate  $\lambda_z = (N - 1)w_z$ . In particular,

$$(1 - w_z)^{N-1} = \exp((N - 1) \log(1 - w_z)) = \exp(-(N - 1)w_z + O((N - 1)w_z^2)), \quad (36)$$

so when  $\max_z w_z \ll 1$  and  $(N - 1) \max_z w_z^2$  is not too large, we may use

$$(1 - w_z)^{N-1} \approx \exp(-(N - 1)w_z). \quad (37)$$

Plugging (37) into (33) yields

$$q_N \approx 1 - \sum_z w_z \exp(-(N - 1)w_z). \quad (38)$$

**Step 2: Collapsing the Mixture to a Single Effective rate.** Let  $W$  be the random variable  $W_{Z_1}$  when  $Z_1 \sim \{w_z\}$ , i.e.  $\mathbb{P}(W = w_z) = w_z$ . Then (38) can be written compactly as

$$\sum_z w_z \exp(-(N-1)w_z) = \mathbb{E}\left[e^{-(N-1)W}\right]. \quad (39)$$

Moreover,

$$\mathbb{E}[W] = \sum_z w_z \cdot w_z = \sum_z w_z^2 = p_{\text{lat}} = \frac{1}{K_{\text{eff}}}.$$

If the mixture has *no heavy modes* (informally:  $w_z \ll 1$  and the distribution of  $W$  is not extremely spread out), we can approximate the expectation in (39) by its mean-field form:

$$\mathbb{E}\left[e^{-(N-1)W}\right] \approx \exp(-(N-1)\mathbb{E}[W]) = \exp\left(- (N-1) \sum_z w_z^2\right) = \exp\left(-\frac{N-1}{K_{\text{eff}}}\right). \quad (40)$$

A standard way to justify (40) is via a cumulant (Taylor) expansion:

$$\log \mathbb{E}[e^{-aW}] = -a \mathbb{E}[W] + \frac{a^2}{2} \text{Var}(W) + O(a^3 \mathbb{E}[|W - \mathbb{E}W|^3]), \quad a \rightarrow 1,$$

so if  $a^2 \text{Var}(W)$  is small compared to  $a \mathbb{E}[W]$  (i.e.  $W$  is concentrated around its mean at the scale relevant for  $a$ ), then  $\log \mathbb{E}[e^{-aW}] \approx -a \mathbb{E}[W]$  and (40) follows.

**Putting the Steps Together.** Combining (38) with (40) yields

$$q_N \approx 1 - \exp\left(- (N-1) \sum_z w_z^2\right) = 1 - \exp\left(-\frac{N-1}{K_{\text{eff}}}\right). \quad (41)$$

This matches Eq. (29) in the main text.

#### D.4. Sanity Check: uniform- $K$ case

If  $w_z = 1/K$  for  $z = 1, \dots, K$ , then (33) becomes

$$q_N = 1 - \sum_{z=1}^K \frac{1}{K} \left(1 - \frac{1}{K}\right)^{N-1} = 1 - \left(1 - \frac{1}{K}\right)^{N-1},$$

and using  $\log(1-x) \approx -x$  gives

$$q_N \approx 1 - \exp\left(-\frac{N-1}{K}\right).$$

Since  $K_{\text{eff}} = K$  in the uniform case, Eq. (41) recovers the standard occupancy approximation exactly up to the usual  $\log(1-x) \approx -x$  step.

#### D.5. Remark: What Breaks when there are Heavy Modes?

If some  $w_z$  are not small (a few “heavy” semantics), then: (i) the Poisson approximation (37) can be inaccurate for those modes, and (ii) the mean-field collapse (40) can be poor because  $W$  is no longer concentrated. In that case, Eq. (33) remains correct and can be used directly, and  $K_{\text{eff}}$  still meaningfully summarizes pairwise collision probability via (35), but the single-exponential approximation to  $q_N$  may systematically overestimate collision probability.

### E. A First-Principles Model of Duplicate-Limited Scaling via Hutter-Style Learning Curves

This appendix derives a collision-aware scaling correction by combining: (i) a Hutter-style learning-curve model in which performance improves as a power law in the number of *independent* training signals, and (ii) a reduction of independent signal due to duplicates/semantic collisions. The goal is not a fully realistic theory of language modeling, but a minimal mechanism that explains why a plane law of the form  $\Delta(C, K) \approx a C^\beta K^{-\gamma}$  arises naturally.

**Step 1: A Hutter-Style “New Information” Learning Curve.** A classic abstraction (learning curve theory) models learning progress as driven by discovering previously unseen “features” or “types” in a heavy-tailed environment. Concretely, let  $z$  denote a latent “type” (semantic class, rule, or pattern) with weights  $\{w_z\}$ . Consider the idealized memorization learner that, upon seeing *one* example of type  $z$ , can thereafter predict  $z$  perfectly, while unseen types incur a fixed excess loss. In this model, the expected excess risk after  $n$  iid draws is proportional to the probability mass of unseen types:

$$\epsilon(n) = \sum_z w_z (1 - w_z)^n, \quad (42)$$

a form that appears in learning-curve theory and is closely related to occupancy/species discovery. (For a detailed treatment and conditions under which heavy tails yield power laws, see [Hutter \(2021\)](#).)

**Step 2: Power Laws from Heavy Tails.** If the type weights follow a Zipf/regularly varying tail,  $\epsilon(n)$  follows a power law:

$$\epsilon(n) \propto n^{-\alpha} \quad \text{for some } \alpha \in (0, 1), \quad (43)$$

with  $\alpha$  determined by the tail index of  $\{w_z\}$  (see [Hutter \(2021\)](#)). We use (43) as a generic “first-principles” justification for a power law dependence of excess loss on the amount of *independent* training signal.

**Step 3: Duplicates Reduce the Effective Number of Independent Signals.** In our setting, training examples are not independent sources of new information: duplicates (exact or semantic) induce correlated gradients and therefore reduce the number of effectively independent update directions. Let  $n$  denote the number of training documents processed. Let  $K$  denote the number of effective semantic classes available (or  $K_{\text{eff}}$  in the main text). Let  $\rho \in [0, 1]$  summarize semantic sensitivity (gradient alignment within a class) as in Eq. (14). Under the correlation model in Eq. (17), Proposition 5.2 implies an effective sample size

$$n_{\text{eff}} = \frac{n}{1 + \rho \frac{n-1}{K}} \approx \frac{n}{1 + r_{\text{eff}}}, \quad r_{\text{eff}} := \rho \frac{n}{K}. \quad (44)$$

Intuitively,  $r_{\text{eff}}$  is an *effective reuse ratio*: when  $r_{\text{eff}} \ll 1$  the stream is mostly novel, and when  $r_{\text{eff}} \gg 1$  the stream is dominated by redundant semantics.

**Step 4: Substitute  $n_{\text{eff}}$  into the Learning Curve.** Assume the excess loss (or excess cross-entropy) is a power law in the *independent* signal count:

$$L(n, K) - L_\star \approx B n_{\text{eff}}^{-\alpha}, \quad (45)$$

where  $L_\star$  is an irreducible floor and  $B > 0$ . For the high-uniqueness baseline (negligible collisions),  $n_{\text{eff}} \approx n$  and  $L_\infty(n) - L_\star \approx B n^{-\alpha}$ . For finite  $K$ , combining (44)–(45) gives

$$L(n, K) - L_\star \approx B n^{-\alpha} (1 + r_{\text{eff}})^\alpha. \quad (46)$$

**Step 5: A Duplicate-Induced Degradation Law.** Define the normalized degradation  $\Delta$  as in Eq. (20):  $\Delta := (L(n, K) - L_\infty(n))/L_\infty(n)$ . Using (46) and  $L_\infty(n) = L_\star + B n^{-\alpha}$ , we obtain

$$\Delta(n, K) \approx \frac{B n^{-\alpha} ((1 + r_{\text{eff}})^\alpha - 1)}{L_\star + B n^{-\alpha}}. \quad (47)$$

In the regime where  $B n^{-\alpha}$  is not negligible relative to  $L_\star$  (typical for the losses in our controlled ladders), the prefactor is slowly varying and (47) is well-approximated by a power law in  $r_{\text{eff}}$ . In particular, when  $r_{\text{eff}} \lesssim 1$  we can linearize:

$$\Delta(n, K) \approx \tilde{\lambda} r_{\text{eff}} = \tilde{\lambda} \rho \frac{n}{K}, \quad (48)$$

where  $\tilde{\lambda}$  absorbs the slowly varying ratio in (47). Equation (48) recovers the main-text intuition that degradation is (approximately) proportional to an effective reuse ratio.

**Step 6: Translating to Compute and the Plane Law.** Let  $C$  denote compute. Over restricted ranges, it is empirically accurate to approximate

$$n(C) \propto C^u, \quad \rho(C) \propto C^v,$$

as in Eq. (23). Substituting into (48) yields

$$\Delta(C, K) \approx a C^{u+v} K^{-1}, \tag{49}$$

which is a *plane law* in  $(\log C, \log K)$  with  $\gamma \approx 1$ . More generally, if one does not linearize (47), the same substitution yields a plane  $\Delta(C, K) \propto C^\beta K^{-\gamma}$  with  $\beta = \eta(u + v)$  and  $\gamma = \eta$  for some effective exponent  $\eta$ , matching Eq. (24).

**Discussion: Why  $\rho(C)$  Should Grow with Scale (and why a Power Law is a Reasonable Local Model).** The parameter  $\rho(C)$  captures the fraction of gradient energy explained by invariances to surface form (Eq. (14)). A growing body of theory and empirical work on *hierarchical/compositional* data suggests that neural networks learn coarse, high-level structure before finer structure, and that deeper invariances require more data/compute. For example, the random hierarchy model (RHM) formalizes language-like hierarchical generation and yields staged learning dynamics where progressively deeper variables become learnable as sample size increases (Cagnetta et al., 2024b;a). Separately, work on formal-language recognition by transformers highlights a connection between model depth/recurrence and the ability to represent hierarchical (context-free) structure, which is a canonical form of compositional invariance (Hahn, 2020a; Jerad et al., 2026). Taken together, these results motivate modeling  $\rho(C)$  as monotone increasing with scale; over the narrow compute ranges used in scaling ladders, a power law approximation  $\rho(C) \propto C^v$  is a parsimonious local model.

**A Reduced-Parameter Variant.** Equation (49) suggests a two-degree-of-freedom correction: fix  $\gamma = 1$  and fit only  $(a, \beta)$  (or even fit  $v$  with  $u$  known from the compute-to-sample mapping). In our controlled ladders, the fitted  $\gamma$  is close to 1, consistent with this linear-reuse regime.

1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619  
1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649

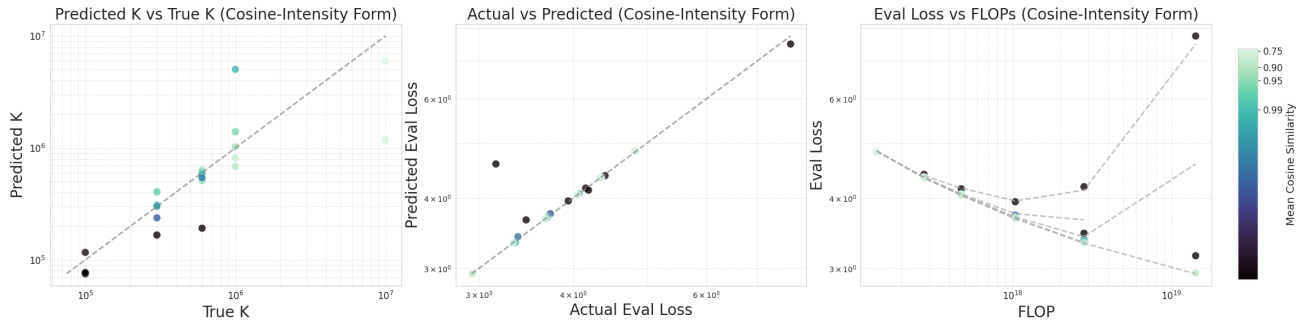


Figure 9. Predictions of eval loss using cosine intensity-based estimation of  $\hat{K}_{\text{eff}}$ .

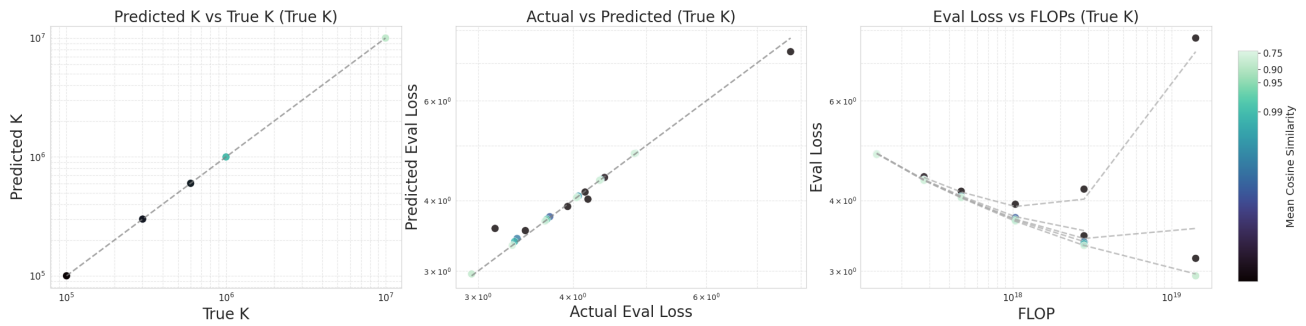


Figure 10. Predictions of eval loss using true  $K$ .