
Internal Data Repetition Destroys Language Models

Anonymous Authors¹

Abstract

Language models are running out of high-quality training data, and even aggressively deduplicated corpora retain some amount of repetition. Earlier controlled studies predated Chinchilla-style scaling laws and could only measure the cost of repetition indirectly. We revisit repetition in the Chinchilla-style scaling regime, using a fitted no-repetition scaling law to report Compute-Equivalent Gain and Compute-Equivalent Loss. We show that repetition damage in this modernized regime is systematic in three ways. First, eval loss is worst at an intermediate repeat count R , so repeating a moderately sized subset many times hurts more than either repeating a large subset a few times or a small subset many times. Second, the location of this peak is well fit by a power law in model size. Finally, when repeated documents make up 10% of training tokens in a controlled exact-document repetition setting, the compute-equivalent loss can be large: on FineWeb-Edu-Dedup, the most damaging repeat count for a Qwen3-style 344M-parameter model at $OT = 1$ matches the loss of a no-repetition run using about 67% of the FLOPs, under our fitted no-repetition scaling law. A misspecified linear regression with verbatim duplicates reproduces the same qualitative non-monotonicity in closed form, suggesting that such peaks can arise from a statistical tradeoff between memorization and generalization. Our findings give practitioners a way to predict which settings waste the most compute before they spend any of it.

1. Introduction

Pretraining has entered a data-constrained regime. The high-quality public text corpora used for frontier training are

¹Anonymous 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.

within a small constant factor of being exhausted (Villalobos et al., 2024; Longpre et al., 2024), and recent flagship corpora such as FineWeb-Edu, DataComp-LM, Dolma, and RedPajama-v2 (Penedo et al., 2024; Li et al., 2024; Soldaini et al., 2024; Weber et al., 2024) have responded with more aggressive deduplication and filtering. Yet aggressive deduplication is not perfect deduplication (Lee et al., 2022; Abbas et al., 2023; Tirumala et al., 2023). Pretraining streams continue to contain near-duplicate documents, paraphrased templates, and semantically redundant web pages, and as scale grows the meaning of “duplicate” itself shifts (Kazdan et al., 2026). Replication and duplicates can vary in type and effect. In this paper, we isolate one controlled case: exact document-level replay of a selected repeated pool. Our question is how much this controlled form can cost.

The closest earlier controlled study, Hernandez et al. (2022), established that repeating a small fraction of training data degrades held-out loss in a non-monotonic way, and framed the damage as a reduction in *effective parameter count*. That framing was natural before Chinchilla-style scaling (Hoffmann et al., 2022) provided a clean compute axis for predictions. The quantity a practitioner allocating a pretraining budget cares about is how many FLOPs a no-repetition run would need to reach the same loss. We replace effective parameter count with *Compute-Equivalent Gain* (CEG), following the compute-equivalent gains framework of Gundlach et al. (2025); Davidson et al. (2023); Meta Superintelligence Labs (2026). For a repeated-data run with loss L and actual compute C_{actual} , CEG is the no-repetition compute required to reach L divided by C_{actual} . We define *Compute-Equivalent Loss* as $CEL = 1 - \text{CEG}$. Thus $CEL = 1$ matches the no-repetition reference, while $CEL < 1$ indicates compute-equivalent loss.

We train Qwen3-style models (Yang et al., 2025a) with $N \in \{34, 48, 63, 93, 153, 344\}$ M parameters on FineWeb-Edu-Dedup (Penedo et al., 2024). We sweep the overtraining multiplier $OT \in \{0.25, 0.5, 1, 2, 4\}$ for all but the larger models due to compute constraints, and sweep the per-document repeat count R on an approximately logarithmic grid from no repeats to as high as $R = 20000$ (we plot up to $R \approx 3000$). Throughout, we fix the repeated-token fraction at $f = 0.1$: 90% of training tokens come from non-repeated documents, and the remaining 10% come from a smaller document pool replayed R times. This fraction is large

enough to produce measurable damage while keeping the bulk of training on unique data, and matches the setup used by Hernandez et al. (2022).

The Chinchilla budget identity $C = 6NT = 120 \cdot OT \cdot N^2$ ties model size, total tokens T , and total compute C together (Gadre et al., 2025; Porian et al., 2024; Kaplan et al., 2020; Sardana et al., 2024). This lets us vary repetition structure inside an otherwise fixed training budget. We find that the most damaging repetition structure is small enough to be encountered many times but large enough to consume meaningful model capacity, and that its location is well fit by a power law in N . Finally, we translate the resulting loss increases into Compute-Equivalent Gain and Loss using a fitted no-repetition scaling law. Under this fitted frontier, the most damaging repeat count at our largest scale has CEG ≈ 0.67 , corresponding to CEL ≈ 0.33 . A misspecified linear regression with verbatim duplicates reproduces the same qualitative non-monotonicity in closed form, suggesting that such peaks can arise from the statistical structure of duplicated samples rather than from attention, depth, or optimizer-specific effects.

Contributions. (i) We measure repeated-data damage in compute-equivalent units by comparing each run against a fitted no-repetition Chinchilla scaling law (§4.3). In our Qwen3-style 344M-parameter model trained on FineWeb-Edu-Dedup at $OT = 1$ (Yang et al., 2025a; Penedo et al., 2024), the most damaging repeat count has CEL ≈ 0.33 . (ii) We identify the intermediate- R regime where damage peaks (§4.1) and show that its location follows the power law $R^{\text{peak}} \propto N^{-0.96}$ (§4.2), making the worst-case configuration predictable using N . (iii) We derive closed-form train and test losses for a misspecified linear regression with verbatim duplicates (§4.4) and recover the same non-monotonic peak in simulation, offering a statistical analogue based on a tradeoff between fitting repeated samples and generalizing beyond repeated samples.

2. Related Work

We position our work against three closely connected lines of research. An extended discussion appears in Appendix A.

Compute-optimal scaling and the Chinchilla scaling law. Kaplan et al. (2020) established the power-law form for transformer loss as a function of compute. Hoffmann et al. (2022) corrected the optimal (N, T) allocation, where N is the parameter count and T is the token count. Porian et al. (2024); Besiroglu et al. (2024) examine the robustness of these fits, Sardana et al. (2024) extend them to inference-aware budgets, and Gadre et al. (2025) show reliable extrapolation through aggressive over-training. We use this functional form, fitted on no-repetition baselines, as the reference curve for CEG and CEL.

Deduplication, memorization, and statistical accounts of overfitting. A large literature studies deduplication and memorization for privacy and contamination reasons (Lee et al., 2022; Kandpal et al., 2022; Carlini et al., 2023; Abbas et al., 2023; Tirumala et al., 2023; Deng et al., 2024; Schaeffer et al., 2026). Our setting is distinct in two ways. First, we exclude the evaluation split from training and from all repeated pools by the train/test split, so any harm we observe must come from a distorted effective training distribution. Second, we hold the fraction of repeated tokens fixed and vary only their concentration, isolating the effect of repetition structure. Our theoretical analysis connects to classical results on double descent and benign overfitting (Belkin et al., 2019; Hastie et al., 2022; Bartlett et al., 2020; Nakkiran et al., 2020), specialized to literal sample duplication. We are not aware of prior work using this block-covariance view to analyze literal duplication.

Repeated data in language model pretraining. Hernandez et al. (2022) repeat a small fraction of training data and observed a non-monotonic test-loss curve. They framed the damage as a reduction in effective parameter count and connected it to degradation of induction heads. Muennighoff et al. (2023) study the complementary regime in which the entire corpus is uniformly repeated, and find that up to roughly four epochs are nearly free. Xue et al. (2023); Maini et al. (2024); Komatsuzaki (2019) examine related repetition, rephrasing, and epoch strategies. Recent work further argues that semantic duplication is itself scale-dependent (Kazdan et al., 2026) and proposes explicit overfitting penalties for data-constrained scaling (Lovelace et al., 2026). We sit between the Hernandez et al. (2022) and Muennighoff et al. (2023) regimes. Previous work (Hernandez et al., 2022) observed that repeated data can produce non-monotonic held-out loss degradation, and reported the damage as a reduction in effective parameter count. We extend that setting by measuring damage as CEG and CEL: for each repeated-data run, we ask how much compute a no-repetition run on a fitted scaling law would need to reach the same loss. We also run the sweep under Chinchilla-style token and compute budgets while holding the repeated-token fraction fixed at $f = 0.1$ and varying its concentration through R , and we fit how the most damaging repeat count shifts with model size.

3. Methods

Repeated data is increasingly hard to avoid, so the practical question becomes which repeat structures are dangerous and how much compute they waste. We fix the repeated-token fraction at $f = 0.1$ and vary only its concentration: the repeated 10% can come from a large pool of data seen a few times or a small pool seen many times. This isolates repetition *structure* from repetition *amount*, letting us identify the

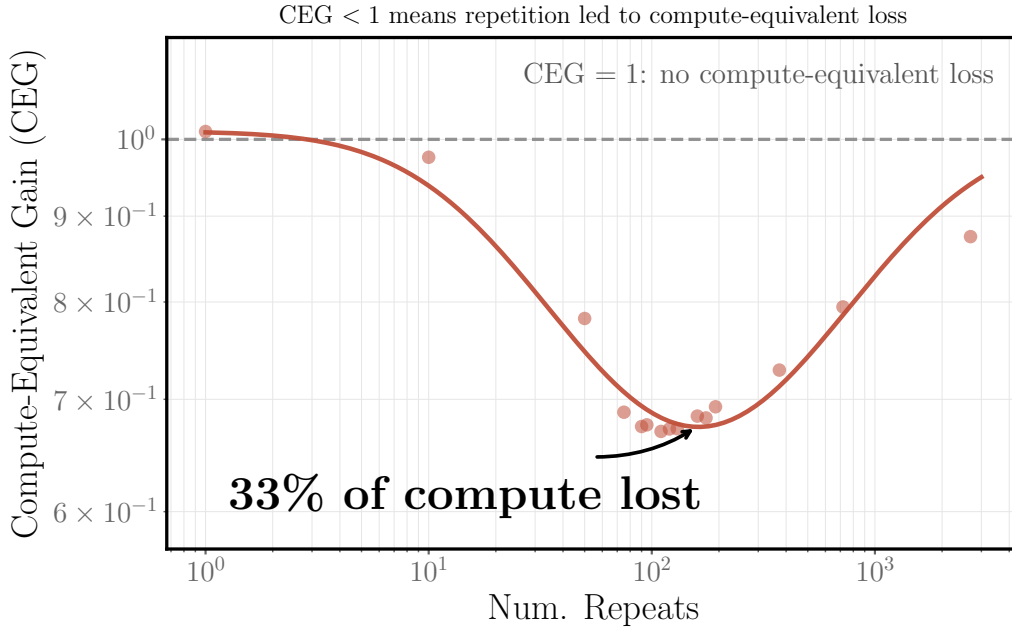
344M model, $OT = 1$ 

Figure 1. Compute-Equivalent Gain (CEG) as a function of the per-document repeat count R , for the Qwen3-style 344M-parameter model trained on FineWeb-Edu-Dedup at the Chinchilla-optimal multiplier $OT = 1$. CEG is the ratio of the no-repetition compute that would reach the achieved loss to the compute actually spent $CEG = 1$ is the no-repetition reference. At R near 100, CEL rises to roughly 0.33, meaning the run reaches the loss of a run without repetitions trained with only two-thirds of the FLOPs (Yang et al., 2025a; Penedo et al., 2024).

most harmful configurations at fixed compute.

Setup. We train Qwen3-style decoder-only transformers (Yang et al., 2025a; Vaswani et al., 2017; Su et al., 2024) on FineWeb-Edu-Dedup (Penedo et al., 2024), using six parameter counts $N \in \{34, 48, 63, 93, 153, 344\}$ M. We define this architecture in more detail in Appendix B. Each run is parameterized by a model size N and an overtraining multiplier $OT \in \{0.25, 0.5, 1, 2, 4\}$ with larger models run on a subset of this grid due to compute constraints. The case $OT = 1$ corresponds to 20 tokens per parameter, following Chinchilla-style scaling (Hoffmann et al., 2022), while $OT > 1$ studies overtraining around that reference point. For a given (N, OT) pair, the total number of training tokens is $T = 20 \cdot OT \cdot N$. Using the standard dense-transformer estimate of $6N$ FLOPs per token (Kaplan et al., 2020; Hoffmann et al., 2022; Sardana et al., 2024; Porian et al., 2024; Gadre et al., 2025), the total training compute is

$$C = 6NT = 120 \cdot OT \cdot N^2. \quad (1)$$

Thus, within each (N, OT) sweep, both the token budget T and compute budget C are fixed.

Repeated-pool construction. For each repeated-data run, $(1 - f)T$ tokens are drawn from non-repeated documents. Let D_r denote the number of unique tokens in the repeated pool, and let R denote the number of times each repeated document is replayed. The repeated-token budget is therefore

$$fT \approx RD_r$$

so

$$D_r \approx \frac{fT}{R} = \frac{2 \cdot OT \cdot N}{R}. \quad (2)$$

Increasing R does not change the amount of repeated material we train on; it concentrates the same 10% repeated-token budget onto a smaller pool. The R view answers how many times repeated documents are replayed, while the D_r view answers how large the repeated corpus is. Both are needed because repetition damage depends on their interaction.

The repeated documents are sampled at document granularity and are disjoint from the non-repeated training documents. Copies of repeated documents are randomly interleaved with the non-repeated stream during training. This setup gives us a controlled setting to study repetition—which is unavoidable in any realistic training corpus—and

identify the configurations that cause the most damage. See Appendix C for more details on the sampling protocol.

Evaluation and baselines. We evaluate every model on a fixed held-out split of approximately 150M tokens, constructed with a fixed train/test seed. This evaluation split is excluded from all training data and from all repeated pools by the fixed train/test split. For each completed (N, OT) sweep, we also train a no-repetition baseline. These baselines serve two purposes. First, the baseline at the same (N, OT) gives the per-sweep reference for measuring the fractional eval-loss increase caused by repetition. Second, the six $OT = 1$ no-repetition baselines calibrate the no-repetition Chinchilla scaling law $L(C) = E + KC^{-\gamma}$ used in §4.3 to convert eval loss into CEG and CEL. Additional implementation and evaluation details are given in Appendix F.

4. Results

We characterize how eval loss depends on the structure of the repeated subset (§4.1), extract a model-size scaling law for the worst-case configuration (§4.2), translate the resulting loss differences into CEG and CEL via the no-repetition Chinchilla scaling law (§4.3), and explain the same non-monotonicity in a closed-form linear-regression model (§4.4–§4.5).

4.1. Eval loss is non-monotonic in the repeat count

For each (N, OT) pair we hold the repeated fraction $f = 0.1$ and the total compute C fixed, varying only R along an iso-FLOP curve. Figure 2 plots eval loss against R . The qualitative pattern is an intermediate-repeat peak with a sharpness that varies between sweeps. Across the completed sweeps, the raw maximum is 1.0 to 4.2% above the corresponding no-repetition baseline, with median 3.1%. Measured relative to the larger of the two endpoints, no repeats and the largest measured R , the peak prominence ranges from 0.7 to 2.7%, with median 1.8%. The peaks are often broad. We therefore use the log-Gaussian fits in §4.2 to estimate peak locations, rather than treating the discrete argmax as exact, because R is sampled on a finite, approximately logarithmic grid and the true maximum may fall between measured repeat counts.

Hernandez et al. (2022) reported a similar non-monotonic dependence in a fixed 100B-token-budget setting. We confirm that the same qualitative intermediate-repeat regime appears under Chinchilla-style budgets, and quantify it using CEG and CEL in §4.3. The implication is that the most damaging repetition structure usually lies away from both extremes of the R range. Configurations with many repeats of a tiny pool, or few repeats of a large pool, generally produce smaller loss increases than configurations in between,

though the recovery at large R is flatter in some sweeps. §4.4 offers a statistical mechanism for why such a peak can arise.

4.2. Peak locations are consistent with a power-law trend in model size

§4.1 established that loss is maximized at some intermediate R . To summarize how this peak shifts with scale, we fit a simple empirical relationship between the estimated peak location and model size. We extract continuous peak locations by fitting a log-Gaussian to eval loss as a function of $\log_{10} R$ for each (N, OT) pair,

$$L(x) = b + A \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right), \quad R^{\text{peak}} = 10^\mu. \quad (3)$$

This fit captures the single peak we observe in log-repeat space. We then fit power laws to the estimated R^{peak} values across the completed (N, OT) grid and convert them to repeated-pool sizes using (2), giving:

$$R^{\text{peak}} = 2.31 \times 10^{10} N^{-0.96}, \quad D_r^{\text{peak}} = 7.58 \times 10^{-10} N^{1.84}, \quad (4)$$

$$R^{\text{peak}} = 1.47 \times 10^7 C^{-0.25}, \quad D_r^{\text{peak}} = 5.49 \times 10^{-12} C^{0.93}. \quad (5)$$

Figure 3 plots these four related views. Because the peak locations are themselves estimated from fitted curves, and because the model sizes span roughly one order of magnitude within a single architecture family, we interpret these fits as within-range empirical summaries rather than validated extrapolative scaling laws.

Observed trend. Within the measured range, larger models tend to reach their largest loss increase at fewer repeats of a larger repeated pool. The Qwen3-style 34M-parameter model peaks at $R \approx 1400$ with $D_r \approx 5 \times 10^4$ tokens, while the Qwen3-style 344M-parameter model peaks at $R \approx 155$ with $D_r \approx 4.5 \times 10^6$ tokens (Yang et al., 2025a). We also do not observe a strong dependence on training duration over the completed grid: the $OT \in \{0.25, 0.5, 1, 2, 4\}$ sweeps largely fall near the same trend in N . Thus, the fitted trend provides a useful heuristic for identifying repetition regimes that may be especially harmful within the studied scale range.

4.3. Repetition can cause an $\mathcal{O}(1)$ Compute-Equivalent Loss

The 2 to 4% loss increases reported in §4.1 and §4.2 translate into different compute gaps depending on where the run sits on the no-repetition Chinchilla curve. Hernandez

Where repetition damage peaks
Dashed curves: log-Gaussian fits per sweep

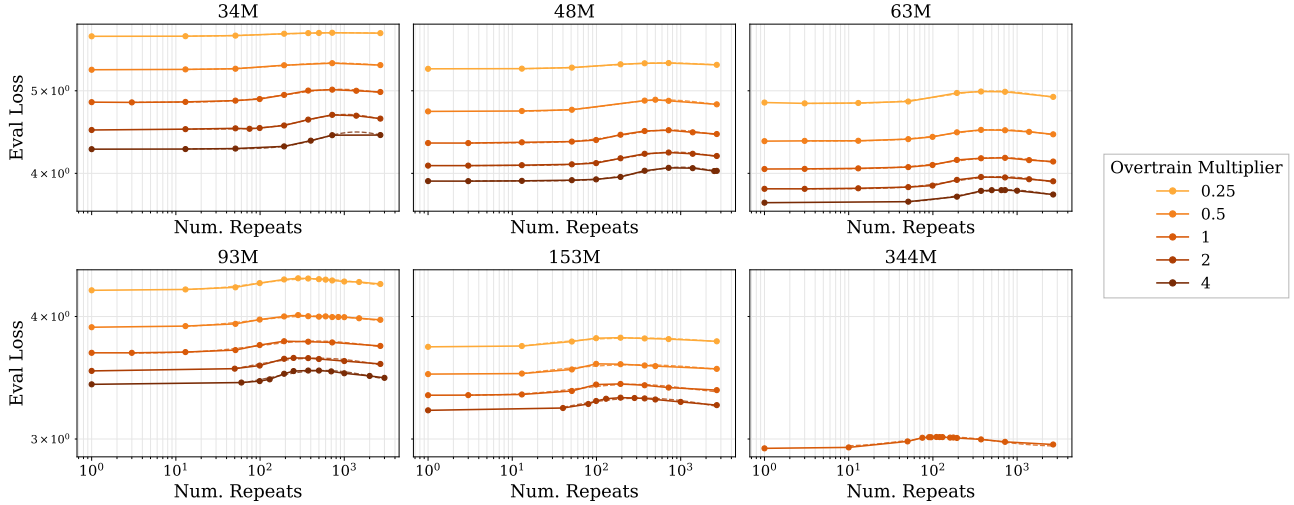


Figure 2. Gaussian fits to eval loss as a function of repeat count. Each panel fixes a model size N , and each curve corresponds to an overtraining multiplier OT . The fitted peak gives R^{peak} , the repeat count at which eval loss is largest for that (N, OT) sweep. Across all model sizes, eval loss is maximized at an intermediate repeat count.

et al. (2022) addressed this by reporting damage as a reduction in effective parameter count, but that framing predates Chinchilla-style scaling and is hard to compare across model sizes and training durations. We replace it with a compute-equivalent ratio. We ask how much compute a no-repetition run would need to match the loss of a repeated-data run.

We fit the no-repetition Chinchilla scaling law $L(C) = E + KC^{-\gamma}$ to the six $OT=1, R=1$ baselines (Hoffmann et al., 2022; Gadre et al., 2025; Porian et al., 2024), obtaining

$$L(C) = 2.365 + 6.647 \times 10^5 C^{-0.317}. \quad (6)$$

The fit is shown in Figure 4. Here E is the irreducible-loss floor that the model approaches in the limit of infinite compute, K sets the overall vertical scale, and $\gamma > 0$ is the rate at which loss decreases with compute. Inverting the law gives $C^*(L) = (K/(L - E))^{1/\gamma}$, the amount of compute a no-repetition run on the law would need to reach a measured loss L . The *Compute-Equivalent Gain* of a repeated-data run is

$$\text{CEG} = \frac{C^*(L)}{C_{\text{actual}}}, \quad \text{CEL} = 1 - \text{CEG}. \quad (7)$$

CEG = 1 matches the no-repetition law. When $\text{CEG} < 1$, the run reaches the loss of a no-repetition run trained with only $\text{CEG} \cdot C_{\text{actual}}$ FLOPs, and $\text{CEL} = 1 - \text{CEG}$ is the Compute-Equivalent Loss.

Figure 5 shows CEG as a function of R . At $OT=1$, the worst repeat settings increase CEL

to 0.19, 0.19, 0.21, 0.21, 0.26, 0.33 for $N \in \{34, 48, 63, 93, 153, 344\}$ M respectively. **At the largest scale, peak repetition produces CEL ≈ 0.33 .** Two patterns are notable. First, the 2 to 4% loss bump translates into CEL values of roughly 0.19 to 0.33 because the no-repetition scaling law is shallow. A small loss gap maps to a large compute gap, and the loss-space view systematically understates the practical cost. Second, varying OT shifts the level of CEG but leaves the location of the peak in R approximately unchanged. The worst-case repetition structure persists across training durations.

Sensitivity to the loss floor. Equation (7) contains the factor $(L - E)^{-1/\gamma}$, which diverges as L approaches the fitted floor E . With $\gamma \approx 0.32$, a 1% shift in $(L - E)$ produces a $\approx 3\%$ relative shift in CEG. The 344M peak run sits at $L - E \approx 0.65$ nats, so the headline 33% number is robust to leave-one-out perturbations of the scaling-law fit. Smaller models live closer to E and inherit larger uncertainty. The scaling law is fit on six points, at the boundary of identifiability for the three-parameter Chinchilla form (Besiroglu et al., 2024; Porian et al., 2024). We therefore treat the absolute CEG and CEL values as point estimates and the qualitative non-monotonicity as the robust finding.

4.4. A statistical model of repetition damage

The non-monotonic peak in eval loss could plausibly be a transformer-specific artifact, a memorization quirk arising from attention, depth, or optimization dynamics. We sug-

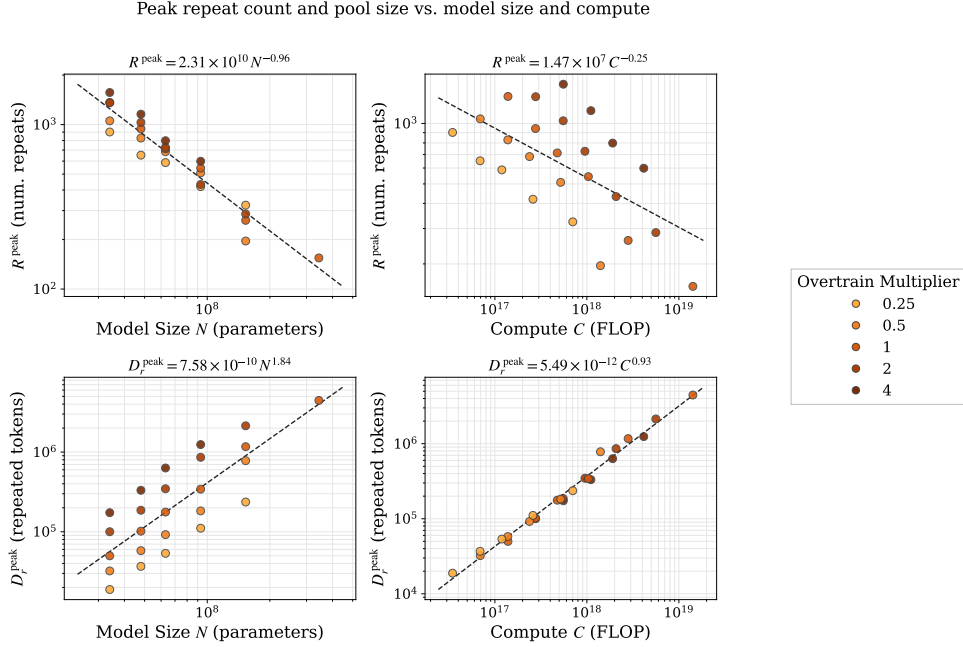


Figure 3. Scaling laws for the peak-damage regime. Top row, the repeat count at peak eval loss R^{peak} decreases with model size (left) and compute (right). Bottom row, the repeated-pool size at peak eval loss D_r^{peak} increases with both. These fits predict the most damaging repetition structure for a given training budget.

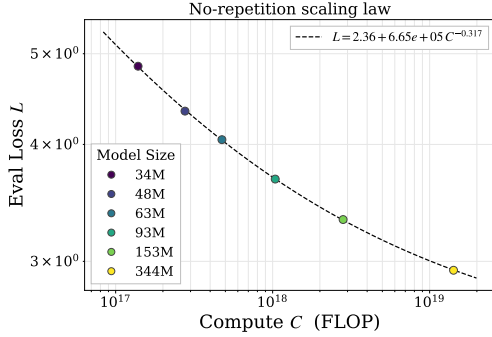


Figure 4. No-repetition scaling law fit. We fit $L(C) = E + KC^{-\gamma}$ to the six OT=1, $R = 1$ baselines. This scaling law converts any eval loss into the equivalent no-repetition compute, enabling the CEG and CEL metrics.

gest that it may reflect a more general statistical effect. A finite-capacity learner trained on duplicated samples from a richer distribution can trade fit on the repeated set against generalization beyond the repeated samples, and this trade-off can produce non-monotonic behavior in the duplication count. We illustrate this possibility in misspecified linear regression, deriving closed-form conditional risks and recovering the same qualitative peak in simulations.

Setup. The data-generating process is a high-dimensional linear model with isotropic Gaussian inputs. Inputs are $x \sim \mathcal{N}(0, I_p)$ in \mathbb{R}^p , with noiseless labels $y = x^\top \beta$ and

a fixed coefficient vector $\beta \in \mathbb{R}^p$. The learner observes only the first $m < p$ coordinates of x , so the model is misspecified. Decompose $x = (x_{\text{in}}, x_{\text{out}})$ and $\beta = (\beta_{\text{in}}, \beta_{\text{out}})$, with $x_{\text{in}}, \beta_{\text{in}} \in \mathbb{R}^m$ observed and $x_{\text{out}}, \beta_{\text{out}} \in \mathbb{R}^{p-m}$ unobserved. The same isotropic distribution governs the test point. The unobserved coordinates carry real predictive signal, so finite-sample correlations between x_{in} and x_{out} can be mistaken for useful structure. The training set contains n unique examples plus a block of d examples each repeated r times, for a total of $n + rd$ rows.

Block-diagonal noise covariance. Let $X_{u,\text{in}}$ and $X_{d,\text{in}}$ stack the observed features for the unique and repeated examples respectively, and let $C_u = X_{u,\text{in}}^\top X_{u,\text{in}}$ and $C_d = X_{d,\text{in}}^\top X_{d,\text{in}}$ be the corresponding observed-feature Gram matrices. The expanded observed-feature Gram of the full training set is $B_r = C_u + rC_d$. The restricted OLS estimator decomposes as $\hat{\beta}_{\text{in}} = \beta_{\text{in}} + a_r$, where a_r is an aliasing term that arises from fitting the unobserved signal through the observed coordinates. The crucial step is that the r copies of each repeated example share a single unobserved-feature realization, so the conditional covariance of $z = X_{\text{out}}\beta_{\text{out}}$ given X_{in} is the block-diagonal direct sum

$$\Sigma_r = I_n \oplus \bigoplus_{i=1}^d \mathbf{1}_r \mathbf{1}_r^\top \in \mathbb{R}^{(n+rd) \times (n+rd)}, \quad (8)$$

where $\mathbf{1}_r \in \mathbb{R}^r$ is the all-ones vector and the symbol \oplus denotes block-diagonal direct sum. The unique block is the

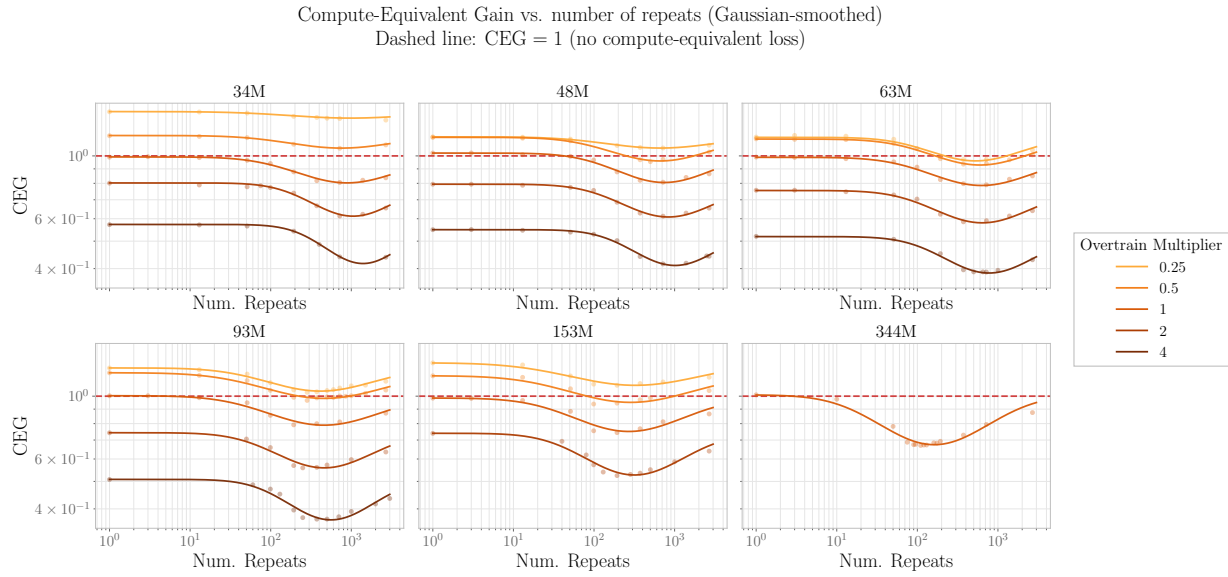


Figure 5. Compute-Equivalent Gain as a function of repeat count, by model size and overtraining multiplier. CEG = 1 (dashed line) matches the no-repetition reference. CEG falls at intermediate repeat counts and partially recovers at large ones, revealing a worst-case repetition structure for each sweep. Damage grows with model size. The Qwen3-style 344M-parameter model reaches CEG ≈ 0.67 , corresponding to CEL ≈ 0.33 .

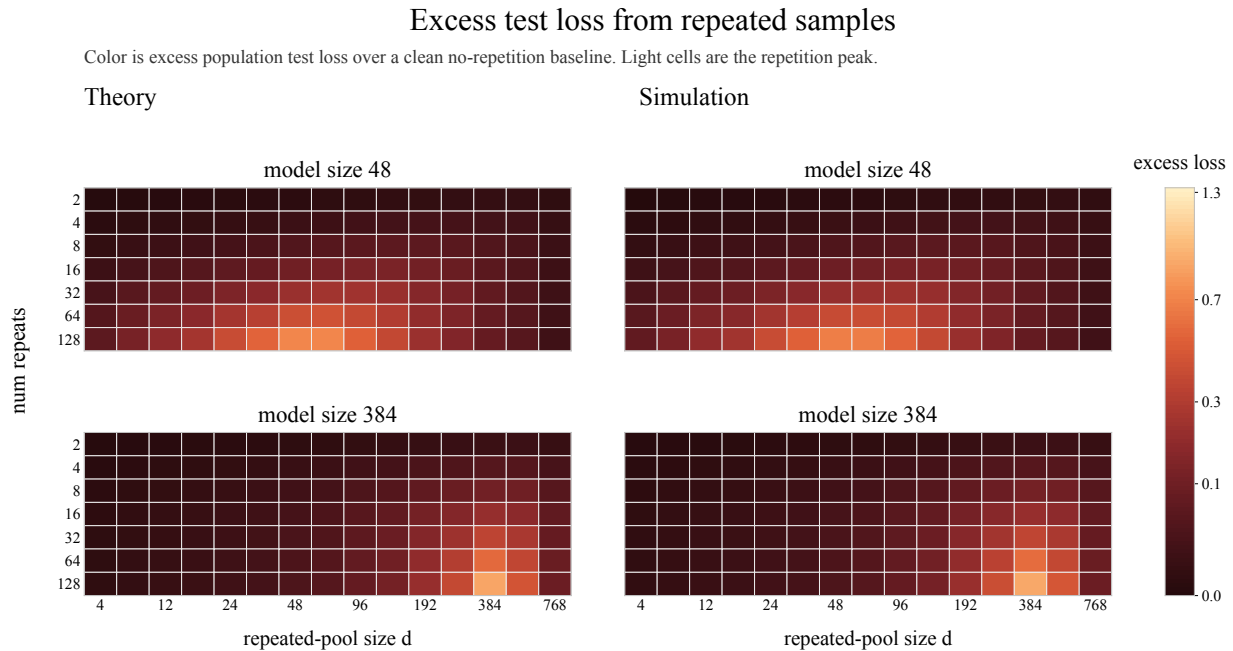


Figure 6. Excess test loss in misspecified linear regression with repeated samples. Rows vary the observed dimension m , and columns compare closed-form theory with direct OLS simulations. Within each panel, the x -axis is the repeated-pool size d , the y -axis is the repeat count r , and color indicates excess population test loss over a no-repetition baseline at the same m . The peak shifts to larger d as m increases, providing a statistical analogue consistent with the $D_r^{\text{peak}}(N)$ trend in Figure 3.

$n \times n$ identity, reflecting independent unobserved features across unique examples. Each repeated block is the rank-one $r \times r$ matrix $\mathbf{1}_r \mathbf{1}_r^\top$, reflecting that the r copies of document i share a single $x_{\text{out},i}$. Distinct repeated documents are uncorrelated. Multiplying through gives $X_{\text{in}}^\top \Sigma_r X_{\text{in}} = C_u + r^2 C_d$. The extra factor of r is what makes duplication qualitatively different from adding more independent samples. Let

$$A_r = C_u + rC_d, \quad G_r = C_u + r^2 C_d.$$

Then, conditioning on X_{in} , the expected train and test errors are

$$\mathbb{E}[L_{\text{train}} | X_{\text{in}}] = \frac{\|\beta_{\text{out}}\|_2^2}{n + rd} [n + rd - \text{tr}(G_r A_r^{-1})], \quad (9)$$

$$\mathbb{E}[L_{\text{test}} | X_{\text{in}}] = \|\beta_{\text{out}}\|_2^2 [1 + \text{tr}(A_r^{-1} G_r A_r^{-1})]. \quad (10)$$

Equations (9)–(10) are closed-form conditional risks. The non-monotonic behavior is a simulation-supported mechanism explored in §4.5. Appendix D gives the derivation.

What the formulas suggest. The analogy to the language-model setting is that m plays the role of model capacity and d corresponds to the unique repeated-data pool D_r . In the regimes we simulate in §4.5, the test loss in (10) is non-monotonic in r at fixed (n, d, m) , as follows. When r is small, the repeated examples carry little extra weight and the predictor is only weakly affected. When r is too large, the repeated block saturates the rank of $C_u + r^2 C_d$ relative to $C_u + rC_d$, and the test loss returns toward a memorize-and-isolate fixed point. The harmful middle regime appears when the repeated pool is both influential and too large to be harmlessly absorbed. The same formulas suggest that increasing m shifts the peak to larger d in the simulated setting, offering a statistical analogue consistent with the empirical signature in (4). They also suggest weak dependence on the unique-sample count n in this toy setting, consistent with the approximate OT-independence observed in §4.2.

4.5. Simulations of the linear model

We validate the closed-form theory with direct simulations. Inputs are sampled as $x \sim \mathcal{N}(0, I_p)$ and labels are generated by $y = x^\top \beta$ with $\|\beta\|_2 = 1$. The learner fits OLS on the first m coordinates of n unique samples plus d samples each repeated r times. For each (m, d, r) we evaluate both (10) and the population test loss of the corresponding OLS solve. We report excess test loss over a no-repetition baseline at the same m . Figure 6 shows that the closed-form and simulation agree to within numerical precision, that excess loss is non-monotonic in the repeated-pool size d at fixed m and r , and that the peak shifts to larger d as m grows. This trend is consistent with the empirical $D_r^{\text{peak}}(N)$ relationship in (4).

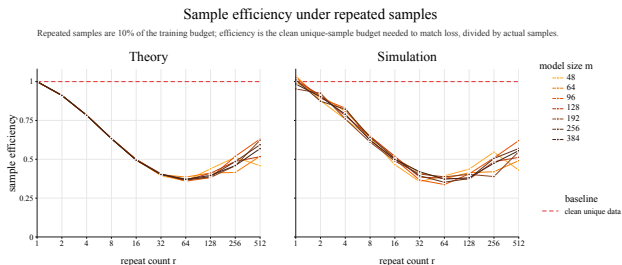


Figure 7. Sample efficiency under repeated samples. The repeated block accounts for 10% of the training budget. For each repeat count r , SE is the unique-sample budget needed to match the repeated-data test loss, divided by the actual sample budget. Theory and simulation both show non-monotonic efficiency loss, mirroring Figure 5.

We also compute a sample-efficiency analogue of CEG. For each repeated-data run, we estimate the no-repetition unique sample budget N_{clean}^* required to match its test loss, and define $\text{SE} = N_{\text{clean}}^*/N_{\text{actual}}$, with the repeated block again accounting for 10% of the training budget. Figure 7 shows that SE falls sharply at intermediate r and partially recovers at extreme r , mirroring the CEG curve in Figure 5 for language models. Figure 8 shows where the worst repeated-pool size occurs for each model size and repeat count. The closed-form linear model thus captures the same qualitative phenomenon, suggesting that the peak is a generic statistical feature of repeated samples in misspecified models.

5. Conclusion

Pretraining is now data-constrained, and this increases risk of residual repetition. We studied exact document-level repetition in a controlled setting, holding the repeated-token fraction fixed at $f = 0.1$ and varying only the repeat count R . At fixed compute, eval loss is worst at an intermediate repeat count: a moderately sized repeated pool replayed many times can hurt more than either a tiny pool replayed many times or a larger pool replayed only a few times. The estimated peak locations follow a clear trend over the model sizes we test: larger models tend to peak at fewer repeats of larger repeated pools. We treat this as an empirical summary of our sweep. Under our fitted no-repetition scaling law, the most damaging repeat setting at our largest scale had $\text{CEG} \approx 0.67$, corresponding to $\text{CEL} \approx 0.33$, meaning it reaches the loss a no-repetition run would reach with about two-thirds of the compute. The linear-regression model gives a simple analogue for this pattern. Our results show that it is not enough to report how much data is duplicated: at the same duplicated-token fraction f , the number of times each duplicate appears can substantially change the compute cost.

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- Abbas, A., Tirumala, K., Simig, D., Ganguli, S., and Morcos, A. S. SemDeDup: Data-efficient learning at web-scale through semantic deduplication. *arXiv preprint arXiv:2303.09540*, 2023. URL <https://arxiv.org/abs/2303.09540>.
- Advani, M. S., Saxe, A. M., and Sompolinsky, H. High-dimensional dynamics of generalization error in neural networks. *Neural Networks*, 132:428–446, 2020. doi: 10.1016/j.neunet.2020.08.022. URL <https://arxiv.org/abs/1710.03667>.
- Albalak, A., Elazar, Y., Xie, S. M., Longpre, S., Lambert, N., Wang, X., Muennighoff, N., Hou, B., Pan, L., Jeong, H., Raffel, C., Chang, S., Hashimoto, T., and Wang, W. Y. A survey on data selection for language models. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=XfHWcNTSHp>. Survey Certification, Featured Certification.
- Allen-Zhu, Z. and Li, Y. Physics of language models: Part 3.3, knowledge capacity scaling laws. In *International Conference on Learning Representations (ICLR)*, 2025. URL <https://openreview.net/forum?id=FxNNiUgtfa>.
- Bahri, Y., Dyer, E., Kaplan, J., Lee, J., and Sharma, U. Explaining neural scaling laws. *Proceedings of the National Academy of Sciences*, 121(27):e2311878121, 2024. doi: 10.1073/pnas.2311878121. URL <https://arxiv.org/abs/2102.06701>.
- Bartlett, P. L., Long, P. M., Lugosi, G., and Tsigler, A. Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences*, 117(48):30063–30070, 2020. doi: 10.1073/pnas.1907378117. URL <https://arxiv.org/abs/1906.11300>.
- Belkin, M., Hsu, D., Ma, S., and Mandal, S. Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2019. doi: 10.1073/pnas.1903070116. URL <https://arxiv.org/abs/1812.11118>.
- Belkin, M., Hsu, D., and Xu, J. Two models of double descent for weak features. *SIAM Journal on Mathematics of Data Science*, 2(4):1167–1180, 2020. doi:

10.1137/20M1336072. URL <https://arxiv.org/abs/1903.07571>.

- Besiroglu, T., Erdil, E., Barnett, M., and You, J. Chinchilla scaling: A replication attempt. *arXiv preprint arXiv:2404.10102*, 2024. URL <https://arxiv.org/abs/2404.10102>.

- Biderman, S., Prashanth, U. S., Sutawika, L., Schoelkopf, H., Anthony, Q., Purohit, S., and Raff, E. Emergent and predictable memorization in large language models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023. URL https://papers.nips.cc/paper_files/paper/2023/hash/59404fb89d6194641c69ae99ecdf8f6d-Abstract-Conference.html.

- Bordelon, B., Atanasov, A., and Pehlevan, C. A dynamical model of neural scaling laws. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 4345–4382, 2024. URL <https://proceedings.mlr.press/v235/bordelon24a.html>.

- Caballero, E., Gupta, K., Rish, I., and Krueger, D. Broken neural scaling laws. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=sckjveqlCZ>.

- 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. In *30th USENIX Security Symposium (USENIX Security 21)*, pp. 2633–2650, 2021. URL <https://www.usenix.org/conference/usenixsecurity21/presentation/carlini-extracting>.

- Carlini, N., Ippolito, D., Jagielski, M., Lee, K., Tramer, F., and Zhang, C. Quantifying memorization across neural language models. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=TatRHT_1cK.

- Davidson, T., Denain, J.-S., Villalobos, P., and Bas, G. AI capabilities can be significantly improved without expensive retraining, 2023. URL <https://arxiv.org/abs/2312.07413>.

- Deng, C., Zhao, Y., Heng, Y., Li, Y., Cao, J., Tang, X., and Cohan, A. Unveiling the spectrum of data contamination in language models: A survey from detection to remediation. In *Findings of the Association for Computational Linguistics: ACL 2024*, pp. 16078–16092, 2024. URL <https://aclanthology.org/2024.findings-acl.951/>.

- 495 Gadre, S. Y., Smyrnis, G., Shankar, V., Gururangan, S.,
496 Wortsman, M., Shao, R., Mercat, J., Fang, A., Li, J., Keh,
497 S., Xin, R., Nezhurina, M., Vasiljevic, I., Soldaini, L.,
498 Jitsev, J., Dimakis, A., Ilharco, G., Koh, P. W., Song,
499 S., Kollar, T., Carmon, Y., Dave, A., Heckel, R., Muen-
500 nighoff, N., and Schmidt, L. Language models scale
501 reliably with over-training and on downstream tasks. In
502 *The Thirteenth International Conference on Learning*
503 *Representations*, 2025. URL [https://openreview.](https://openreview.net/forum?id=iZeQBqJamf)
504 [net/forum?id=iZeQBqJamf](https://openreview.net/forum?id=iZeQBqJamf).
505
- 506 Goyal, S., Maini, P., Lipton, Z. C., Raghunathan, A., and
507 Kolter, J. Z. Scaling laws for data filtering – data
508 curation cannot be compute agnostic. *arXiv preprint*
509 *arXiv:2404.07177*, 2024. URL [https://arxiv.](https://arxiv.org/abs/2404.07177)
510 [org/abs/2404.07177](https://arxiv.org/abs/2404.07177).
- 511 Gundlach, H., Fogelson, A., Lynch, J., Trisovic, A., Rosen-
512 feld, J., Sandhu, A., and Thompson, N. On the ori-
513 gin of algorithmic progress in ai, 2025. URL [https:](https://arxiv.org/abs/2511.21622)
514 [//arxiv.org/abs/2511.21622](https://arxiv.org/abs/2511.21622).
515
- 516 Hastie, T., Montanari, A., Rosset, S., and Tibshirani,
517 R. J. Surprises in high-dimensional ridgeless least
518 squares interpolation. *The Annals of Statistics*, 50
519 (2):949–986, 2022. doi: 10.1214/21-AOS2133.
520 URL [https://projecteuclid.org/](https://projecteuclid.org/journals/annals-of-statistics/volume-50/issue-2/Surprises-in-high-dimensional-ridgeless-least-squares-interpolation/10.1214/21-AOS2133.full)
521 [journals/annals-of-statistics/](https://projecteuclid.org/journals/annals-of-statistics/volume-50/issue-2/Surprises-in-high-dimensional-ridgeless-least-squares-interpolation/10.1214/21-AOS2133.full)
522 [volume-50/issue-2/](https://projecteuclid.org/journals/annals-of-statistics/volume-50/issue-2/Surprises-in-high-dimensional-ridgeless-least-squares-interpolation/10.1214/21-AOS2133.full)
523 [Surprises-in-high-dimensional-ridgeless-least-squares-interpolation/](https://projecteuclid.org/journals/annals-of-statistics/volume-50/issue-2/Surprises-in-high-dimensional-ridgeless-least-squares-interpolation/10.1214/21-AOS2133.full)
524 [10.1214/21-AOS2133.full](https://projecteuclid.org/journals/annals-of-statistics/volume-50/issue-2/Surprises-in-high-dimensional-ridgeless-least-squares-interpolation/10.1214/21-AOS2133.full).
525
- 526 Henighan, T., Kaplan, J., Katz, M., Chen, M., Hesse, C.,
527 Jackson, J., Jun, H., Brown, T. B., Dhariwal, P., Gray, S.,
528 Hallacy, C., Mann, B., Radford, A., Ramesh, A., Ryder,
529 N., Ziegler, D. M., Schulman, J., Amodei, D., and Mc-
530 Candlish, S. Scaling laws for autoregressive generative
531 modeling. *arXiv preprint arXiv:2010.14701*, 2020. URL
532 <https://arxiv.org/abs/2010.14701>.
533
- 534 Hernandez, D., Brown, T., Conerly, T., DasSarma, N.,
535 Drain, D., El-Showk, S., Elhage, N., Hatfield-Dodds,
536 Z., Henighan, T., Hume, T., Johnston, S., Mann, B.,
537 Olah, C., Olsson, C., Amodei, D., Joseph, N., Kap-
538 lan, J., and McCandlish, S. Scaling laws and in-
539 terpretability of learning from repeated data. *arXiv*
540 *preprint arXiv:2205.10487*, 2022. URL [https://](https://arxiv.org/abs/2205.10487)
541 arxiv.org/abs/2205.10487.
542
- 543 Hoffmann, J., Borgeaud, S., Mensch, A., Buchatskaya, E.,
544 Cai, T., Rutherford, E., de Las Casas, D., Hendricks,
545 L. A., Welbl, J., Clark, A., Hennigan, T., Noland,
546 E., Millican, K., van den Driessche, G., Damoc,
547 B., Guy, A., Osindero, S., Simonyan, K., Elsen,
548 E., Rae, J. W., Vinyals, O., and Sifre, L. Training
549 compute-optimal large language models. In *Advances*
in Neural Information Processing Systems (NeurIPS),
2022. URL [https://proceedings.neurips.](https://proceedings.neurips.cc/paper_files/paper/2022/hash/c1e2faff6f588870935f114e04a3e5-Abstract-Conference.html)
[cc/paper_files/paper/2022/hash/](https://proceedings.neurips.cc/paper_files/paper/2022/hash/c1e2faff6f588870935f114e04a3e5-Abstract-Conference.html)
[c1e2faff6f588870935f114e04a3e5-Abstract-Conference](https://proceedings.neurips.cc/paper_files/paper/2022/hash/c1e2faff6f588870935f114e04a3e5-Abstract-Conference.html)
[.html](https://proceedings.neurips.cc/paper_files/paper/2022/hash/c1e2faff6f588870935f114e04a3e5-Abstract-Conference.html).
- Jacovi, A., Caciularu, A., Goldman, O., and Gold-
berg, Y. Stop uploading test data in plain text:
Practical strategies for mitigating data contamination
by evaluation benchmarks. In *Proceedings of the*
*2023 Conference on Empirical Methods in Natu-
ral Language Processing (EMNLP)*, pp. 5075–5084,
2023. URL [https://aclanthology.org/2023.](https://aclanthology.org/2023.emnlp-main.308/)
[emnlp-main.308/](https://aclanthology.org/2023.emnlp-main.308/).
- Kandpal, N., Wallace, E., and Raffel, C. Deduplicating
training data mitigates privacy risks in language mod-
els. In *Proceedings of the 39th International Confer-
ence on Machine Learning*, volume 162 of *Proceed-
ings of Machine Learning Research*, pp. 10697–10707,
2022. URL [https://proceedings.mlr.press/](https://proceedings.mlr.press/v162/kandpal22a.html)
[v162/kandpal22a.html](https://proceedings.mlr.press/v162/kandpal22a.html).
- Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B.,
Chess, B., Child, R., Gray, S., Radford, A., Wu, J., and
Amodei, D. Scaling laws for neural language models.
arXiv preprint arXiv:2001.08361, 2020. URL [https:](https://arxiv.org/abs/2001.08361)
[//arxiv.org/abs/2001.08361](https://arxiv.org/abs/2001.08361).
- Kazdan, J., Levi, N., Schaeffer, R., Chudnovsky, J., Puri,
A., He, B., Donmez, M., Koyejo, S., and Donoho,
D. Scale dependent data duplication. *arXiv preprint*
arXiv:2603.06603, 2026. URL [https://arxiv.](https://arxiv.org/abs/2603.06603)
[org/abs/2603.06603](https://arxiv.org/abs/2603.06603).
- Kingma, D. P. and Ba, J. Adam: A method for stochastic op-
timization, 2017. URL [https://arxiv.org/abs/](https://arxiv.org/abs/1412.6980)
[1412.6980](https://arxiv.org/abs/1412.6980).
- Komatsuzaki, A. One epoch is all you need. *arXiv*
preprint arXiv:1906.06669, 2019. URL [https://](https://arxiv.org/abs/1906.06669)
arxiv.org/abs/1906.06669.
- Lee, K., Ippolito, D., Nystrom, A., Zhang, C., Eck, D.,
Callison-Burch, C., and Carlini, N. Deduplicating train-
ing data makes language models better. In *Proceedings*
*of the 60th Annual Meeting of the Association for Com-
putational Linguistics (ACL)*, 2022. doi: 10.18653/v1/
2022.acl-long.577. URL [https://aclanthology.](https://aclanthology.org/2022.acl-long.577/)
[org/2022.acl-long.577/](https://aclanthology.org/2022.acl-long.577/).
- Lesci, P., Meister, C., Hofmann, T., Vlachos, A., and Pi-
mentel, T. Causal estimation of memorisation profiles.
In *Proceedings of the 62nd Annual Meeting of the Asso-
ciation for Computational Linguistics (ACL)*, pp. 15616–
15635, 2024. URL [https://aclanthology.org/](https://aclanthology.org/2024.acl-long.834/)
[2024.acl-long.834/](https://aclanthology.org/2024.acl-long.834/).

- 550 Li, J., Fang, A., Smyrnis, G., Ivgi, M., Jordan, M., Gadre,
551 S., Bansal, H., Guha, E., Keh, S., Arora, K., Garg,
552 S., Xin, R., Muennighoff, N., Heckel, R., Mercat, J.,
553 Chen, M., Gururangan, S., Wortsman, M., Albalak, A.,
554 Bitton, Y., Nezhurina, M., Abbas, A., Hsieh, C.-Y.,
555 Ghosh, D., Gardner, J., Kilian, M., Zhang, H., Shao, R.,
556 Pratt, S., Sanyal, S., Ilharco, G., Daras, G., Marathe,
557 K., Gokaslan, A., Zhang, J., Chandu, K., Nguyen,
558 T., Vasiljevic, I., Kakade, S., Song, S., Sanghavi, S.,
559 Faghri, F., Oh, S., Zettlemoyer, L., Lo, K., El-Nouby,
560 A., Pouransari, H., Toshev, A., Wang, S., Groeneveld,
561 D., Soldaini, L., Koh, P. W., Jitsev, J., Kollar, T.,
562 Dimakis, A. G., Carmon, Y., Dave, A., Schmidt, L.,
563 and Shankar, V. DataComp-LM: In search of the next
564 generation of training sets for language models. In
565 *Advances in Neural Information Processing Systems*
566 (*NeurIPS*), *Datasets and Benchmarks Track*, 2024. doi:
567 10.52202/079017-0455. URL [https://papers.
568 nips.cc/paper_files/paper/2024/hash/
569 19e4ea30dded58259665db375885e412-Abstract-Dataset-
570 and_Benchmarks_Track.html](https://papers.nips.cc/paper_files/paper/2024/hash/19e4ea30dded58259665db375885e412-Abstract-Dataset-and-Benchmarks_Track.html).
- 571 Longpre, S., Yauney, G., Reif, E., Lee, K., Roberts, A.,
572 Zoph, B., Zhou, D., Wei, J., Robinson, K., Mimno, D.,
573 and Ippolito, D. A pretrainer’s guide to training data:
574 Measuring the effects of data age, domain coverage, qual-
575 ity, and toxicity. In *Proceedings of the 2024 Conference*
576 *of NAACL: Human Language Technologies*, pp. 3245–
577 3276, 2024. URL [https://aclanthology.org/
578 2024.naacl-long.179/](https://aclanthology.org/2024.naacl-long.179/).
- 580 Lovelace, J., Belardi, C., Kundurthy, S., Sudhakar, S., and
581 Weinberger, K. Q. Prescriptive scaling laws for data con-
582 strained training. *arXiv preprint arXiv:2605.01640*, 2026.
583 URL <https://arxiv.org/abs/2605.01640>.
- 584 Magar, I. and Schwartz, R. Data contamination: From
585 memorization to exploitation. In *Proceedings of the 60th*
586 *Annual Meeting of the Association for Computational*
587 *Linguistics (ACL), Volume 2: Short Papers*, pp. 157–165,
588 2022. URL [https://aclanthology.org/2022.
589 acl-short.18/](https://aclanthology.org/2022.acl-short.18/).
- 591 Maini, P., Seto, S., Bai, H., Grangier, D., Zhang, Y., and
592 Jaitly, N. Rephrasing the web: A recipe for compute
593 and data-efficient language modeling. *arXiv preprint*
594 *arXiv:2401.16380*, 2024. URL [https://arxiv.
595 org/abs/2401.16380](https://arxiv.org/abs/2401.16380).
- 597 Marion, M., Üstün, A., Pozzobon, L., Wang, A., Fadaee,
598 M., and Hooker, S. When less is more: Investigat-
599 ing data pruning for pretraining llms at scale. *arXiv*
600 *preprint arXiv:2309.04564*, 2023. URL [https://
601 arxiv.org/abs/2309.04564](https://arxiv.org/abs/2309.04564).
- 602 Mei, S. and Montanari, A. The generalization error of
603 random features regression: Precise asymptotics and the
604 double descent curve. *Communications on Pure and Applied Mathematics*, 75(4):667–766, 2022. doi: 10.1002/cpa.22008. URL [https://arxiv.org/abs/
1908.05355](https://arxiv.org/abs/1908.05355).
- Meta Superintelligence Labs. Introducing Muse Spark: Scaling towards personal superintelligence, April 2026. URL [https://ai.meta.com/blog/
introducing-muse-spark-msl/](https://ai.meta.com/blog/introducing-muse-spark-msl/). Accessed: 2026-05-06.
- Muennighoff, N., Rush, A. M., Barak, B., Scao, T. L., Tazi, N., Piktus, A., Pyysalo, S., Wolf, T., and Raffel, C. Scaling data-constrained language models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL [https://openreview.net/forum?
id=j5BuTrEj35](https://openreview.net/forum?id=j5BuTrEj35).
- Nakkiran, P., Kaplun, G., Bansal, Y., Yang, T., Barak, B., and Sutskever, I. Deep double descent: Where bigger models and more data hurt. In *International Conference on Learning Representations (ICLR)*, 2020. URL [https://openreview.net/forum?
id=Blg5sA4twr](https://openreview.net/forum?id=Blg5sA4twr).
- Oren, Y., Meister, N., Chatterji, N., Ladhak, F., and Hashimoto, T. B. Proving test set contamination in black-box language models. In *The Twelfth International Conference on Learning Representations (ICLR)*, 2024. URL [https://openreview.net/forum?
id=KS8mIvetg2](https://openreview.net/forum?id=KS8mIvetg2).
- Penedo, G., Malartic, Q., Hesslow, D., Cojocar, R., Alobeidli, H., Cappelli, A., Pannier, B., Almazrouei, E., and Lounay, J. The RefinedWeb dataset for Falcon LLM: Outperforming curated corpora with web data only. In *Advances in Neural Information Processing Systems (NeurIPS), Datasets and Benchmarks Track*, 2023. URL [https://papers.nips.cc/paper/2023/hash/
fa3ed726cc5073b9c31e3e49a807789c-Abstract-Dataset-
and-Benchmarks.html](https://papers.nips.cc/paper/2023/hash/fa3ed726cc5073b9c31e3e49a807789c-Abstract-Dataset-and-Benchmarks.html).
- Penedo, G., Kydlíček, H., Ben Allal, L., Lozhkov, A., Mitchell, M., Raffel, C., Von Werra, L., and Wolf, T. The FineWeb datasets: Decanting the web for the finest text data at scale. In *Advances in Neural Information Processing Systems (NeurIPS), Datasets and Benchmarks Track*, 2024. doi: 10.52202/079017-0970. URL [https://papers.
nips.cc/paper_files/paper/2024/hash/
370df50ccfd8bde18f8f9c2d9151bda-Abstract-Dataset-
and-Benchmarks_Track.html](https://papers.nips.cc/paper_files/paper/2024/hash/370df50ccfd8bde18f8f9c2d9151bda-Abstract-Dataset-and-Benchmarks_Track.html).
- Porian, T., Wortsman, M., Jitsev, J., Schmidt, L., and Carmon, Y. Resolving discrepancies in compute-optimal scaling of language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*,

2024. URL <https://openreview.net/forum?id=4fSSqpk1sM>.
- Sardana, N., Portes, J., Doubov, S., and Frankle, J. Beyond chinchilla-optimal: Accounting for inference in language model scaling laws. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=0bmXrtTDUu>.
- Schaeffer, R., Kazdan, J., Abbasi, B., Liu, K. Z., Miranda, B., Ahmed, A., Berez, F., Puri, A., Biderman, S., Miresghallah, N., and Koyejo, S. Quantifying the effect of test set contamination on generative evaluations. *arXiv preprint arXiv:2601.04301*, 2026. URL <https://arxiv.org/abs/2601.04301>.
- Soldaini, L., Kinney, R., Bhagia, A., Schwenk, D., Atkinson, D., Authur, R., Bogin, B., Chandu, K., Dumas, J., Elazar, Y., Hofmann, V., Jha, A. H., Kumar, S., Lucy, L., Lyu, X., Lambert, N., Magnusson, I., Morrison, J., Muennighoff, N., Naik, A., Nam, C., Peters, M. E., Ravichander, A., Richardson, K., Shen, Z., Strubell, E., Subramani, N., Taffjord, O., Walsh, P., Zettlemoyer, L., Smith, N. A., Hajishirzi, H., Beltagy, I., Groeneveld, D., Dodge, J., and Lo, K. Dolma: An open corpus of three trillion tokens for language model pretraining research. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL)*, pp. 15725–15788, 2024. URL <https://aclanthology.org/2024.acl-long.840/>.
- Sorscher, B., Geirhos, R., Shekhar, S., Ganguli, S., and Morcos, A. S. Beyond neural scaling laws: Beating power law scaling via data pruning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/hash/7b75da9b61eda40fa35453ee5d077df6-Abstract-Conference.html.
- Su, J., Lu, Y., Pan, S., Murtadha, A., Wen, B., and Liu, Y. RoFormer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024. doi: 10.1016/j.neucom.2023.127063. URL <https://doi.org/10.1016/j.neucom.2023.127063>.
- Tay, Y., Dehghani, M., Abnar, S., Chung, H. W., Fedus, W., Rao, J., Narang, S., Tran, V. Q., Yogatama, D., and Metzler, D. Scaling laws vs model architectures: How does inductive bias influence scaling? *arXiv preprint arXiv:2207.10551*, 2022. URL <https://arxiv.org/abs/2207.10551>.
- Tirumala, K., Markosyan, A. H., Zettlemoyer, L., and Aghajanyan, A. Memorization without overfitting: Analyzing the training dynamics of large language models. In *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 38274–38290, 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/hash/fa0509f4dab6807e2cb465715bf2d249-Abstract-Conference.html.
- Tirumala, K., Simig, D., Aghajanyan, A., and Morcos, A. S. D4: Improving llm pretraining via document de-duplication and diversification. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023. URL https://papers.nips.cc/paper_files/paper/2023/hash/a8f8cbd7f7a5fb2c837e578c75e5b615-Abstract-Dataset-and-Benchmarks.html.
- Tsigler, A. and Bartlett, P. L. Benign overfitting in ridge regression. *Journal of Machine Learning Research*, 24(123):1–76, 2023. URL <https://jmlr.org/papers/v24/22-1398.html>.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 5998–6008, 2017. URL <https://arxiv.org/abs/1706.03762>.
- Villalobos, P., Ho, A., Sevilla, J., Besiroglu, T., Heim, L., and Hobbhahn, M. Position: Will we run out of data? limits of LLM scaling based on human-generated data. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=ViZcgDQjyG>.
- Weber, M., Fu, D. Y., Anthony, Q., Oren, Y., Adams, S., Alexandrov, A., Lyu, X., Nguyen, H., Yao, X., Adams, V., Athiwaratkun, B., Chalamala, R., Chen, K., Ryabinin, M., Dao, T., Liang, P., Ré, C., Rish, I., and Zhang, C. RedPajama: An open dataset for training large language models. In *Advances in Neural Information Processing Systems (NeurIPS), Datasets and Benchmarks Track*, 2024. doi: 10.52202/079017-3697. URL https://proceedings.neurips.cc/paper_files/paper/2024/hash/d34497330b1fd6530f7afd86d0df9f76-Abstract-Dataset-and-Benchmarks_Track.html.
- Xie, S. M., Pham, H., Dong, X., Du, N., Liu, H., Lu, Y., Liang, P., Le, Q. V., Ma, T., and Yu, A. W. DoReMi: Optimizing data mixtures speeds up language model pretraining. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023. URL https://papers.nips.cc/paper_files/paper/2023/hash/dcba6be91359358c2355cd920da3fcbd-Abstract-Conference.html.
- Xue, F., Fu, Y., Zhou, W., Zheng, Z., and You, Y. To repeat or not to repeat: Insights from scaling

660 llm under token-crisis. In *Advances in Neural*
661 *Information Processing Systems (NeurIPS)*, 2023.
662 URL [https://proceedings.neurips.](https://proceedings.neurips.cc/paper_files/paper/2023/hash/b9e472cd579c83e2f6aa3459f46aac28-Abstract-Conference.html)
663 [cc/paper_files/paper/2023/hash/](https://proceedings.neurips.cc/paper_files/paper/2023/hash/b9e472cd579c83e2f6aa3459f46aac28-Abstract-Conference.html)
664 [b9e472cd579c83e2f6aa3459f46aac28-Abstract-Conference.](https://proceedings.neurips.cc/paper_files/paper/2023/hash/b9e472cd579c83e2f6aa3459f46aac28-Abstract-Conference.html)
665 [html](https://proceedings.neurips.cc/paper_files/paper/2023/hash/b9e472cd579c83e2f6aa3459f46aac28-Abstract-Conference.html).

666 Yang, A., Yang, B., Hui, B., Zheng, B., Yu, B., Zhou,
667 C., et al. Qwen2 technical report. *arXiv preprint*
668 *arXiv:2407.10671*, 2024. URL [https://arxiv.](https://arxiv.org/abs/2407.10671)
669 [org/abs/2407.10671](https://arxiv.org/abs/2407.10671).

671 Yang, A., Li, A., Yang, B., Zhang, B., Hui, B., Zheng, B.,
672 Yu, B., Gao, C., Huang, C., Lv, C., Zheng, C., Liu, D.,
673 Zhou, F., Huang, F., Hu, F., Ge, H., Wei, H., Lin, H., Tang,
674 J., Yang, J., Tu, J., Zhang, J., Yang, J., Yang, J., Zhou,
675 J., Zhou, J., Lin, J., Dang, K., Bao, K., Yang, K., Yu, L.,
676 Deng, L., Li, M., Xue, M., Li, M., Zhang, P., Wang, P.,
677 Zhu, Q., Men, R., Gao, R., Liu, S., Luo, S., Li, T., Tang,
678 T., Yin, W., Ren, X., Wang, X., Zhang, X., Ren, X., Fan,
679 Y., Su, Y., Zhang, Y., Zhang, Y., Wan, Y., Liu, Y., Wang,
680 Z., Cui, Z., Zhang, Z., Zhou, Z., and Qiu, Z. Qwen3
681 technical report. *arXiv preprint arXiv:2505.09388*, 2025a.
682 URL <https://arxiv.org/abs/2505.09388>.

684 Yang, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu,
685 B., et al. Qwen2.5 technical report. *arXiv preprint*
686 *arXiv:2412.15115*, 2025b. URL [https://arxiv.](https://arxiv.org/abs/2412.15115)
687 [org/abs/2412.15115](https://arxiv.org/abs/2412.15115).

688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714

A. Comprehensive related work

We expand here the discussion sketched in §2, organized along five threads.

Repeated data in language model pretraining. Our closest predecessor is [Hernandez et al. \(2022\)](#), who train transformers with a small fraction of repeated data and observe a non-monotonic test-loss curve. They frame the damage as a reduction in effective parameter count and connect it to mechanistic changes in induction heads ([Biderman et al., 2023](#); [Allen-Zhu & Li, 2025](#)). The complementary regime, in which the entire corpus is uniformly repeated, is studied by [Muennighoff et al. \(2023\)](#). They find that up to roughly four epochs are nearly as useful as fresh data and that an additive overfitting term in the Chinchilla loss captures the rest. [Xue et al. \(2023\)](#) reach a similar four-epoch threshold from a different angle. Other work in this family includes [Komatsuzaki \(2019\)](#), who studies one-pass-vs-multi-pass tradeoffs at smaller scale, and [Maini et al. \(2024\)](#), who replaces literal repeats with paraphrases generated by a teacher model. [Tirumala et al. \(2022\)](#) and [Lesci et al. \(2024\)](#) study how memorization grows with the number of times an example is seen during training. [Kazdan et al. \(2026\)](#) argue that semantic duplication is itself scale-dependent, in the sense that larger models recognize more documents as duplicates. [Lovelace et al. \(2026\)](#) fit a one-parameter overfitting penalty within Chinchilla scaling that is closely related in spirit to our CEG/CEL metric. Our setting differs from all of these in three ways. First, we hold a fixed minority fraction $f = 0.1$ of training tokens repeated, mimicking the residue of imperfect deduplication and complementing the all-or-nothing regimes of [Muennighoff et al. \(2023\)](#). Second, we measure damage in compute-equivalent units derived from a fitted Chinchilla scaling law, sharpening the effective-parameter view of [Hernandez et al. \(2022\)](#) into a quantity practitioners directly allocate. Third, we extract a closed-form scaling law for the worst-case configuration as a function of model size, which earlier work has left open.

Compute-optimal scaling. The power-law form for transformer loss was established by [Kaplan et al. \(2020\)](#) and refined by [Hoffmann et al. \(2022\)](#), whose Chinchilla fit is the basis for our no-repetition reference curve. Subsequent work has tested the robustness and generalizability of the Chinchilla functional form. [Besiroglu et al. \(2024\)](#) reanalyze the original Chinchilla fits, [Porian et al. \(2024\)](#) reconcile competing exponents across studies, [Sardana et al. \(2024\)](#) extend the framework to inference-aware budgets, and [Gadre et al. \(2025\)](#) show that Chinchilla-style fits extrapolate reliably under aggressive over-training in the regime we use. [Caballero et al. \(2023\)](#); [Bahri et al. \(2024\)](#); [Bordelon et al. \(2024\)](#); [Henighan et al. \(2020\)](#) provide alternative parametric forms and theoretical accounts. [Tay et al. \(2022\)](#) document the model-family dependence of fitted exponents, supporting our caveat in the Limitations that exponents may shift across architecture families.

Deduplication, memorization, and benchmark contamination. Deduplication has been a central tool of pretraining-corpus construction since at least [Lee et al. \(2022\)](#), who show that exact and near-duplicate removal improves training efficiency and reduces verbatim regurgitation. [Kandpal et al. \(2022\)](#) document a superlinear amplification of privacy risk by duplicates, and [Carlini et al. \(2023; 2021\)](#) relate memorization to model scale and duplicate count. Semantic deduplication was advanced by [Abbas et al. \(2023\)](#) and [Tirumala et al. \(2023\)](#), who use embedding clusters to remove near-duplicate documents that elude exact-hashing pipelines. Aggregate surveys include [Deng et al. \(2024\)](#) and the data-selection survey of [Albalak et al. \(2024\)](#). A related but distinct concern is benchmark contamination, where evaluation data leaks into training ([Deng et al., 2024](#); [Schaeffer et al., 2026](#); [Oren et al., 2024](#); [Magar & Schwartz, 2022](#); [Jacovi et al., 2023](#)). [Schaeffer et al. \(2026\)](#) show that even a single test-set replica can drive measured loss below the no-contamination irreducible floor. We deliberately exclude evaluation documents from training and repeated pools by the fixed train/test split, so the harm we observe comes from a distorted effective training distribution.

Pretraining corpora and data selection. Our experiments use FineWeb-Edu-Dedup ([Penedo et al., 2024](#)). Closely related curated web corpora include DataComp-LM ([Li et al., 2024](#)), Dolma ([Soldaini et al., 2024](#)), RedPajama-v2 ([Weber et al., 2024](#)), and RefinedWeb ([Penedo et al., 2023](#)). [Longpre et al. \(2024\)](#) survey the broader pretrainer’s-data landscape. On the data-selection side, [Sorscher et al. \(2022\)](#) show that careful pruning can break power-law scaling, [Marion et al. \(2023\)](#) study perplexity-based pruning, [Xie et al. \(2023\)](#) optimize domain mixtures, and [Goyal et al. \(2024\)](#) extend scaling laws to data-quality interventions. Our results complement this view by studying a *negative* form of data selection, asking which residual repetition structures to avoid. The two perspectives are quantitatively connected through the same Chinchilla scaling law, since both ultimately translate dataset choices into a position on a no-intervention scaling curve.

Statistical accounts of overfitting. The closed-form analysis in §4.4 sits in the literature on benign overfitting and double descent ([Belkin et al., 2019](#); [Nakkiran et al., 2020](#); [Hastie et al., 2022](#); [Bartlett et al., 2020](#); [Belkin et al., 2020](#); [Advani et al.,](#)

Internal Data Repetition Destroys Language Models

Table 1. Qwen3-style model configurations used in the experiments.

| Model | Layers | d_{model} | d_{ff} | Heads | KV heads | d_{head} | Train ctx. | Vocab | Non-emb. params | Total params |
|-------|--------|--------------------|-----------------|-------|----------|-------------------|------------|--------|-----------------|--------------|
| 34M | 3 | 96 | 256 | 32 | 32 | 128 | 2048 | 151670 | 4,941,216 | 34,061,856 |
| 48M | 4 | 128 | 512 | 32 | 32 | 128 | 2048 | 151670 | 9,177,216 | 48,004,736 |
| 63M | 5 | 160 | 512 | 32 | 32 | 128 | 2048 | 151670 | 14,339,040 | 62,873,440 |
| 93M | 6 | 224 | 768 | 32 | 32 | 128 | 2048 | 151670 | 25,121,120 | 93,069,280 |
| 153M | 9 | 320 | 1024 | 32 | 32 | 128 | 2048 | 151670 | 56,041,664 | 153,110,464 |
| 344M | 14 | 576 | 1536 | 32 | 32 | 128 | 2048 | 151670 | 169,299,776 | 344,023,616 |

2020; Tsigler & Bartlett, 2023; Mei & Montanari, 2022), which studies non-monotonic generalization in over-parameterized linear and kernel models. Most of this literature varies model dimension, sample size, or interpolation. Literal verbatim duplication of a subset of observations and the resulting block-diagonal noise covariance (8) are, to our knowledge, novel in this setting. Our derivation isolates the harm from duplication itself, controlling for the change in dataset size that would otherwise confound the effect.

B. Architecture.

We instantiate all models from scratch using the Qwen3 decoder architecture family (Yang et al., 2025a). Across model sizes, we vary depth, hidden width, and feed-forward width, while holding the remaining architectural settings fixed. All models use rotary position embeddings (Su et al., 2024), RMSNorm, SwiGLU feed-forward layers, grouped-query attention, untied input/output embeddings, BF16 training, and FlashAttention-2. The training sequence length is 2048 tokens. The maximum position length is 32768, the vocabulary size is 151670, the attention head dimension is 128, and all models use 32 attention heads and 32 key-value heads. Table 1 reports the exact configurations and parameter counts. Non-embedding parameters exclude both the token embedding matrix and the untied output LM head; total parameters include both.

C. Repeated-pool construction details

For each run, we first split FineWeb-Edu-Dedup into training and held-out evaluation documents using train/test split seed 0. The evaluation split is constructed before any repeated pools are selected, so evaluation documents are excluded from both the repeated and non-repeated training streams.

Documents are tokenized with the Qwen3 (Yang et al., 2025a) tokenizer, truncated to the training sequence length, and assigned EOS tokens before token counts are computed. Let T be the target number of training tokens for a run, $f = 0.1$ the target repeated-token fraction, and R the number of times each repeated document is replayed. The target unique repeated-pool size is

$$D_r^* = \frac{fT}{R},$$

and the target non-repeated budget is $(1 - f)T$.

Documents are selected without replacement from the training split. We form a seeded random ordering of training documents using the run’s shuffle seed. The repeated pool is the first prefix of this ordering whose cumulative token count reaches D_r^* . The non-repeated pool is then selected from the immediately following documents until its cumulative token count reaches $(1 - f)T$. Thus selection is uniform over documents through a seeded shuffle, not token-weighted sampling, and the repeated documents are disjoint from the non-repeated documents.

Because cutoffs occur at document boundaries, the realized unique repeated-pool size \hat{D}_r and realized repeated-token fraction \hat{f} are approximate. If L_{max} is the maximum tokenized document length after truncation and EOS insertion, then

$$D_r^* \leq \hat{D}_r < D_r^* + L_{\text{max}}.$$

The non-repeated pool has the same one-document overshoot bound. In our preprocessing, $L_{\text{max}} \leq 2049$ because documents are truncated to length 2048 and an EOS token may be appended after truncation. Each document in the repeated pool is then inserted exactly R times, each non-repeated document is inserted once, and the resulting document-index list is shuffled before training.

Throughout the paper, $D_r = fT/R$ denotes the target unique repeated-pool size. The approximation $fT \approx RD_r$ reflects this document-boundary rounding.

D. Linear regression with repeated samples

We derive the formulas used in Section 4.4. Let the original training set contain n unique examples and d repeatable examples. The repeatable block is duplicated r times, so the expanded training set has

$$N = n + rd$$

rows. The learner observes only the first m coordinates of each input. Write the expanded observed-feature matrix as X_{in} and the unobserved-feature matrix as X_{out} . Labels are noiseless:

$$y = X_{\text{in}}\beta_{\text{in}} + X_{\text{out}}\beta_{\text{out}}.$$

The restricted OLS estimator is

$$\hat{\beta}_{\text{in}} = (X_{\text{in}}^\top X_{\text{in}})^{-1} X_{\text{in}}^\top y = \beta_{\text{in}} + a_r,$$

where

$$a_r = (X_{\text{in}}^\top X_{\text{in}})^{-1} X_{\text{in}}^\top X_{\text{out}}\beta_{\text{out}}.$$

Thus a_r is the aliasing term caused by fitting the unobserved part of the signal using the observed coordinates.

Let

$$C_u = X_{u,\text{in}}^\top X_{u,\text{in}}, \quad C_d = X_{d,\text{in}}^\top X_{d,\text{in}}.$$

Then

$$X_{\text{in}}^\top X_{\text{in}} = C_u + rC_d.$$

Now condition on X_{in} and take expectation over the unobserved features. Let

$$z = X_{\text{out}}\beta_{\text{out}}.$$

Because repeated rows share the same unobserved coordinates, the covariance of z over the expanded training set is not proportional to the identity. Instead,

$$\mathbb{E}[zz^\top | X_{\text{in}}] = \|\beta_{\text{out}}\|_2^2 \Sigma_r, \quad \Sigma_r = I_n \oplus \bigoplus_{i=1}^d \mathbf{1}_r \mathbf{1}_r^\top,$$

where $\mathbf{1}_r \in \mathbb{R}^r$ is the all-ones vector and the i -th repeated block is the rank-one $r \times r$ matrix $\mathbf{1}_r \mathbf{1}_r^\top$. The unique block contributes the $n \times n$ identity (independent unobserved features); each repeated block has every entry equal to every other (the r copies of document i share a single $x_{\text{out},i}$). The full matrix is $(n + rd) \times (n + rd)$ and block-diagonal; in particular, distinct repeated documents are uncorrelated. Therefore,

$$X_{\text{in}}^\top \Sigma_r X_{\text{in}} = C_u + r^2 C_d.$$

For training loss, let

$$H = X_{\text{in}} (X_{\text{in}}^\top X_{\text{in}})^{-1} X_{\text{in}}^\top$$

be the projection matrix onto the observed-feature span. The residual is $(I - H)z$, so

$$\mathbb{E}[L_{\text{train}} | X_{\text{in}}] = \frac{1}{N} \mathbb{E}[z^\top (I - H)z | X_{\text{in}}].$$

Using the covariance above,

$$\mathbb{E}[L_{\text{train}} | X_{\text{in}}] = \frac{\|\beta_{\text{out}}\|_2^2}{N} \text{tr}((I - H)\Sigma_r).$$

Repeated-pool size at peak excess test loss

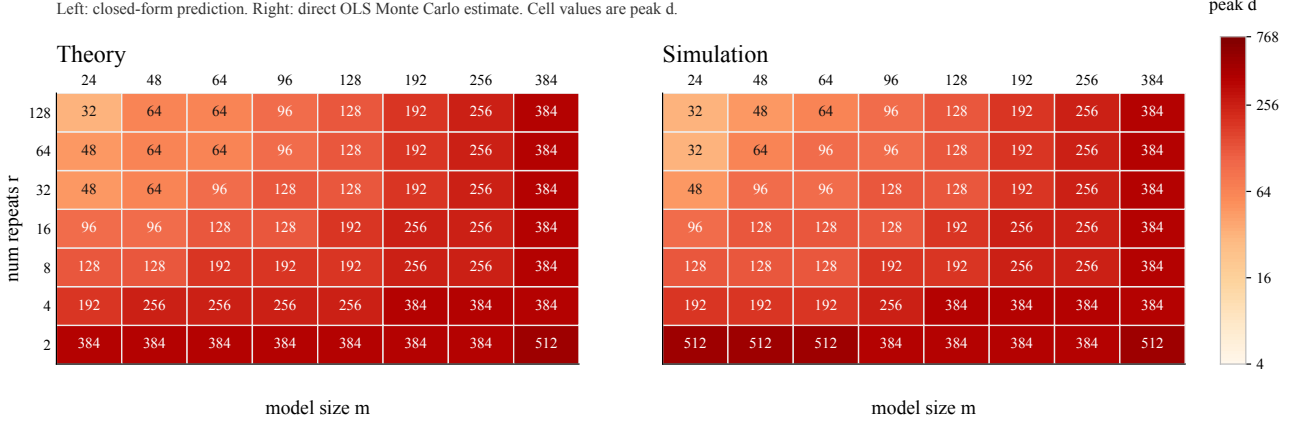


Figure 8. To isolate the peak location, we fix (m, r) and report the repeated-pool size d that maximizes excess test loss. The same trend appears in both the closed-form risk and direct OLS simulations: larger-capacity models peak at larger repeated pools, while higher repeat counts shift the peak toward smaller pools. We note that theory and simulation agree up to the resolution of the d grid.

Since $\text{tr}(\Sigma_r) = N$ and

$$\text{tr}(H\Sigma_r) = \text{tr}((C_u + r^2C_d)(C_u + rC_d)^{-1}),$$

we obtain

$$\mathbb{E}[L_{\text{train}} | X_{\text{in}}] = \frac{\|\beta_{\text{out}}\|_2^2}{N} [N - \text{tr}((C_u + r^2C_d)(C_u + rC_d)^{-1})].$$

For test loss, take a fresh example $(x_{\text{in}}, x_{\text{out}}) \sim \mathcal{N}(0, I_p)$ from the population, with x_{in} and x_{out} independent under the isotropic assumption. The prediction error is

$$x_{\text{out}}^\top \beta_{\text{out}} - x_{\text{in}}^\top a_r.$$

The two terms are independent and mean zero (the cross term vanishes by independence and zero mean of x_{in}), so

$$\mathbb{E}[L_{\text{test}} | X_{\text{in}}] = \|\beta_{\text{out}}\|_2^2 + \mathbb{E}[\|a_r\|_2^2 | X_{\text{in}}].$$

Substituting the expression for a_r gives

$$\mathbb{E}[\|a_r\|_2^2 | X_{\text{in}}] = \|\beta_{\text{out}}\|_2^2 \text{tr}((C_u + rC_d)^{-1}(C_u + r^2C_d)(C_u + rC_d)^{-1}),$$

and therefore

$$\mathbb{E}[L_{\text{test}} | X_{\text{in}}] = \|\beta_{\text{out}}\|_2^2 [1 + \text{tr}((C_u + rC_d)^{-1}(C_u + r^2C_d)(C_u + rC_d)^{-1})].$$

E. Theory: Repetition Peak Heatmap

F. Training and evaluation details

Architecture. We use Qwen3-style decoder-only transformers (Yang et al., 2025a; 2024; 2025b) with rotary position embeddings (Su et al., 2024), RMSNorm, SwiGLU feed-forward layers, grouped-query attention, and the Qwen3 tokenizer. Six parameter counts are used, $N \in \{34, 48, 63, 93, 153, 344\}$ M, obtained by scaling depth and width approximately uniformly. We compute $C = 6NT$ FLOPs per the standard dense-transformer estimate (Kaplan et al., 2020; Hoffmann et al., 2022).

Optimizer and schedule. All runs use AdamW (Kingma & Ba, 2017) with the fused PyTorch implementation (adamw.torch.fused), $\beta_1 = 0.9$, $\beta_2 = 0.95$, weight decay 0.01, and gradient clipping at 1.0. The learning rate follows a cosine schedule with warmup ratio 0.2. The peak learning rate is derived from a base learning rate of 10^{-6} and the computed optimizer-step token count, rather than tuned separately per model size. Sequence length is 2048 throughout. Training uses BF16 weights, FlashAttention-2, and torch compilation.

Data and evaluation. The source corpus is FineWeb-Edu-Dedup from the HuggingFaceTB SmolLM corpus [Penedo et al. \(2024\)](#). We split the corpus once with train/test split seed 0, holding out approximately 150M tokens for evaluation before constructing any repeated pools. Thus evaluation documents are excluded from both the non-repeated training stream and the repeated pool. Documents are tokenized with the Qwen3 tokenizer ([Yang et al., 2025a](#)), truncated to length 2048, and assigned EOS tokens. For each repeated-data run, the repeated pool is selected at document granularity to contribute approximately fT/R unique tokens. These documents are then replayed R times, combined with non-repeated documents contributing approximately $(1 - f)T$ tokens, and shuffled into the final training stream. This makes the repeated-token fraction approximately $f = 0.1$ while varying the concentration of those repeated tokens.

Sweep grid. We train no-repetition baselines and repeated-data sweeps for each completed (N, OT) cell. The completed analysis grid contains 25 cells: all five overtraining multipliers for 34M, 48M, 63M, and 93M; four multipliers through $OT = 2$ for 153M; and $OT = 1$ for 344M. Repeat counts are swept on an irregular, approximately logarithmic grid spanning from $R = 1$ to as high as $R = 20000$, subject to the repeated pool containing at least one document. Each completed cell is a single training run.

Peak fitting. For each completed (N, OT) sweep, we first compute the fractional eval-loss increase relative to the corresponding no-repetition baseline, $\eta(R) = (L(R) - L_{\text{base}})/L_{\text{base}}$. We then fit a three-parameter Gaussian in $\log_{10} R$ to the non-baseline points ($R > 1$):

$$\eta(R) = A \exp\left(-\frac{(\log_{10} R - \log_{10} \mu)^2}{2\sigma^2}\right).$$

The fitted peak repeat count is $R^{\text{peak}} = \mu$, and the corresponding repeated-pool size is computed from $D_{\text{r}}^{\text{peak}} = 2OTN/R^{\text{peak}}$. Power-law fits in (4)–(5) are ordinary least-squares regressions in log-log space over the completed sweeps.

Scaling-law fitting. The no-repetition scaling law $L(C) = E + KC^{-\gamma}$ in §4.3 is fit on the six $OT = 1, R = 1$ baselines. We fit in log-loss space using nonlinear least squares: first estimating an initial K and γ from a log-log linear fit without a loss floor, then refitting E, K , and γ jointly with E free. This gives (6). The fitted curve matches the six no-repetition baselines to within about 0.015 nats, and we discuss sensitivity to this three-parameter fit in §4.3.

Reported evaluation loss. All reported eval losses are final-checkpoint losses. After each run reaches its target token budget $T = 20 \cdot OT \cdot N$, up to document-boundary rounding from the sampling procedure, we evaluate the final checkpoint once on the fixed held-out split and use that value in all figures, peak fits, and scaling-law fits. We use the same rule for repeated-data runs and no-repetition baselines. We do not report the best validation checkpoint or an average over checkpoints; intermediate evaluations are used only for monitoring.

G. Interpreting $CEG > 1$

The CEG ratio in (7) is defined relative to our fitted no-repetition scaling law. It is not an oracle optimum over all possible model sizes, token budgets, optimizers, and data mixtures. Therefore $CEG > 1$ can occur when a run achieves lower loss than the fitted no-repetition reference curve predicts at its actual compute. We interpret such values as being above this fitted reference, not as a universal claim that the run is optimally compute saving.

One reason this can happen is that $OT = 1$ uses the Chinchilla-style rule of 20 tokens per parameter. That rule was estimated in a different setting, with a different model family, optimizer stack, tokenizer, and data mixture. The optimal token-per-parameter ratio for our Qwen3-style ([Yang et al., 2025a](#)) models on FineWeb-Edu-Dedup ([Penedo et al., 2024](#)) may therefore differ from 20. We do not attempt to find the globally optimal (N, T) allocation for this model family, because the goal of the paper is to compare repetition structures at fixed budget and to convert their losses through a consistent no-repetition reference.

For the repeated-data comparisons, what matters is that all runs are evaluated against the same fitted no-repetition reference, and that each repeated-data run is also compared to the no-repetition baseline at the same (N, OT) budget. Values below one imply positive CEL relative to this reference. Values above one imply negative CEL and should be read as a calibration effect of the fitted reference curve.

990 H. Limitations

991 Our experiments are small relative to frontier pretraining: the largest model has 344M parameters, and all runs use one
992 Qwen3-style architecture family, one tokenizer, one corpus, and a fixed repeated-token fraction $f = 0.1$. Each (N, OT, R)
993 cell is a single training run, so we do not estimate seed-level variance; the evidence for robustness comes from consistency
994 across sweeps, not repeated random restarts. The larger-model grid is incomplete because of compute constraints, with
995 153M run through $OT = 2$ and 344M run only at $OT = 1$. The repeat-count sweep extends to $R = 20000$; the log-Gaussian
996 peak fits use the full range, though plots display only up to $R \approx 3000$ for readability. Finally, CEG and CEL depend on a
997 three-parameter no-repetition scaling law fit to six $OT = 1$ baselines, so their absolute values should be treated as point
998 estimates. The linear-regression model explains one mechanism by which repetition can hurt, but it leaves out many features
999 of real language models, including attention, depth, optimization dynamics, and discrete tokens. It should be read as an
1000 explanatory toy model, not a quantitative model of transformer pretraining.
1001
1002

1003 I. Broader Impacts

1004 This work studies how repeated training data can reduce the compute efficiency of language-model pretraining. Better
1005 measurement of repetition damage may help practitioners spend less compute and energy on runs whose data mixture is
1006 poorly structured. The same analysis could also be used to justify more aggressive data filtering, so care is needed to avoid
1007 removing useful minority-domain or low-resource-language data solely because it appears repetitive.
1008
1009

1010 J. Experiments Compute Resource

1011 All experiments were run on GPU clusters using BF16 training and FlashAttention-2. The main experimental grid consists
1012 of Qwen3-style models from 34M to 344M parameters, trained across repeated-data sweeps and no-repetition baselines
1013 as described in Appendix F. The largest individual runs are the 344M-parameter models at $OT = 1$, with total training
1014 compute estimated by $C = 6NT$.
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044

Excess test loss from repeated samples

Color is excess population test loss over a clean no-repetition baseline. Light cells are the repetition peak.

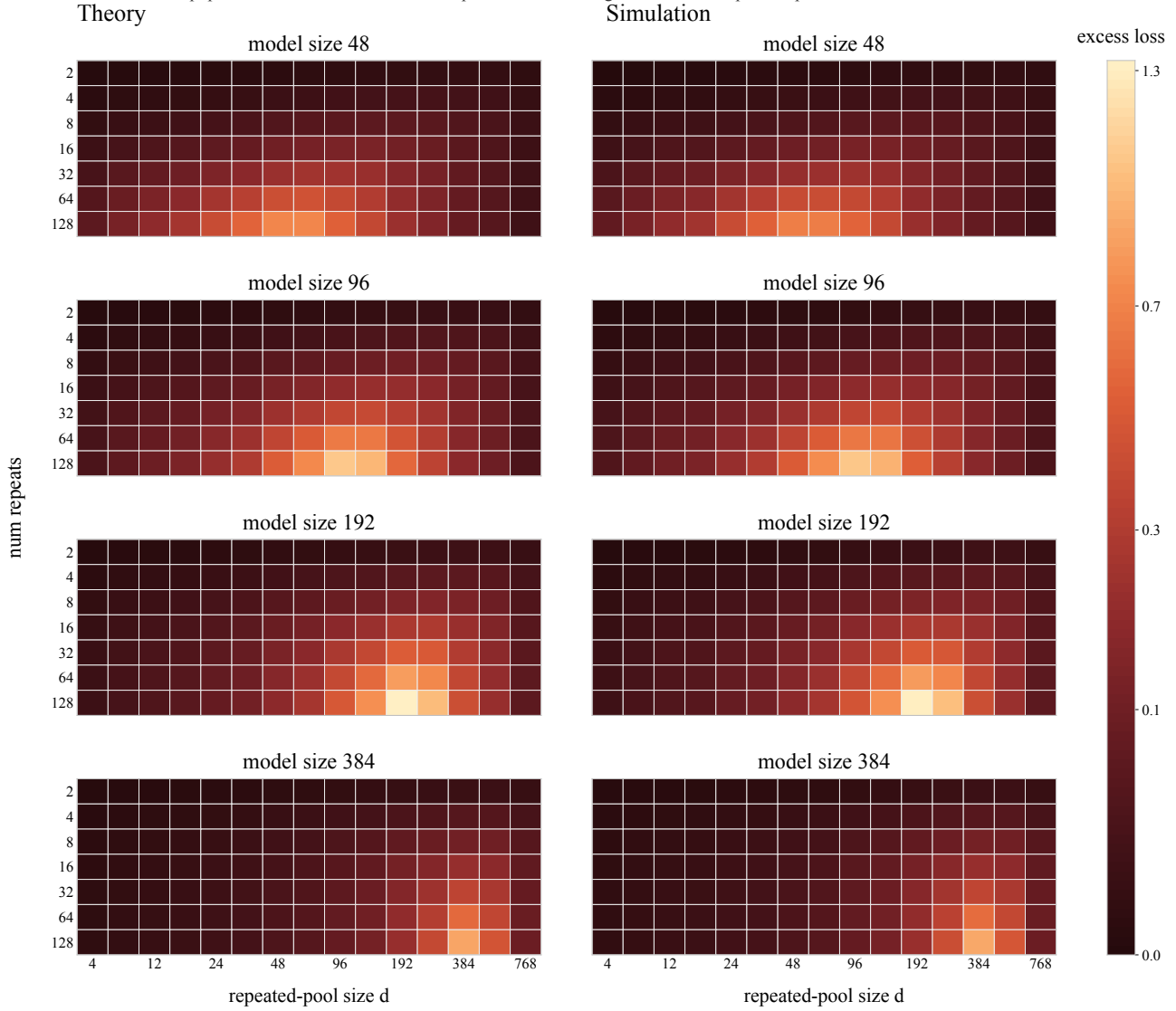


Figure 9. Full visualization of Figure 6