

Stochastic Gradient Descent on the Linear Bigram Model: Bias-Variance Scaling and Critical Batch Size

author names withheld

Under Review for the Workshop on High-dimensional Learning Dynamics, 2026

Abstract

The critical batch size, beyond which data parallelism gives diminishing returns, is a key factor in language-model pretraining. Existing finite-time least-squares theory provides useful templates for predicting it, but does not directly give finite-vocabulary Zipf rates with explicit dependence on the vocabulary size. The linear bigram model, which fits a next-token transition matrix under squared loss with power-law token frequencies, provides a tractable setting for this regime. Its one-hot sampling structure also departs from prior vector analyses, and so far it has only been studied in the deterministic full-batch case. We give a finite-time analysis of mini-batch SGD on this model under power-law token distributions. Our main result is an exact, closed-form bias–variance decomposition of the expected loss, in which the bias term equals the loss of deterministic gradient descent and the variance term captures the cost of mini-batch noise. From this decomposition we obtain scaling laws for the bias and the variance, governed by a frequency cutoff that separates rows with enough effective updates from rarer under-trained rows. When this cutoff reaches the full vocabulary, the learning curve changes phase. Balancing bias and variance yields the scaling of the critical batch size. We confirm the predicted scaling on simulated bigrams and on bigram statistics estimated from OpenWebText.

1. Introduction

Large language models are trained in a regime where token frequency, optimizer choice, and data parallelism are inseparable. Token frequencies follow Zipf’s law [13]; more generally, we write $\pi_i \propto i^{-\alpha}$ for the power-law exponent α , with Zipfian text corresponding to $\alpha \approx 1$. This gives embedding and unembedding layers highly nonuniform updates: common rows are sampled many times, while rare rows move on a much slower clock. This row-wise imbalance is one reason SGD is inefficient for language-model embeddings, while Adam-type adaptive methods are often preferred in practice [3, 4, 11]. Recent theory captures the deterministic part of this story: in the linear bigram model, rare tokens produce slow active-vocabulary growth and make gradient descent especially costly in the Zipfian case $\alpha = 1$ [3].

Practical pretraining, however, uses mini-batches. Mini-batching adds a stochastic variance term on top of the deterministic rare-token bias. In our decomposition, bias is the squared error of the mean averaged iterate, while variance is the fluctuation induced by mini-batch sampling. The active-vocabulary cutoff separates the dominant scaling contributions: rows beyond the cutoff remain bias-limited, whereas activated rows contribute the usual mini-batch variance floor; in heavier-tailed regimes, the inactive tail also contributes to variance. This distinction is central for predicting the

critical batch size (CBS), the point beyond which increasing data parallelism gives diminishing returns [1, 5, 6, 12].

General least-squares SGD theory provides sharp bias–variance decompositions and excess-risk bounds in terms of the full covariance spectrum [2, 7, 8, 14]. However, existing scaling-law analyses of stochastic least squares commonly turn these tools into clean power-law rates under trace-class or fast-decay regimes, and do not directly yield finite- V , normalized Zipf rates with explicit vocabulary dependence at the boundary $\alpha = 1$ or in heavier-tailed regimes $\alpha \leq 1$. The obstacle is not merely that the objective is quadratic. In standard linear regression, the second-moment dynamics involve fourth-moment operators such as $\mathbb{E}[xx^\top \Sigma xx^\top]$, so eigendirections are coupled through aggregate covariance or effective-dimension quantities. Thus current theory does not explain how mini-batch noise, finite-time optimization, and batch-size gains scale in the regime most relevant to language-model embeddings.

We address this gap in the linear bigram model of Kunstner and Bach [3]. The model maps the current one-hot token $x = e_i$ to a next-token prediction $W^\top x$, with token frequencies $\pi_i \propto i^{-\alpha}$ and exponent $\alpha > 0$. Its sampling structure closes the stochastic dynamics row-by-row, reducing the mean and second-moment evolution to one-dimensional recursions for each token row. These scalar recursions let us obtain an exact Polyak bias–variance formula and then split the risk at the active-vocabulary cutoff. We extend the deterministic picture to mini-batch SGD with Polyak averaging [9, 10] and compute the finite-time stochastic risk exactly.

In this work, we make the following contributions:

1. **Exact bias–variance decomposition for averaged SGD** (Theorem 3). We derive a closed-form expression for the expected excess risk of the Polyak-averaged iterate, separating contributions from a deterministic bias term and a stochastic variance term.
2. **Two-phase scaling laws** (Theorem 4). We derive scaling laws for the bias and variance terms, organized by the active-vocabulary cutoff that reflects the effective dimension and the critical time at which this cutoff reaches the full vocabulary. Beyond this critical time, the bias and variance terms recover the statistically optimal $\Theta(1/T^2)$ and $\Theta(1/(bT))$ rates.
3. **Fixed-compute critical batch size** (Proposition 6). We derive the scaling of $b_{\text{crit}}(N)$ in V , α , and compute budget $N = bT$. Across all α , the compute dependence is $b_{\text{crit}}(N) \propto \sqrt{N}$, matching the empirically observed square-root law for language-model pretraining [12].

2. Model and SGD Dynamics

Let $[V] = \{1, \dots, V\}$. Let W^* denote the population optimum, and define the iterate error $\Delta_t := W_t - W^*$, with (i, j) -coordinate $\Delta_{ij,t}$. We assume power laws for the marginal token frequencies π_i and the sorted conditional next-token frequencies $\pi_{j|i}$:

Assumption 1 (Power-law marginals and conditionals) *Given vocabulary size V , for all $i, j \in [V]$, we have $\pi_i = z^{-1}i^{-\alpha}$, $\pi_{\rho_i(j)|i} = z^{-1}j^{-\alpha}$, where $z = \sum_{k=1}^V k^{-\alpha}$ is the power-law normalizing constant, and ρ_i is a sorting permutation that ranks the conditional frequencies in decreasing order.*

We remark that the assumption does not require all conditionals to follow the same distribution. It only states that the sorted frequencies follow a power law with the same exponent α . Our experiments on OpenWebText (Figure 3 and Figure 4) suggest this provides a reasonable coarse-grained approximation of real-world language data.

Lemma 2 (Population structure) *Let $W \in \mathbb{R}^{V \times V}$, and let x, y be one-hot vectors in \mathbb{R}^V . Then the linear predictor $f_W(x) = W^\top x$ trained with squared loss $\mathcal{L}(W) = \frac{1}{2} \mathbb{E} \|W^\top x - y\|_2^2$, where $\mathbb{P}(x = e_i) = \pi_i$ and $\mathbb{P}(y = e_j \mid x = e_i) = \pi_{j|i}$, has population minimizer $W_{ij}^* = \pi_{j|i}$, diagonal second moment $D = \text{diag}(\pi_1, \dots, \pi_V)$, and excess risk*

$$\mathcal{E}(W) = \mathcal{L}(W) - \mathcal{L}(W^*) = \frac{1}{2} \sum_{i=1}^V \pi_i \|W_{i,:} - W_{i,:}^*\|_2^2. \quad (1)$$

Thus the problem is a quadratic objective in the row-vectorized parameter W , with Hessian matrix $D \otimes I_V$, and the loss decouples row-wise with weight π_i .

Let $\{W_t\}_{t=0}^T$ be iterates generated by online mini-batch SGD with batch size b : $W_{t+1} = W_t - \eta G_t$ where $G_t = \frac{1}{b} \sum_{\ell=1}^b x_t^{(\ell)} (W_t^\top x_t^{(\ell)} - y_t^{(\ell)})^\top$. We study the Polyak average $\bar{W}_T = T^{-1} \sum_{t=0}^{T-1} W_t$ with initialization at $W_0 = 0$.

3. Main Results

Our first result gives the bias–variance decomposition of the expected excess risk. The derivation has two steps. First, conditional on the current iterate W_t , the mini-batch gradient is unbiased, so the mean dynamics recover gradient descent: $\mathbb{E}[\Delta_{ij,t+1} \mid W_t] = \phi_i \Delta_{ij,t}$, where $\phi_i := 1 - \eta \pi_i$ is the per-row iterate contraction factor. Second, the mini-batch noise enters through the second moment of the iterate error, $\mathbb{E}[\Delta_{ij,t+1}^2 \mid W_t] = a_i \Delta_{ij,t}^2 + \{\text{noise floor}\}$, where $a_i := \phi_i^2 + \frac{\eta^2}{b} \pi_i (1 - \pi_i)$ is the stochastic second-moment contraction factor.

Theorem 3 (Exact bias–variance decomposition) *Under Assumption 1, suppose $W_0 = 0$ and choose a stable step size $\eta < 2bz/(b+z-1)$ for batch size b . Define the per-row contraction factor $\phi_i = 1 - \eta \pi_i$ and second-moment contraction factor $a_i = \phi_i^2 + \frac{\eta^2}{b} \pi_i (1 - \pi_i)$. Then for every $T \geq 1$,*

$$\mathbb{E}[\mathcal{E}(\bar{W}_T)] = \underbrace{\frac{\mathcal{E}(0)}{T^2} \sum_{i=1}^V \pi_i \Gamma_i(T)}_{\mathcal{E}_{\text{bias}}(T)} + \underbrace{\frac{\mathcal{E}(0)}{T^2} \sum_{i=1}^V \pi_i (S_i^{(a)}(T) - \Gamma_i(T)) + \frac{\eta^2(1-2\mathcal{E}(0))}{2bT^2} \sum_{i=1}^V \frac{\pi_i^2}{1-a_i} S_i^{(1-a)}(T)}_{\mathcal{L}_{\text{var}}(T)},$$

where $\Gamma_i(T) := \left(\frac{1-\phi_i^T}{1-\phi_i}\right)^2$, and $S_i^{(a)}(T)$ and $S_i^{(1-a)}(T)$ are the weighted temporal sums w.r.t. a_i :

$$S_i^{(a)}(T) := \sum_{s=0}^{T-1} w_{i,s} a_i^s, \quad S_i^{(1-a)}(T) := \sum_{s=0}^{T-1} w_{i,s} (1 - a_i^s) \quad \text{with weights } w_{i,s} := 1 + 2 \sum_{k=1}^{T-1-s} \phi_i^k.$$

The advantage of the linear bigram model is that its learning dynamics separate cleanly across frequency modes. Each row-coordinate has a scalar mean contraction ϕ_i and scalar second-moment contraction a_i , so the excess loss can be written exactly as a sum of modal contributions. This is the technical resolution of the least-squares coupling challenge: in standard SGD linear regression, the second-moment recursion is governed by fourth-moment operators and each eigendirection is affected by aggregate covariance across directions [7]. The row-wise recursion here lets us directly identify which frequencies have been learned, which remain underfit, and how different frequency ranges shape the learning curve.

Building on this decomposition, we study how the bias and variance terms scale with large vocabulary size V . We report the relative quantities $r_{\text{bias}}(T) := \frac{\mathcal{E}_{\text{bias}}(T)}{\mathcal{E}(0)}$ and $r_{\text{var}}(T) := \frac{\mathcal{L}_{\text{var}}(T)}{\mathcal{E}(0)}$, which allow comparisons across different V and α .

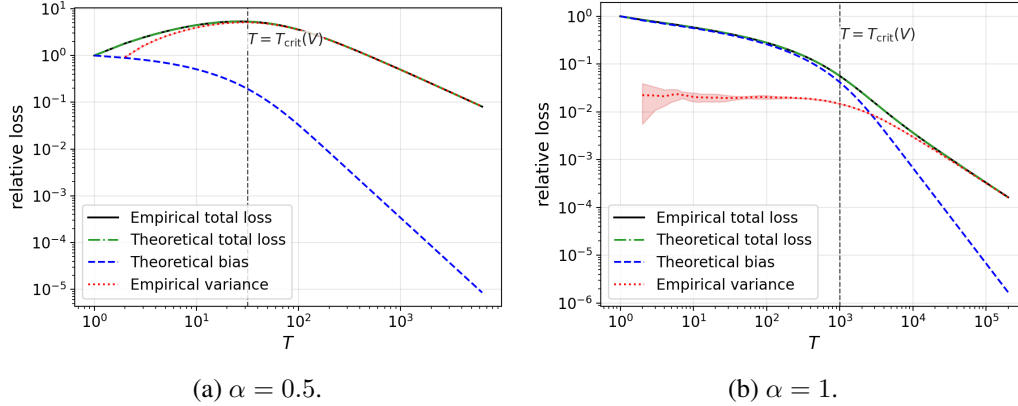


Figure 1: **Exact bias–variance decomposition.** Online SGD follows the exact theoretical decomposition closely for $V = 1000$, $b = 1024$ at $\alpha = 0.5$ and $\alpha = 1$. Simulations use five independent seeds and $\eta = z$. The dashed vertical line denotes $T_{\text{crit}}(V) = V^\alpha$.

Theorem 4 (Bias and variance scaling) *Let the learning-rate rescaling factor be $\kappa = \eta/z \leq 1$, the row-wise effective learning rate be $\delta_i = \eta\pi_i = \kappa i^{-\alpha}$, the batch size be b , and the active-vocabulary cutoff be $i_{\kappa,*}(T) = \min\{V, \lfloor (\kappa T)^{1/\alpha} \rfloor\}$. For $\eta \leq z$, the relative bias satisfies*

$$r_{\text{bias}}(T) \asymp \frac{1}{z} \left(\frac{1}{\kappa^2 T^2} \sum_{i \leq i_{\kappa,*}(T)} i^\alpha + \sum_{i > i_{\kappa,*}(T)} i^{-\alpha} \right). \quad (2)$$

For $\eta \leq \min\{z, b/2\}$, the relative variance satisfies

$$r_{\text{var}}(T) \asymp \frac{1}{b\mathcal{E}(0)} \left(\frac{i_{\kappa,*}(T)}{T} + \kappa^2 T \sum_{i > i_{\kappa,*}(T)} i^{-2\alpha} \right). \quad (3)$$

The two terms in (3) refine the activation interpretation. Activated rows contribute the mini-batch variance floor $i_{\kappa,*}(T)/T$, while inactive rows contribute the tail term $\kappa^2 T \sum_{i > i_{\kappa,*}(T)} i^{-2\alpha}$. Thus inactive rows are not only a source of residual bias; in sufficiently heavy-tailed regimes, their stochastic fluctuations can also control the variance scaling.

The active-vocabulary cutoff $i_{\kappa,*}(T)$ can be interpreted as the effective dimension at time T . High-frequency rows with indices below $i_{\kappa,*}(T)$ have been activated by time T and are in their asymptotic regime, while lower-frequency rows remain under-activated and drive the residual bias through the tail sums. This cutoff reaches the full vocabulary at the *critical time* $T_{\text{crit}}^{(\kappa)} = \kappa^{-1} V^\alpha$, the smallest time at which all V tokens have been activated.

Corollary 5 (Zipfian total risk) *For $\alpha = 1$, with learning rate $\eta = z$ and batch size $b \geq 2z$, the relative total loss $r_{\text{total}}(T) := r_{\text{bias}}(T) + r_{\text{var}}(T)$ has the following scaling:*

$$r_{\text{total}}(T) \asymp \begin{cases} 1 - \frac{\log T}{\log V} + \frac{(\log V)^2}{b}, & T \lesssim V \\ \frac{V^2}{T^2 \log V} + \frac{V(\log V)^2}{bT}, & T \gg V \end{cases}$$

Corollary 5 recovers the worst-case $\Theta(1 - \log T / \log V)$ scaling [3] in the bias term, while the variance term stays flat in the pre-critical regime. A further consequence is that for heavier-tailed regimes $\alpha < 1$, optimization difficulty can stem from the increasing variance term (Figure 1 panel

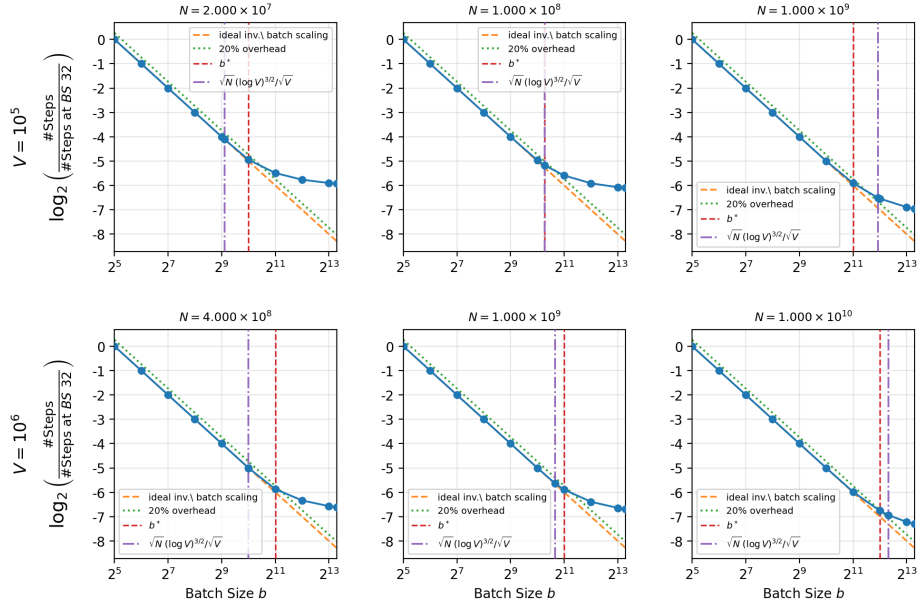


Figure 2: **Critical-batch-size validation for $\alpha = 1$.** The measured critical batch size b^* , defined by a 20% overhead in T_{hit} relative to ideal inverse-batch scaling, aligns with the theoretical value $b_{\text{crit}}(N)$ at the predicted scale.

(a) and Figure 8), even though the critical time is smaller and the bias term scales better than in the Zipfian case. We define the *critical batch size* (CBS) as the batch size that balances the asymptotic rates of the bias and variance terms, and validate it experimentally (Figure 2) against the empirical CBS definition used by [12].

Proposition 6 (Critical batch size) *In the supercritical regime $T \gg V^\alpha$, $b_{\text{crit}}(T) \asymp \frac{z}{\mathcal{E}(0)} \frac{T}{V^\alpha}$. At fixed sample compute $N = bT$, the Zipfian $\alpha = 1$ case satisfies $b_{\text{crit}}(N) \asymp \sqrt{N} \frac{(\log V)^{3/2}}{\sqrt{V}}$.*

The full regime-wise critical-batch scaling, together with the bias and variance scalings before and after the critical time, is summarized in Appendix A.

4. Experiments

We validate the exact decomposition on synthetic power-law bigrams (Figure 1) and test the CBS prediction with a batch-size sweep (Figure 2; details are deferred to Appendix B). We also validate the Zipf power-law assumption on OpenWebText in Appendix B.1.

5. Conclusion

We analyzed mini-batch SGD with Polyak averaging in the linear bigram model under power-law token frequencies. The exact bias–variance decomposition exposes row-wise frequency-mode dynamics that are hidden in standard linear regression analyses. From this decomposition, we derive bias and variance scaling laws across two phases and all power-law exponent regimes. We then use the supercritical scaling to predict the critical batch size and validate its empirical scale.

References

- [1] Shane Bergsma, Nolan Dey, Gurpreet Gosal, Gavia Gray, Daria Soboleva, and Joel Hestness. Power Lines: Scaling Laws for Weight Decay and Batch Size in LLM Pre-training, November 2025.
- [2] Prateek Jain, Sham M. Kakade, Rahul Kidambi, Praneeth Netrapalli, Venkata Krishna Pillutla, and Aaron Sidford. A Markov Chain Theory Approach to Characterizing the Minimax Optimality of Stochastic Gradient Descent (for Least Squares). *LIPIcs, Volume 93, FSTTCS 2017*, 93:2:1–2:10, 2018. ISSN 1868-8969. doi: 10.4230/LIPIcs.FSTTCS.2017.2.
- [3] Frederik Kunstner and Francis Bach. Scaling Laws for Gradient Descent and Sign Descent for Linear Bigram Models under Zipf’s Law, May 2025.
- [4] Binghui Li, Kaifei Wang, Han Zhong, Pinyan Lu, and Liwei Wang. Muon in Associative Memory Learning: Training Dynamics and Scaling Laws, February 2026.
- [5] Sam McCandlish, Jared Kaplan, Dario Amodei, and OpenAI Dota Team. An Empirical Model of Large-Batch Training, December 2018.
- [6] William Merrill, Shane Arora, Dirk Groeneveld, and Hannaneh Hajishirzi. Critical Batch Size Revisited: A Simple Empirical Approach to Large-Batch Language Model Training, November 2025.
- [7] Alexandru Meterez, Depen Morwani, Costin-Andrei Oncescu, Jingfeng Wu, Cengiz Pehlevan, and Sham Kakade. A Simplified Analysis of SGD for Linear Regression with Weight Averaging, June 2025.
- [8] Alexandru Meterez, Pranav Ajit Nair, Depen Morwani, Cengiz Pehlevan, and Sham Kakade. Anytime Pretraining: Horizon-Free Learning-Rate Schedules with Weight Averaging, February 2026.
- [9] Boris T. Polyak. New method of stochastic approximation type. *Automation and Remote Control*, 51:937–946, 1990.
- [10] Boris T. Polyak and Anatoli B. Juditsky. Acceleration of stochastic approximation by averaging. *SIAM Journal on Control and Optimization*, 30(4):838–855, 1992.
- [11] Guodong Zhang, Lala Li, Zachary Nado, James Martens, Sushant Sachdeva, George E. Dahl, Christopher J. Shallue, and Roger Grosse. Which Algorithmic Choices Matter at Which Batch Sizes? Insights From a Noisy Quadratic Model, October 2019.
- [12] Hanlin Zhang, Depen Morwani, Nikhil Vyas, Jingfeng Wu, Difan Zou, Udaya Ghai, Dean Foster, and Sham Kakade. How Does Critical Batch Size Scale in Pre-training?, April 2025.
- [13] George Kingsley Zipf. *Human Behavior and the Principle of Least Effort*. Addison-Wesley, 1949.
- [14] Difan Zou, Jingfeng Wu, Vladimir Braverman, Quanquan Gu, and Sham M. Kakade. Benign overfitting of constant-stepsize SGD for linear regression. In *Conference on Learning Theory (COLT)*, 2021.

AI Assistance Disclosure

The authors used AI assistance for editing, formatting, and submission preparation. All scientific content, claims, proofs, and final text were reviewed and approved by the authors.

Appendix Contents

A Scaling Tables **8**

B Additional Experimental Details **9**

 B.1 Zipf Law Validation 11

 B.2 Exact Decomposition Figures 11

 B.3 Surrogate Scaling Verification 11

C Proofs for Preliminaries **18**

 C.1 Proof of Lemma 2 (Population structure) 18

 C.2 Asymptotics of the normalizing constant 18

 C.3 Asymptotics of $\mathcal{E}(0)$ 18

D Stochastic Gradient Analysis **19**

 D.1 Stochastic gradient moments 19

 D.2 Coordinate recursion 20

 D.3 Mean trajectory and bias 20

 D.4 Second moments and temporal covariances 21

E Proof of the Bias–Variance Decomposition (Theorem 3) **23**

F Proof of the Bias Scaling Law (Theorem 4) **24**

G Proof of the Variance Scaling Law (Theorem 4) **26**

H Proof of the Total Scaling and Critical Batch Size **31**

 H.1 Proof of Corollary 5 31

 H.2 Proof of Proposition 6 (Critical batch size) 31

Appendix A. Scaling Tables

Regime	$b_{\text{crit}}(N)$ at fixed sample compute $N = bT$
$0 < \alpha < 1/2$	$\sqrt{N} V^{1-\alpha}$
$\alpha = 1/2$	$\sqrt{NV/\log V}$
$1/2 < \alpha < 1$	$\sqrt{N} V^{3/2-2\alpha}$
$\alpha = 1$	$\sqrt{N} (\log V)^{3/2}/\sqrt{V}$
$\alpha > 1$	$\sqrt{N}/V^{\alpha/2}$

Table 1: Regime-wise critical-batch-size scaling at fixed sample compute $N = bT$.

Regime	Pre-critical $T \lesssim V^\alpha$		Post-critical $T \gg V^\alpha$	
	$r_{\text{bias}}(T)$	$r_{\text{var}}(T)$	$r_{\text{bias}}(T)$	$r_{\text{var}}(T)$
$0 < \alpha < 1/2$	$1 - \frac{2\alpha}{1+\alpha} \frac{T^{(1-\alpha)/\alpha}}{V^{1-\alpha}}$	$TV^{2-2\alpha}/b$	$V^{2\alpha}/T^2$	$V^2/(bT)$
$\alpha = 1/2$	$1 - \frac{2\alpha}{1+\alpha} \frac{T^{(1-\alpha)/\alpha}}{V^{1-\alpha}}$	TV/b	$V^{2\alpha}/T^2$	$V^2/(bT \log V)$
$1/2 < \alpha < 1$	$1 - \frac{2\alpha}{1+\alpha} \frac{T^{(1-\alpha)/\alpha}}{V^{1-\alpha}}$	$T^{1/\alpha-1}V^{2-2\alpha}/b$	$V^{2\alpha}/T^2$	$V^{3-2\alpha}/(bT)$
$\alpha = 1$	$1 - \log T / \log V$	$(\log V)^2/b$	$V^2/(T^2 \log V)$	$V(\log V)^2/(bT)$
$\alpha > 1$	$T^{(1-\alpha)/\alpha}$	$T^{1/\alpha-1}/b$	$V^{\alpha+1}/T^2$	$V/(bT)$

Table 2: Relative bias and variance scalings across power-law regimes, before and after the critical time $T_{\text{crit}}(V) = V^\alpha$ for the canonical step size $\eta = z$. The variance column includes its $1/b$ batch-size dependence, while the bias is batch-size independent.

Appendix B. Additional Experimental Details

Critical-batch sweep. The experiment in Figure 2 uses the exact Polyak bias and variance formulas. For each (V, α, N, b) , it evaluates relative risk at the nominal horizon $T = \lfloor N/b \rfloor$. We fix a reference batch size b_{ref} that sits under the linear scaling regime, and defines the reference loss L_{ref} at $T_{\text{ref}} = \lfloor N/b_{\text{ref}} \rfloor$. For every other batch size, the experiment records the smallest T_{hit} for which the closed-form Polyak risk is at most L_{ref} , and compares it with the ideal inverse-batch prediction $T_{\text{ideal}} = T_{\text{ref}}b_{\text{ref}}/b$. The plotted knee is where the overhead $T_{\text{hit}}/T_{\text{ideal}} - 1$ becomes non-negligible. We follow the 20% overhead convention as in [12] to define empirical critical batch size b^* . The vertical theory reference is $b_{\text{crit}}(N)$, e.g. $\sqrt{N}(\log V)^{3/2}/\sqrt{V}$ for $\alpha = 1$.

Learning-rate sweep for Figure 2. The batch-size sweep uses the synthetic Zipf–Markov linear bigram model in the $\alpha = 1$ regime. For each plotted triple (V, N, b) , we set $T = \lfloor N/b \rfloor$, and sweep the rescaled learning rate $\kappa = \eta/z$ on a dense log grid over $(0, 1]$, followed by local refinement around the best point. Each candidate η is scored by the $\alpha = 1$ surrogate Polyak risk

$$r_{\text{total}}(\eta) = \frac{1}{z} \left(\frac{i_*(i_* + 1)}{2\kappa^2 T^2} + H_V - H_{i_*} \right) + \frac{1}{b\mathcal{E}(0)} \left(\frac{i_*}{T} + \kappa^2 T \left(H_V^{(2)} - H_{i_*}^{(2)} \right) \right), \quad i_* = \min\{V, \lfloor \kappa T \rfloor\}.$$

The optimum η_* is reported in Table 3. The main pattern is simple: in the lowest-compute panels, η_* grows almost linearly with batch size so that the active cutoff i_* stays nearly fixed; once the budget is comfortably supercritical, the sweep usually saturates at the canonical rate $\eta_* = z$.

Table 3: Learning-rate sweep used for the critical-batch experiment.
 Here $z = H_V$, $\kappa_* = \eta_*/z$, and $i_* = \lfloor \kappa_* T \rfloor$ after clipping at V .
 Numeric columns use grouped digits for readability.

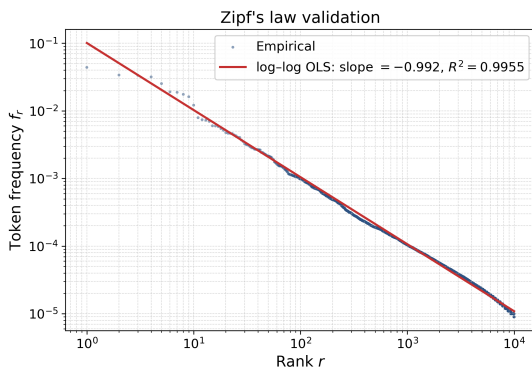
V	N (samples)	b	$T = N/b$	η_*	i_*
$V = 10^5$, $z = 12.090146$					
100,000	20,000,000	32	625,000	0.094585	4,889
100,000	20,000,000	64	312,500	0.189579	4,900
100,000	20,000,000	128	156,250	0.378586	4,892
100,000	20,000,000	256	78,125	0.756032	4,885
100,000	20,000,000	512	39,062	1.515335	4,895
100,000	20,000,000	1,024	19,531	3.031661	4,897
100,000	20,000,000	2,048	9,765	6.061716	4,895
100,000	20,000,000	4,096	4,882	12.090146	4,882
100,000	20,000,000	8,192	2,441	12.090146	2,441
100,000	100,000,000	32	3,125,000	12.090146	100,000
100,000	100,000,000	64	1,562,500	12.090146	100,000
100,000	100,000,000	128	781,250	12.090146	100,000
100,000	100,000,000	256	390,625	12.090146	100,000
100,000	100,000,000	512	195,312	12.090146	100,000
100,000	100,000,000	1,024	97,656	12.090146	97,656
100,000	100,000,000	2,048	48,828	9.115626	36,814
100,000	100,000,000	4,096	24,414	12.090146	24,414
100,000	100,000,000	8,192	12,207	12.090146	12,207
100,000	1,000,000,000	32	31,250,000	12.090146	100,000
100,000	1,000,000,000	64	15,625,000	12.090146	100,000
100,000	1,000,000,000	128	7,812,500	12.090146	100,000
100,000	1,000,000,000	256	3,906,250	12.090146	100,000
100,000	1,000,000,000	512	1,953,125	12.090146	100,000
100,000	1,000,000,000	1,024	976,562	12.090146	100,000
100,000	1,000,000,000	2,048	488,281	12.090146	100,000
100,000	1,000,000,000	4,096	244,140	12.090146	100,000
100,000	1,000,000,000	8,192	122,070	12.090146	100,000
$V = 10^6$, $z = 14.392727$					
1,000,000	400,000,000	32	12,500,000	0.067521	58,641
1,000,000	400,000,000	64	6,250,000	0.135085	58,660
1,000,000	400,000,000	128	3,125,000	0.269764	58,572
1,000,000	400,000,000	256	1,562,500	0.539705	58,591
1,000,000	400,000,000	512	781,250	1.077784	58,503
1,000,000	400,000,000	1,024	390,625	2.160231	58,629
1,000,000	400,000,000	2,048	195,312	4.313955	58,541
1,000,000	400,000,000	4,096	97,656	8.640061	58,623
1,000,000	400,000,000	8,192	48,828	14.392727	48,828

Continued on next page

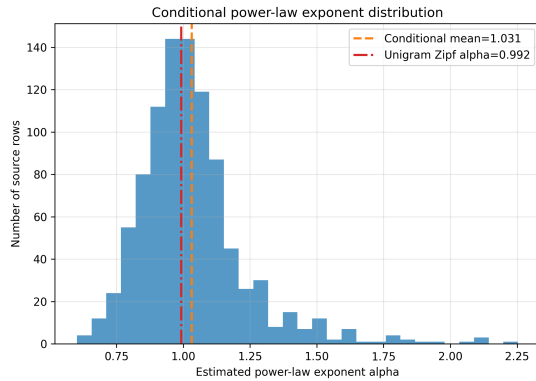
V	N (samples)	b	$T = N/b$	η_*	i_*
1,000,000	1,000,000,000	32	31,250,000	0.076208	165,466
1,000,000	1,000,000,000	64	15,625,000	0.152187	165,217
1,000,000	1,000,000,000	128	7,812,500	0.304474	165,271
1,000,000	1,000,000,000	256	3,906,250	0.609147	165,325
1,000,000	1,000,000,000	512	1,953,125	1.216459	165,076
1,000,000	1,000,000,000	1,024	976,562	2.438180	165,433
1,000,000	1,000,000,000	2,048	488,281	4.864101	165,017
1,000,000	1,000,000,000	4,096	244,140	9.747727	165,348
1,000,000	1,000,000,000	8,192	122,070	14.392727	122,070
1,000,000	10,000,000,000	32	312,500,000	14.392727	1,000,000
1,000,000	10,000,000,000	64	156,250,000	14.392727	1,000,000
1,000,000	10,000,000,000	128	78,125,000	14.392727	1,000,000
1,000,000	10,000,000,000	256	39,062,500	14.392727	1,000,000
1,000,000	10,000,000,000	512	19,531,250	14.392727	1,000,000
1,000,000	10,000,000,000	1,024	9,765,625	14.392727	1,000,000
1,000,000	10,000,000,000	2,048	4,882,812	14.392727	1,000,000
1,000,000	10,000,000,000	4,096	2,441,406	14.392727	1,000,000
1,000,000	10,000,000,000	8,192	1,220,703	14.392727	1,000,000

B.1. Zipf Law Validation

We estimate bigram statistics from 5×10^6 OpenWebText tokens using the `tf5-small` SentencePiece tokenizer, restricted to the top $V = 10,000$ types. The empirical marginal is Zipfian, and among 948 well-supported conditional rows the median fitted exponent is $\hat{\alpha} \approx 1.01$.



(a) Marginal Zipf law.



(b) Conditional-row fitted exponents.

Figure 3: OpenWebText validation of the power-law bigram assumption.

B.2. Exact Decomposition Figures

B.3. Surrogate Scaling Verification

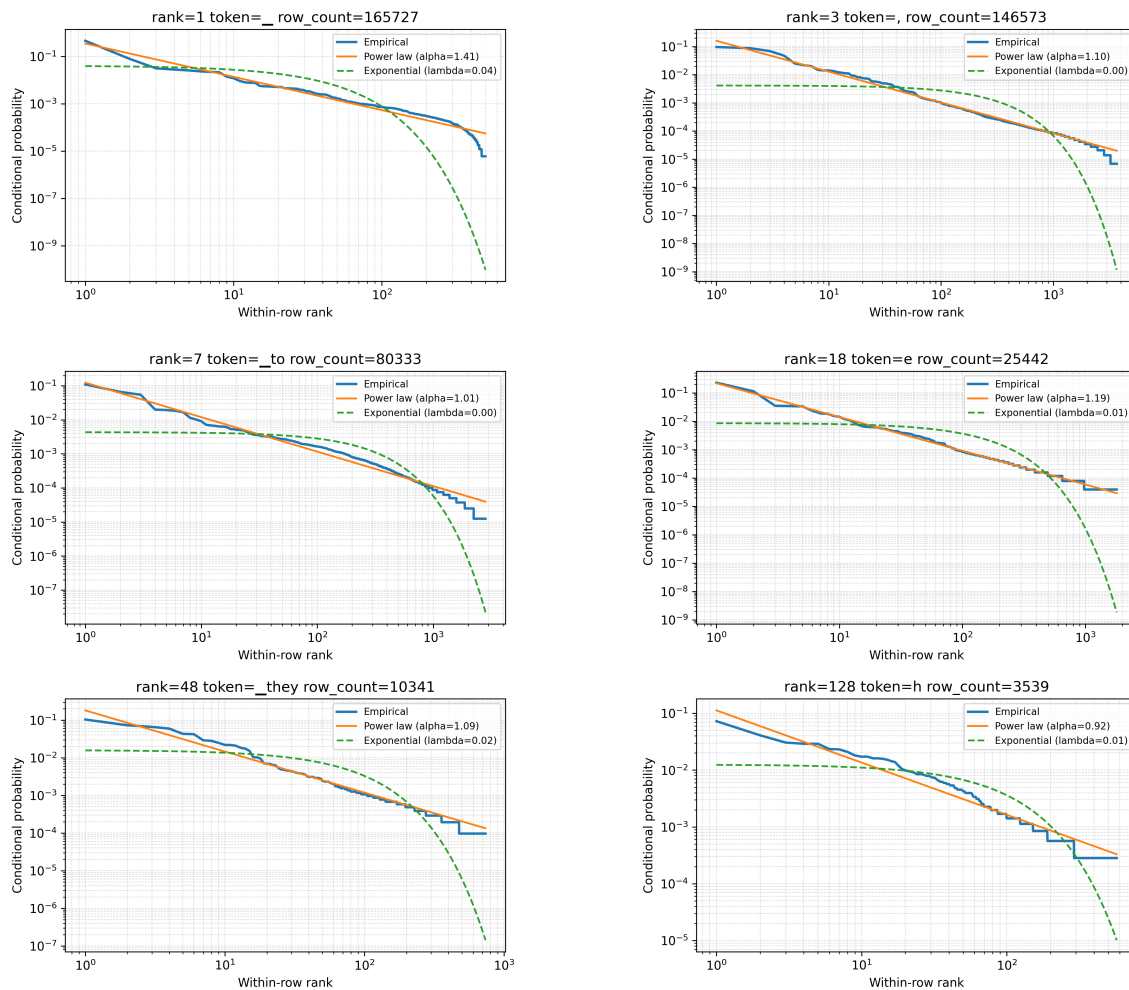


Figure 4: **Per-row power-law overlays.** Sorted-tail log–log plots of six representative rows of the empirical bigram transition matrix, with fitted finite-support discrete power-law overlays.

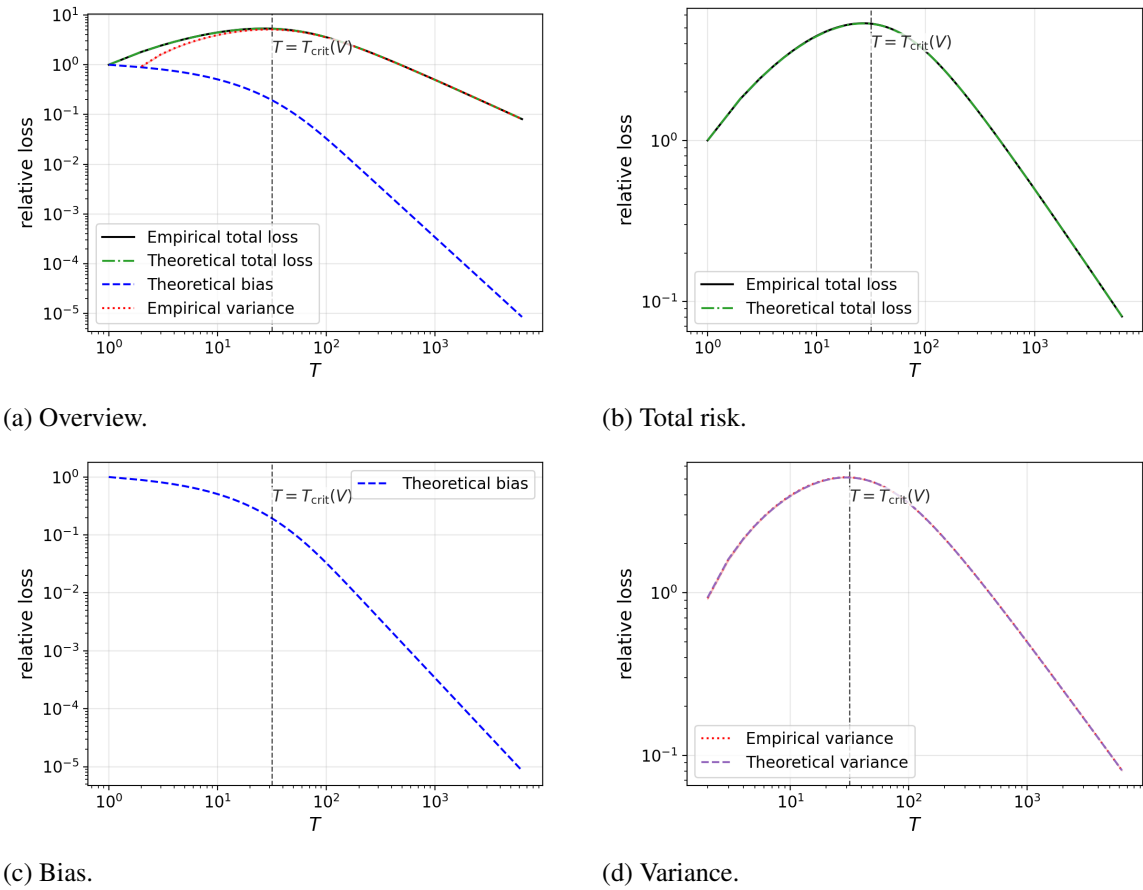


Figure 5: **Exact bias–variance decomposition**, $\alpha = 0.5$. Detailed views for $V = 1000$, $b = 1024$, and $K = 5$ seeds.

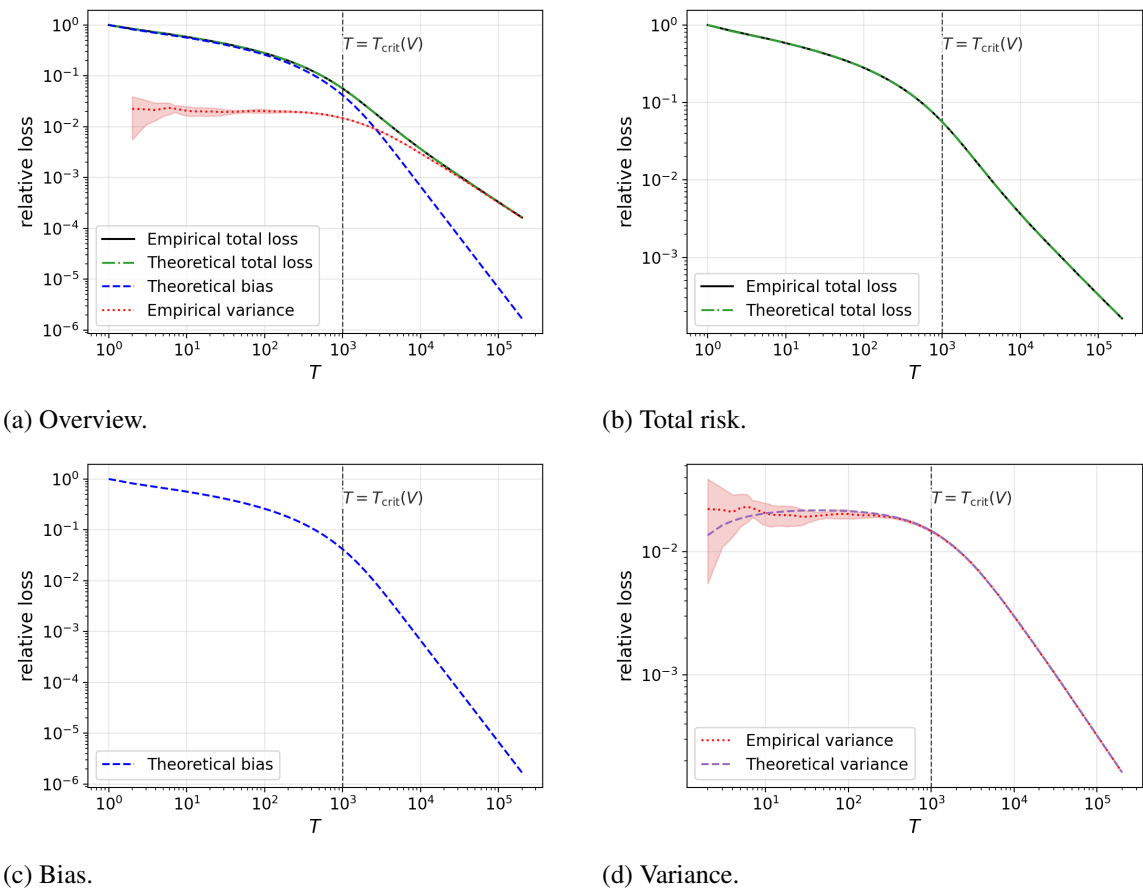


Figure 6: **Exact bias–variance decomposition**, $\alpha = 1$. Detailed views for $V = 1000$, $b = 1024$, and $K = 5$ seeds.

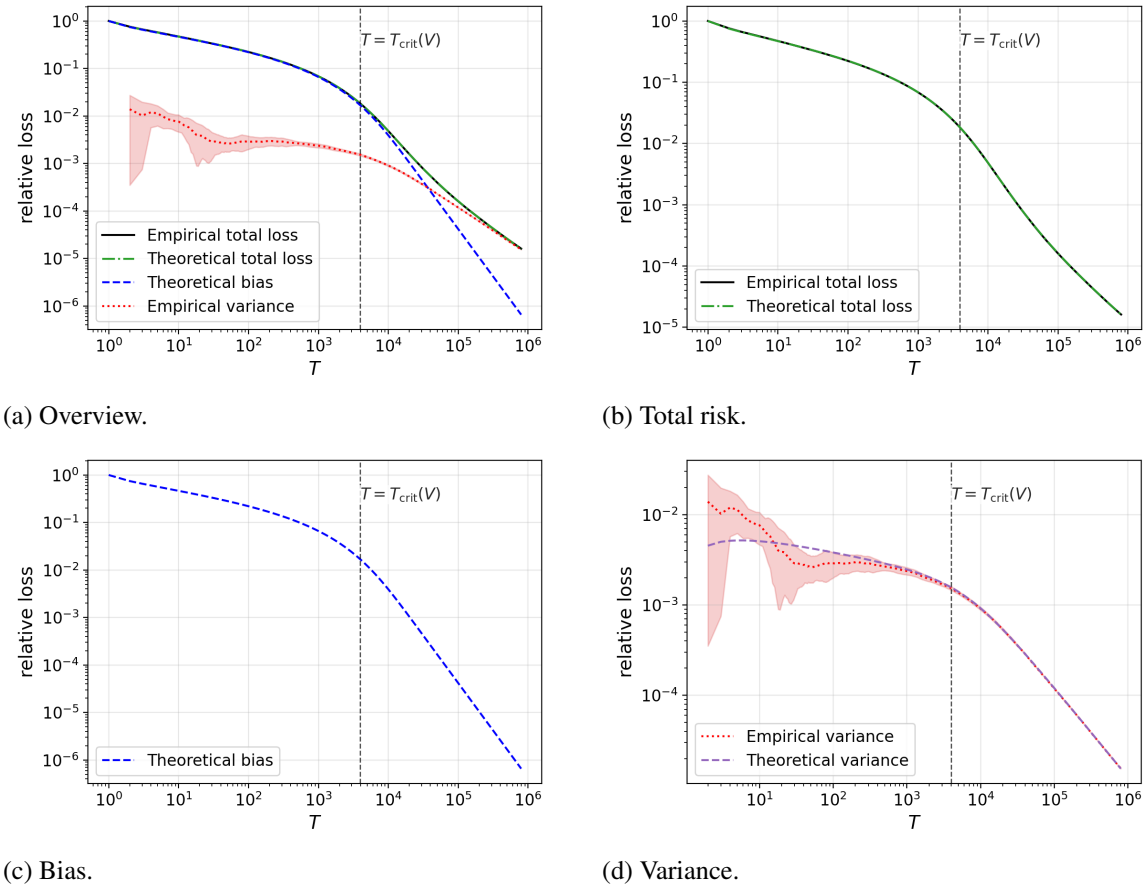


Figure 7: **Exact bias–variance decomposition**, $\alpha = 1.2$. Detailed views for $V = 1000$, $b = 1024$, and $K = 5$ seeds.

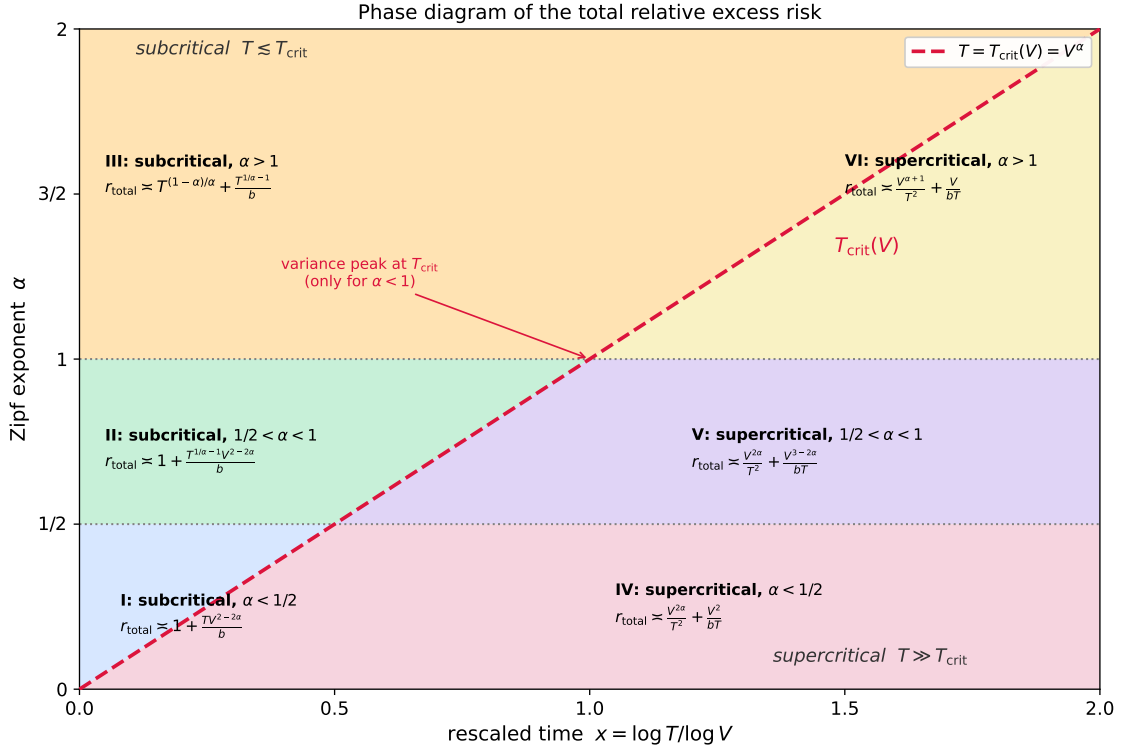


Figure 8: Phase diagram of total relative excess risk over $x = \log T / \log V$ and α . The red diagonal marks $T_{\text{crit}}(V) = V^\alpha$. For $\alpha < 1$, the variance term increases during subcritical phrase.

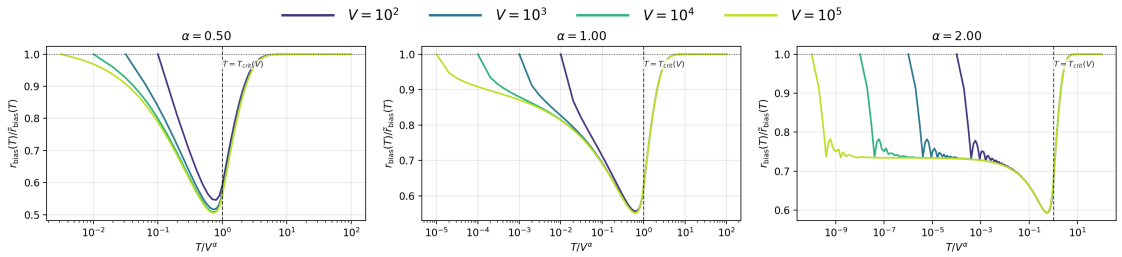


Figure 9: **Bias surrogate verification across vocabulary sizes and power-law exponents.** Ratio $r_{\text{bias}}(T) / \tilde{r}_{\text{bias}}(T)$ versus T/V^α for $V \in \{10^2, 10^3, 10^4, 10^5\}$ and $\alpha \in \{0.5, 1, 2\}$, using the surrogate in (2). The ratio stays bounded across all V, α regime.

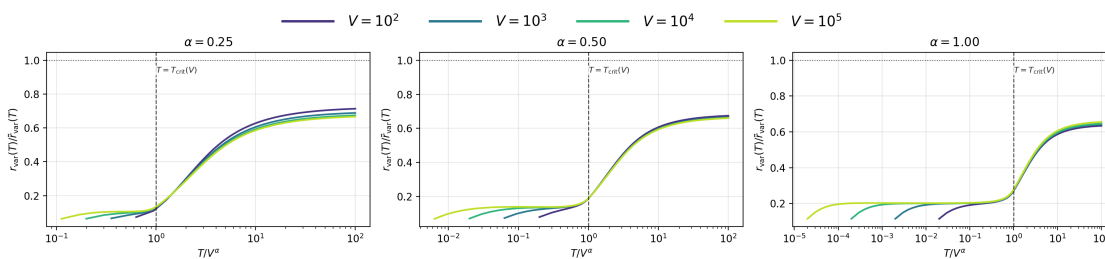


Figure 10: **Variance surrogate verification across vocabulary sizes and power-law exponents.** For $V \in \{10^2, 10^3, 10^4, 10^5\}$ and $\alpha \in \{0.25, 0.5, 1\}$, using the surrogate in 3. The ratio stays bounded across all V, α regime.

Appendix C. Proofs for Preliminaries

C.1. Proof of Lemma 2 (Population structure)

Proof Since $x = e_I$, we have $xx^\top = e_I e_I^\top$, and therefore

$$D = \mathbb{E}[xx^\top] = \sum_{i=1}^V \mathbb{P}(I = i) e_i e_i^\top = \text{diag}(\pi_1, \dots, \pi_V).$$

Next, conditioning on I gives

$$\mathcal{L}(W) = \frac{1}{2} \sum_{i=1}^V \pi_i \mathbb{E} \left[\|W_{i,:}^\top - y\|_2^2 \mid I = i \right].$$

Therefore the minimization over W decouples row-by-row. For fixed i , the unique minimizer is $\mathbb{E}[y \mid I = i] = \sum_{j=1}^V \pi_{j|i} e_j$, hence $[W^*]_{ij} = \pi_{j|i}$.

For the quadratic expansion, we apply the conditional bias–variance identity:

$$\mathbb{E}[\|a - y\|_2^2 \mid I = i] = \|a - \mathbb{E}[y \mid I = i]\|_2^2 + \mathbb{E}[\|y - \mathbb{E}[y \mid I = i]\|_2^2 \mid I = i].$$

Substituting $a = W_{i,:}^\top$ and summing gives $\mathcal{E}(W) = \frac{1}{2} \sum_{i=1}^V \pi_i \|W_{i,:} - W_{i,:}^*\|_2^2 = \frac{1}{2} \text{Tr}((W - W^*)^\top D (W - W^*))$. ■

C.2. Asymptotics of the normalizing constant

Let $\zeta(z) = \sum_{i=1}^{\infty} i^{-z}$ be the Riemann zeta function.

Lemma 7 (Asymptotics of normalization factor z) As $V \rightarrow \infty$,

$$z = H_V^{(\alpha)} = \begin{cases} \frac{V^{1-\alpha}}{1-\alpha} + O(1), & 0 < \alpha < 1, \\ \log V + O(1), & \alpha = 1, \\ \zeta(\alpha) + O(V^{1-\alpha}), & \alpha > 1. \end{cases}$$

C.3. Asymptotics of $\mathcal{E}(0)$

Lemma 8 (Asymptotics of $\mathcal{E}(0)$) Under Assumption 1, as $V \rightarrow \infty$,

$$\mathcal{E}(0) = \frac{1}{2} \frac{H_V^{(2\alpha)}}{(H_V^{(\alpha)})^2} = \begin{cases} \frac{(1-\alpha)^2}{2(1-2\alpha)} V^{-1} + O(V^{-2(1-\alpha)}), & 0 < \alpha < \frac{1}{2}, \\ \frac{\log V}{8V} + O(V^{-1}), & \alpha = \frac{1}{2}, \\ \frac{(1-\alpha)^2 \zeta(2\alpha)}{2} V^{-2(1-\alpha)} (1 + o(1)), & \frac{1}{2} < \alpha < 1, \\ \frac{\zeta(2)}{2(\log V)^2} + O((\log V)^{-3}), & \alpha = 1, \\ \frac{\zeta(2\alpha)}{2\zeta(\alpha)^2} + O(V^{1-\alpha}), & \alpha > 1. \end{cases}$$

Appendix D. Stochastic Gradient Analysis

D.1. Stochastic gradient moments

Proposition 9 (Expectation and covariance of the stochastic gradient) *Let $v \in \mathbb{R}^{V^2}$ denote the row-vectorized stochastic gradient associated with a single sampled pair (I, J) , partitioned into V consecutive blocks $v^{(i)} \in \mathbb{R}^V$. Define $q_i := \mathbb{E}[e_J | I = i] = W_{i,:}^*$, and its j th entry be q_{ij} . Then*

$$\mathbb{E}[v]^{(i)} = \pi_i(W_{i,:} - q_i).$$

The second-moment matrix $\mathbb{E}[vv^\top]$ is block diagonal with

$$\mathbb{E}[vv^\top]_{i,i} = \sum_{j=1}^V \pi_{ij}(W_{i,:} - e_j)(W_{i,:} - e_j)^\top, \quad \mathbb{E}[vv^\top]_{i,j} = 0 \quad (i \neq j).$$

The diagonal entries of the covariance diagonal blocks satisfy

$$[\Sigma_{i,i}]_{jj} = \pi_i q_{ij}(1 - q_{ij}) + \pi_i(1 - \pi_i)(W_{ij} - q_{ij})^2, \quad (4)$$

separating the irreducible sampling noise $\pi_i q_{ij}(1 - q_{ij})$ from the state-dependent contribution $\pi_i(1 - \pi_i)(W_{ij} - q_{ij})^2$.

Proof The row-vectorized gradient for one-hot inputs $x = e_I, y = e_J$ is

$$v(W; I, J) = \overline{\text{vec}}(e_I(W_{I,:} - e_J)^\top) = e_I \otimes (W_{I,:} - e_J) \in \mathbb{R}^{V^2}.$$

The i -th block (V entries corresponding to row i) is

$$v^{(i)} = \mathbf{1}_{\{I=i\}}(W_{i,:} - e_J).$$

Taking expectation we get:

$$\mathbb{E}[v^{(i)}] = \mathbb{P}(I = i) \mathbb{E}[W_{i,:} - e_J | I = i] = \pi_i(W_{i,:} - q_i).$$

For the (i, i) diagonal block:

$$\begin{aligned} \mathbb{E}[v^{(i)}(v^{(i)})^\top] &= \mathbb{E}[\mathbf{1}_{\{I=i\}}(W_{i,:} - e_J)(W_{i,:} - e_J)^\top] \\ &= \pi_i \sum_{j=1}^V \pi_{j|i}(W_{i,:} - e_j)(W_{i,:} - e_j)^\top = \sum_{j=1}^V \pi_{ij}(W_{i,:} - e_j)(W_{i,:} - e_j)^\top. \end{aligned}$$

For the (i, k) off-diagonal block with $i \neq k$:

$$\mathbb{E}[v^{(i)}(v^{(k)})^\top] = \mathbb{E}[\mathbf{1}_{\{I=i\}}\mathbf{1}_{\{I=k\}}(\cdots)] = 0,$$

since I cannot equal both i and k simultaneously.

Next we calculate the covariance diagonal entries which will be used in 10. The (j, j) -th diagonal entry of $\Sigma_{i,i} = \mathbb{E}[v^{(i)}(v^{(i)})^\top] - \mathbb{E}[v^{(i)}]\mathbb{E}[v^{(i)}]^\top$ is

$$\begin{aligned} [\Sigma_{i,i}]_{jj} &= \sum_{k=1}^V \pi_{ik}(W_{ij} - \mathbf{1}_{\{k=j\}})^2 - \pi_i^2(W_{ij} - q_{ij})^2 \\ &= \pi_i[q_{ij}(W_{ij} - 1)^2 + (1 - q_{ij})W_{ij}^2] - \pi_i^2(W_{ij} - q_{ij})^2. \end{aligned}$$

Expanding and simplifying with $W_{ij} = q_{ij} + \Delta_{ij}$:

$$\begin{aligned} q_{ij}(W_{ij} - 1)^2 + (1 - q_{ij})W_{ij}^2 &= q_{ij}(1 - q_{ij}) + (W_{ij} - q_{ij})^2 \\ &= q_{ij}(1 - q_{ij}) + \Delta_{ij}^2. \end{aligned}$$

Therefore

$$[\Sigma_{i,i}]_{jj} = \pi_i q_{ij}(1 - q_{ij}) + \pi_i \Delta_{ij}^2 - \pi_i^2 \Delta_{ij}^2 = \pi_i q_{ij}(1 - q_{ij}) + \pi_i(1 - \pi_i) \Delta_{ij}^2. \quad \blacksquare$$

D.2. Coordinate recursion

Lemma 10 (Coordinate recursion and conditional second moment) *For every $t \geq 0$ and every pair (i, j) ,*

$$\Delta_{ij,t+1} = \phi_i \Delta_{ij,t} - \eta \epsilon_{ij,t},$$

where $\epsilon_{ij,t} := [G_t]_{ij} - \pi_i \Delta_{ij,t}$ satisfies $\mathbb{E}[\epsilon_{ij,t} \mid \mathcal{F}_t] = 0$ and

$$\mathbb{E}[\epsilon_{ij,t}^2 \mid \mathcal{F}_t] = \frac{1}{b} [\pi_i q_{ij}(1 - q_{ij}) + \pi_i(1 - \pi_i) \Delta_{ij,t}^2].$$

Proof For one sample (I, J) , the (i, j) coordinate of the sample gradient is $g_{ij}(W; I, J) = \mathbf{1}_{\{I=i\}}(W_{ij} - \mathbf{1}_{\{J=j\}})$. Conditioning on W yields

$$\mathbb{E}[g_{ij}(W; I, J) \mid W] = \pi_i(W_{ij} - q_{ij}).$$

Hence $\mathbb{E}[[G_t]_{ij} \mid \mathcal{F}_t] = \pi_i \Delta_{ij,t}$, and $\mathbb{E}[\epsilon_{ij,t} \mid \mathcal{F}_t] = 0$. The SGD update becomes $\Delta_{ij,t+1} = \phi_i \Delta_{ij,t} - \eta \epsilon_{ij,t}$.

For the conditional second moment, the b terms in the mini-batch average are conditionally i.i.d., so the conditional variance of the batch average is the one-sample conditional variance divided by b . For one sample,

$$\mathbb{E}[g_{ij}(W; I, J)^2 \mid W] = \pi_i [q_{ij}(W_{ij} - 1)^2 + (1 - q_{ij})W_{ij}^2].$$

Writing $W_{ij} = q_{ij} + \Delta_{ij}$ gives $q_{ij}(W_{ij} - 1)^2 + (1 - q_{ij})W_{ij}^2 = q_{ij}(1 - q_{ij}) + \Delta_{ij}^2$. Therefore,

$$\text{Var}(g_{ij}(W; I, J) \mid W) = \pi_i q_{ij}(1 - q_{ij}) + \pi_i(1 - \pi_i) \Delta_{ij}^2.$$

Dividing by b gives the stated formula. \blacksquare

D.3. Mean trajectory and bias

Proposition 11 (Mean trajectory is deterministic GD) *Let W_t^{gd} be the full-gradient descent iterates defined by $W_{t+1}^{\text{gd}} = W_t^{\text{gd}} - \eta D(W_t^{\text{gd}} - W^*)$ with $W_0^{\text{gd}} = W_0$. Then for every $t \geq 0$, $\mathbb{E}[W_t] = W_t^{\text{gd}}$. Moreover, for every (i, j) ,*

$$\mathbb{E}[\Delta_{ij,t}] = \phi_i^t \Delta_{ij,0}.$$

Let $\bar{\Delta}_{ij,T} = \frac{1}{T} \sum_{t=0}^{T-1} \Delta_{ij,t}$, then

$$\mathbb{E}[\bar{\Delta}_{ij,T}] = \frac{1}{T} \sum_{t=0}^{T-1} \phi_i^t \Delta_{ij,0} = \frac{1 - \phi_i^T}{T(1 - \phi_i)} \Delta_{ij,0}.$$

Proof Taking conditional expectation in Lemma 10 gives

$$\mathbb{E}[\Delta_{ij,t+1} \mid \mathcal{F}_t] = \phi_i \Delta_{ij,t}.$$

Taking expectation once more yields

$$\mathbb{E}[\Delta_{ij,t+1}] = \phi_i \mathbb{E}[\Delta_{ij,t}],$$

and induction gives $\mathbb{E}[\Delta_{ij,t}] = \phi_i^t \Delta_{ij,0}$. This is exactly the coordinate recursion of deterministic gradient descent. Averaging over $t = 0, \dots, T-1$ and using the geometric series formula:

$$\mathbb{E}[\bar{\Delta}_{ij,T}] = \frac{1}{T} \sum_{t=0}^{T-1} \phi_i^t \Delta_{ij,0} = \frac{1 - \phi_i^T}{T(1 - \phi_i)} \Delta_{ij,0}. \quad \blacksquare$$

Proposition 12 (Bias term) Define $\Gamma_i(T) := \left(\frac{1 - \phi_i^T}{1 - \phi_i}\right)^2$. Then

$$\mathcal{E}_{\text{bias}}(T) := \mathcal{L}(\mathbb{E}[\bar{W}_T]) - \mathcal{L}(W^*) = \frac{1}{2T^2} \sum_{i=1}^V \pi_i \Gamma_i(T) \sum_{j=1}^V \Delta_{ij,0}^2.$$

Proof By Proposition 11,

$$\mathbb{E}[\bar{\Delta}_{ij,T}] = \frac{1 - \phi_i^T}{T(1 - \phi_i)} \Delta_{ij,0}.$$

Squaring and summing with weights $\pi_i/2$:

$$\begin{aligned} \mathcal{E}_{\text{bias}}(T) &= \frac{1}{2} \sum_{i=1}^V \pi_i \sum_{j=1}^V (\mathbb{E}[\bar{\Delta}_{ij,T}])^2 = \frac{1}{2} \sum_{i=1}^V \pi_i \sum_{j=1}^V \left(\frac{1 - \phi_i^T}{T(1 - \phi_i)}\right)^2 \Delta_{ij,0}^2 \\ &= \frac{1}{2T^2} \sum_{i=1}^V \pi_i \Gamma_i(T) \sum_{j=1}^V \Delta_{ij,0}^2. \quad \blacksquare \end{aligned}$$

D.4. Second moments and temporal covariances

Proposition 13 (Second moments) Define $a_i := \phi_i^2 + \frac{\eta^2}{b} \pi_i(1 - \pi_i)$ and $c_{ij} := \frac{\eta^2}{b} \pi_i q_{ij}(1 - q_{ij})$. If $a_i < 1$, then

$$\mathbb{E}[\Delta_{ij,t}^2] = a_i^t \Delta_{ij,0}^2 + \frac{c_{ij}}{1 - a_i} (1 - a_i^t).$$

Proof Squaring the recursion of Lemma 10: $\Delta_{ij,t+1}^2 = \phi_i^2 \Delta_{ij,t}^2 - 2\eta\phi_i \Delta_{ij,t} \epsilon_{ij,t} + \eta^2 \epsilon_{ij,t}^2$. Taking conditional expectation and using $\mathbb{E}[\epsilon_{ij,t} \mid \mathcal{F}_t] = 0$ to get:

$$\mathbb{E}[\Delta_{ij,t+1}^2 \mid \mathcal{F}_t] = a_i \Delta_{ij,t}^2 + c_{ij}.$$

Taking expectation gives the linear recursion $\mathbb{E}[\Delta_{ij,t+1}^2] = a_i \mathbb{E}[\Delta_{ij,t}^2] + c_{ij}$. Since $a_i < 1$, solving yields the stated formula. \blacksquare

Lemma 14 (Cross-time covariance) For $0 \leq s < t$, $\text{Cov}(\Delta_{ij,t}, \Delta_{ij,s}) = \phi_i^{t-s} \text{Var}(\Delta_{ij,s})$. Consequently,

$$\text{Var}(\bar{\Delta}_{ij,T}) = \frac{1}{T^2} \sum_{s=0}^{T-1} w_{i,s} \text{Var}(\Delta_{ij,s}),$$

where $w_{i,s} := 1 + 2 \sum_{k=1}^{T-1-s} \phi_i^k = \frac{1+\phi_i}{1-\phi_i} - \frac{2\phi_i^{T-s}}{1-\phi_i}$.

Proof Unrolling the recursion from time s to time t gives

$$\Delta_{ij,t} = \phi_i^{t-s} \Delta_{ij,s} - \eta \sum_{k=s}^{t-1} \phi_i^{t-1-k} \epsilon_{ij,k}.$$

Multiplying by $\Delta_{ij,s}$ and taking expectations yields

$$\mathbb{E}[\Delta_{ij,t} \Delta_{ij,s}] = \phi_i^{t-s} \mathbb{E}[\Delta_{ij,s}^2] - \eta \sum_{k=s}^{t-1} \phi_i^{t-1-k} \mathbb{E}[\epsilon_{ij,k} \Delta_{ij,s}].$$

For every $k \geq s$, the variable $\Delta_{ij,s}$ is \mathcal{F}_k -measurable, hence

$$\mathbb{E}[\epsilon_{ij,k} \Delta_{ij,s}] = \mathbb{E}[\mathbb{E}[\epsilon_{ij,k} | \mathcal{F}_k] \Delta_{ij,s}] = 0.$$

Therefore

$$\mathbb{E}[\Delta_{ij,t} \Delta_{ij,s}] = \phi_i^{t-s} \mathbb{E}[\Delta_{ij,s}^2].$$

Since $\mathbb{E}[\Delta_{ij,t}] \mathbb{E}[\Delta_{ij,s}] = \phi_i^{t-s} (\mathbb{E}[\Delta_{ij,s}])^2$ by Proposition 11, subtraction gives

$$\text{Cov}(\Delta_{ij,t}, \Delta_{ij,s}) = \phi_i^{t-s} \text{Var}(\Delta_{ij,s}).$$

The variance of the average follows from

$$\begin{aligned} \text{Var}\left(\frac{1}{T} \sum_{t=0}^{T-1} \Delta_{ij,t}\right) &= \frac{1}{T^2} \left(\sum_{s=0}^{T-1} \text{Var}(\Delta_{ij,s}) + 2 \sum_{0 \leq s < t \leq T-1} \text{Cov}(\Delta_{ij,t}, \Delta_{ij,s}) \right) \\ &= \frac{1}{T^2} \sum_{s=0}^{T-1} \left(1 + 2 \sum_{k=1}^{T-1-s} \phi_i^k \right) \text{Var}(\Delta_{ij,s}). \end{aligned}$$

For the closed form of $w_{i,s}$, note that

$$w_{i,s} = 1 + 2 \sum_{k=1}^{T-1-s} \phi_i^k = \frac{1 + \phi_i}{1 - \phi_i} - \frac{2\phi_i^{T-s}}{1 - \phi_i}. \square$$

■

Corollary 15 (Exact variance of Polyak-averaged coordinate) Define $A_i := \frac{1+\phi_i}{1-\phi_i}$ and $B_i := \frac{2\phi_i}{1-\phi_i}$. Then for every (i, j) ,

$$\text{Var}(\bar{\Delta}_{ij,T}) = \frac{\Delta_{ij,0}^2}{T^2} (S_i^{(a)}(T) - \Gamma_i(T)) + \frac{c_{ij}}{(1-a_i)T^2} S_i^{(1-a)}(T),$$

where $w_{i,s} = A_i - B_i\phi_i^{T-1-s}$, and

$$S_i^{(1-a)}(T) = \left(A_i T - B_i \frac{1-\phi_i^T}{1-\phi_i} \right) - S_i^{(a)}(T).$$

Proof Substituting the variance formula from Proposition 13 into Lemma 14 gives

$$\text{Var}(\bar{\Delta}_{ij,T}) = \frac{1}{T^2} \sum_{s=0}^{T-1} w_{i,s} \left[(a_i^s - \phi_i^{2s}) \Delta_{ij,0}^2 + \frac{c_{ij}}{1-a_i} (1-a_i^s) \right].$$

It therefore remains to evaluate $\sum_{s=0}^{T-1} w_{i,s} \phi_i^{2s}$. Using the definition of $w_{i,s}$,

$$\begin{aligned} \sum_{s=0}^{T-1} w_{i,s} \phi_i^{2s} &= \sum_{s=0}^{T-1} \phi_i^{2s} + 2 \sum_{0 \leq s < t \leq T-1} \phi_i^{t-s} \phi_i^{2s} \\ &= \sum_{s=0}^{T-1} \phi_i^{2s} + 2 \sum_{0 \leq s < t \leq T-1} \phi_i^{s+t} = \left(\sum_{t=0}^{T-1} \phi_i^t \right)^2 = \Gamma_i(T). \end{aligned}$$

This proves the variance formula.

For the closed forms, note that $w_{i,s} = A_i - B_i\phi_i^{T-1-s}$. Therefore,

$$S_i^{(a)}(T) = \sum_{s=0}^{T-1} (A_i - B_i\phi_i^{T-1-s}) a_i^s = A_i \frac{1-a_i^T}{1-a_i} - B_i \sum_{s=0}^{T-1} \phi_i^{T-1-s} a_i^s.$$

Also,

$$\sum_{s=0}^{T-1} w_{i,s} = A_i T - B_i \sum_{s=0}^{T-1} \phi_i^{T-1-s} = A_i T - B_i \frac{1-\phi_i^T}{1-\phi_i},$$

whence

$$S_i^{(1-a)}(T) = \sum_{s=0}^{T-1} w_{i,s} - S_i^{(a)}(T) = \left(A_i T - B_i \frac{1-\phi_i^T}{1-\phi_i} \right) - S_i^{(a)}(T). \square$$

Appendix E. Proof of the Bias–Variance Decomposition (Theorem 3)

Proof By Lemma 2, $\mathcal{L}(W) - \mathcal{L}(W^*) = \frac{1}{2} \|D^{1/2}(W - W^*)\|_F^2$. Therefore $\mathbb{E}[\mathcal{L}(\bar{W}_T)] - \mathcal{L}(W^*) = \frac{1}{2} \mathbb{E}[\|D^{1/2} \bar{\Delta}_T\|_F^2]$. Decompose $\bar{\Delta}_T = \mathbb{E}[\bar{\Delta}_T] + (\bar{\Delta}_T - \mathbb{E}[\bar{\Delta}_T])$. Expanding the squared Frobenius

norm and using $\mathbb{E}[\bar{\Delta}_T - \mathbb{E}[\bar{\Delta}_T]] = 0$, the cross term vanishes:

$$\begin{aligned} \mathbb{E}[\mathcal{L}(\bar{W}_T)] - \mathcal{L}(W^*) &= \frac{1}{2} \|D^{1/2} \mathbb{E}[\bar{\Delta}_T]\|_F^2 + \frac{1}{2} \mathbb{E}[\|D^{1/2}(\bar{\Delta}_T - \mathbb{E}[\bar{\Delta}_T])\|_F^2] \\ &= \frac{1}{2} \sum_{i=1}^V \pi_i \sum_{j=1}^V (\mathbb{E}[\bar{\Delta}_{ij,T}])^2 + \frac{1}{2} \sum_{i=1}^V \pi_i \sum_{j=1}^V \text{Var}(\bar{\Delta}_{ij,T}). \end{aligned}$$

The first term is $\mathcal{E}_{\text{bias}}(T)$ by Proposition 12 and the second term is expressed explicitly by Corollary 15. Adding them yields the total-risk formula. Under initialization $W_0 = 0$ and assumption 1 we get $\sum_{j=1}^V \Delta_{ij,0}^2 = 2\mathcal{E}(0)$ and $\sum_{j=1}^V c_{ij} = \frac{\eta^2}{b} \pi_i (1 - 2\mathcal{E}(0))$. ■

Appendix F. Proof of the Bias Scaling Law (Theorem 4)

Proof Denote $\kappa := \eta/z$ and $\delta_i := \eta\pi_i = \kappa i^{-\alpha}$ the effective learning rate of the i -th row. Define

$$b_T(x) := \frac{1 - (1-x)^T}{Tx}, \quad 0 < x \leq 1. \quad (5)$$

Then the relative bias term from Proposition 12 can be written as

$$r_{\text{bias}}(T) = \sum_{i=1}^V \pi_i b_T(\delta_i)^2.$$

We organize the proof into three steps.

Step 1: Two-sided bound on $b_T(x)^2$. We prove

$$b_T(x)^2 \asymp \min \left\{ 1, \frac{1}{T^2 x^2} \right\}. \quad (6)$$

Upper bound. The elementary inequality $1 - (1-x)^T \leq \min\{1, Tx\}$ gives

$$b_T(x)^2 \leq \min \left\{ 1, \frac{1}{T^2 x^2} \right\}.$$

Lower bound. There are two cases. If $Tx \leq 1$, then

$$1 - (1-x)^T \geq Tx - \binom{T}{2} x^2 \geq Tx - \frac{T^2 x^2}{2} \geq \frac{Tx}{2},$$

so $b_T(x)^2 \geq 1/4$. If $Tx \geq 1$, then $(1-x)^T \leq e^{-Tx} \leq e^{-1}$, hence

$$1 - (1-x)^T \geq 1 - e^{-1}, \quad b_T(x)^2 \geq \frac{(1 - e^{-1})^2}{T^2 x^2}.$$

This proves (6).

Step 2: Splitting into active and inactive rows. Applying (6) to $x = \delta_i = \kappa i^{-\alpha}$ yields

$$r_{\text{bias}}(T) \asymp \frac{1}{z} \sum_{i=1}^V i^{-\alpha} \min \left\{ 1, \frac{i^{2\alpha}}{\kappa^2 T^2} \right\}.$$

If $i \leq i_{\eta,*}(T)$, then $T\delta_i = \kappa T i^{-\alpha} \geq 1$, so the summand is comparable to $i^\alpha / (\kappa^2 T^2)$. If $i > i_{\eta,*}(T)$, then $T\delta_i < 1$, so the summand is comparable to $i^{-\alpha}$. Therefore

$$r_{\text{bias}}(T) \asymp \frac{1}{z} \left(\sum_{i > i_{\eta,*}(T)} i^{-\alpha} + \frac{1}{\kappa^2 T^2} \sum_{i \leq i_{\eta,*}(T)} i^\alpha \right). \quad (7)$$

For the displayed scaling law with $\kappa = 1$, write $i_*(T) = \min\{V, \lfloor T^{1/\alpha} \rfloor\}$ and $T_{\text{crit}} = V^\alpha$. For general fixed $\kappa \in (0, 1]$, the same expressions hold after replacing T by κT and T_{crit} by $\kappa^{-1} V^\alpha$.

Step 3: Regime-wise asymptotics via integral test. We use the standard integral-test estimates: for $\beta > -1$,

$$\sum_{i \leq m} i^\beta \asymp m^{\beta+1};$$

for $\beta \neq 1$,

$$\sum_{i > m} i^{-\beta} \asymp \begin{cases} V^{1-\beta} - m^{1-\beta}, & 0 < \beta < 1, \\ m^{1-\beta}, & \beta > 1, \end{cases}$$

and for $\beta = 1$, $\sum_{i > m}^V i^{-1} \asymp \log(V/m)$.

Case 1: $0 < \alpha < 1$. If $T \leq T_{\text{crit}}$, then $i_*(T) \asymp T^{1/\alpha}$ and

$$\sum_{i > i_*}^V i^{-\alpha} = \frac{1}{1-\alpha} (V^{1-\alpha} - T^{(1-\alpha)/\alpha}), \quad \frac{1}{T^2} \sum_{i \leq i_*} i^\alpha = \frac{1}{\alpha+1} \frac{T^{(\alpha+1)/\alpha}}{T^2} = \frac{T^{(1-\alpha)/\alpha}}{\alpha+1}.$$

Adding the two contributions:

$$r_{\text{bias}}(T) \asymp \frac{1}{z} \left(\frac{V^{1-\alpha}}{1-\alpha} - \frac{2\alpha}{1-\alpha^2} T^{(1-\alpha)/\alpha} \right) \asymp 1 - \frac{2\alpha}{1+\alpha} \frac{T^{(1-\alpha)/\alpha}}{V^{1-\alpha}}.$$

If $T \geq V^\alpha$, then $i_*(T) = V$ and

$$r_{\text{bias}}(T) \asymp \frac{1}{z T^2} \sum_{i=1}^V i^\alpha \asymp \frac{V^{\alpha+1}}{V^{1-\alpha} T^2} = \frac{V^{2\alpha}}{T^2}.$$

Case 2: $\alpha = 1$. If $T \leq T_{\text{crit}}$, then $i_*(T) \asymp T$ and

$$r_{\text{bias}}(T) \asymp \frac{1}{\log V} \left(\log \frac{V}{T} + \frac{1}{2} \right) = 1 - \frac{\log T}{\log V} + \frac{1}{2 \log V} \asymp 1 - \frac{\log T}{\log V}.$$

If $T \geq V$, then $i_*(T) = V$ and

$$r_{\text{bias}}(T) \asymp \frac{1}{z T^2} \sum_{i=1}^V i \asymp \frac{V^2}{T^2 \log V}.$$

Case 3: $\alpha > 1$. If $T \leq T_{\text{crit}}$, then $i_*(T) \asymp T^{1/\alpha}$ and, since $z \asymp \zeta(\alpha) = \Theta(1)$,

$$r_{\text{bias}}(T) \asymp T^{(1-\alpha)/\alpha}.$$

If $T \geq V^\alpha$, then $i_*(T) = V$ and

$$r_{\text{bias}}(T) \asymp \frac{1}{T^2} \sum_{i=1}^V i^\alpha \asymp \frac{V^{\alpha+1}}{T^2}.$$

This completes the proof. ■

Appendix G. Proof of the Variance Scaling Law (Theorem 4)

Proof Write $\mathcal{L}_{\text{var}}(T) = F_V(T) + G_V(T)$, where

$$F_V(T) := \frac{\mathcal{E}(0)}{T^2} \sum_{i=1}^V \pi_i (S_i^{(a)}(T) - \Gamma_i(T))$$

and

$$G_V(T) := \frac{\eta^2(1 - 2\mathcal{E}(0))}{2bT^2} \sum_{i=1}^V \frac{\pi_i^2}{1 - a_i} S_i^{(1-a)}(T).$$

Recall the row-wise effective learning rate $\delta_i = \eta\pi_i = \kappa i^{-\alpha}$, where $\kappa := \eta/z$, and the second-moment contraction factor $a_i = \phi_i^2 + \frac{\eta^2}{b} \pi_i(1 - \pi_i)$.

We organize the proof into six steps.

Step 1: Reduction to a one-row quantity. Define

$$u_i(T) := \frac{\delta_i^2}{(1 - a_i)T^2} S_i^{(1-a)}(T) \quad (1 \leq i \leq V). \quad (8)$$

Since $\delta_i = \eta\pi_i$, the second term is exactly

$$G_V(T) = \frac{1 - 2\mathcal{E}(0)}{2b} \sum_{i=1}^V u_i(T). \quad (9)$$

We next show that $F_V(T)$ is controlled by the same sum. Fix i . Because $a_i \geq \phi_i^2$ and all weights $w_{i,s}$ are nonnegative,

$$0 \leq a_i^s - \phi_i^{2s} = (a_i - \phi_i^2) \sum_{r=0}^{s-1} a_i^r \phi_i^{2(s-1-r)} \leq (a_i - \phi_i^2) \sum_{r=0}^{s-1} a_i^r = \frac{a_i - \phi_i^2}{1 - a_i} (1 - a_i^s).$$

After multiplying by $w_{i,s}$ and summing over s , we obtain

$$0 \leq S_i^{(a)}(T) - \Gamma_i(T) \leq \frac{a_i - \phi_i^2}{1 - a_i} S_i^{(1-a)}(T). \quad (10)$$

Using $a_i - \phi_i^2 = \eta^2 b^{-1} \pi_i (1 - \pi_i) \leq \eta^2 b^{-1} \pi_i$, we get

$$0 \leq \frac{\pi_i}{T^2} (S_i^{(a)}(T) - \Gamma_i(T)) \leq \frac{\eta^2 \pi_i^2}{b(1 - a_i) T^2} S_i^{(1-a)}(T) = \frac{1}{b} u_i(T).$$

Therefore

$$0 \leq F_V(T) \leq \frac{\mathcal{E}(0)}{b} \sum_{i=1}^V u_i(T). \quad (11)$$

Combining (9) and (11) gives

$$\frac{1 - 2\mathcal{E}(0)}{2b} \sum_{i=1}^V u_i(T) \leq \mathcal{L}_{\text{var}}(T) \leq \frac{1}{2b} \sum_{i=1}^V u_i(T). \quad (12)$$

By Lemma 8, $2\mathcal{E}(0) = H_V^{(2\alpha)} / (H_V^{(\alpha)})^2 \rightarrow 0$ for $\alpha \leq 1$, while for $\alpha > 1$, $2\mathcal{E}(0) \rightarrow \zeta(2\alpha) / \zeta(\alpha)^2 < 1$. Therefore (12) yields

$$\mathcal{L}_{\text{var}}(T) \asymp \frac{1}{b} \sum_{i=1}^V u_i(T). \quad (13)$$

It remains to estimate $\sum_i u_i(T)$.

Step 2: Comparison of $1 - a_i$ with δ_i . Since

$$a_i = (1 - \delta_i)^2 + \frac{\eta \delta_i}{b} (1 - \pi_i),$$

we have

$$1 - a_i = 2\delta_i - \delta_i^2 - \frac{\eta \delta_i}{b} (1 - \pi_i).$$

Because $0 < \delta_i \leq 1$ and $b \geq 2\eta$,

$$1 - a_i \geq 2\delta_i - \delta_i - \frac{\delta_i}{2} = \frac{\delta_i}{2}, \quad 1 - a_i \leq 2\delta_i. \quad (14)$$

Substituting into (8):

$$u_i(T) \asymp \frac{\delta_i}{T^2} S_i^{(1-a)}(T). \quad (15)$$

Step 3: Active rows ($\delta_i T \geq 1$). *Upper bound.* For every s ,

$$w_{i,s} = 1 + 2 \sum_{k=1}^{T-1-s} \phi_i^k \leq 1 + 2 \sum_{k=1}^{\infty} (1 - \delta_i)^k \leq \frac{2}{\delta_i},$$

so

$$S_i^{(1-a)}(T) \leq \sum_{s=0}^{T-1} \frac{2}{\delta_i} = \frac{2T}{\delta_i}. \quad (16)$$

Lower bound. Fix $s \in [T/4, T/2]$. Then $T - 1 - s \geq T/2 - 1$, hence

$$w_{i,s} \geq \sum_{k=0}^{\lfloor T/2 \rfloor - 1} (1 - \delta_i)^k = \frac{1 - (1 - \delta_i)^{\lfloor T/2 \rfloor}}{\delta_i} \geq \frac{1 - e^{-\delta_i \lfloor T/2 \rfloor}}{\delta_i} \geq \frac{c}{\delta_i},$$

where $c > 0$ is an absolute constant because $\delta_i T \geq 1$. Also, by (14),

$$s(1 - a_i) \geq \frac{T}{4} \cdot \frac{\delta_i}{2} \geq \frac{1}{8},$$

so

$$1 - a_i^s \geq 1 - e^{-s(1-a_i)} \geq 1 - e^{-1/8} =: c' > 0.$$

Since there are $\asymp T$ indices s in $[T/4, T/2]$, we conclude

$$S_i^{(1-a)}(T) \gtrsim \frac{T}{\delta_i}. \quad (17)$$

Combining (16), (17), and (15):

$$\delta_i T \geq 1 \implies u_i(T) \asymp \frac{1}{T}. \quad (18)$$

Step 4: Inactive rows ($\delta_i T \leq 1/2$). *Upper bound.* For all s ,

$$w_{i,s} \leq 1 + 2(T - 1 - s) \leq 2T.$$

Using $1 - a_i^s \leq s(1 - a_i) \lesssim s\delta_i$, we get

$$S_i^{(1-a)}(T) \leq \sum_{s=0}^{T-1} 2T \cdot C s \delta_i \lesssim \delta_i T^3. \quad (19)$$

Lower bound. Let again $s \in [T/4, T/2]$. If $0 \leq k \leq \lfloor T/2 \rfloor - 1$, then $k\delta_i \leq T\delta_i/2 \leq 1/4$, hence

$$(1 - \delta_i)^k \geq 1 - k\delta_i \geq \frac{3}{4}.$$

Therefore

$$w_{i,s} \geq \sum_{k=0}^{\lfloor T/2 \rfloor - 1} (1 - \delta_i)^k \gtrsim T.$$

Moreover, still for $s \in [T/4, T/2]$, we have by (14)

$$s(1 - a_i) \leq \frac{T}{2} \cdot 2\delta_i \leq \frac{1}{2}.$$

Hence the quadratic remainder estimate gives

$$1 - a_i^s \geq s(1 - a_i) - \binom{s}{2} (1 - a_i)^2 \geq \frac{1}{2} s(1 - a_i) \gtrsim s\delta_i \gtrsim T\delta_i.$$

Summing over $s \in [T/4, T/2]$, of which there are $\asymp T$:

$$S_i^{(1-a)}(T) \gtrsim \delta_i T^3. \quad (20)$$

Combining (19), (20), and (15):

$$\delta_i T \leq \frac{1}{2} \implies u_i(T) \asymp \delta_i^2 T = \kappa^2 T i^{-2\alpha}. \quad (21)$$

Step 5: Summation over i . Let

$$I_1 := \{i : \delta_i T \geq 1\} = \{i : i \leq i_{\kappa,*}(T)\}, \quad I_0 := \{i : \delta_i T \leq 1/2\} = \{i : i \geq 2^{1/\alpha} i_{\kappa,*}(T)\} \cap [V].$$

By (18) and (21),

$$\sum_{i \in I_1} u_i(T) \asymp \frac{|I_1|}{T} \asymp \frac{i_{\kappa,*}(T)}{T}, \quad \sum_{i \in I_0} u_i(T) \asymp \kappa^2 T \sum_{i \in I_0} i^{-2\alpha}. \quad (22)$$

Also, the upper bounds from Steps 3 and 4 imply that for every i ,

$$u_i(T) \lesssim \min \left\{ \frac{1}{T}, \kappa^2 T i^{-2\alpha} \right\}.$$

Therefore,

$$\sum_{i=1}^V u_i(T) \lesssim \frac{i_{\kappa,*}(T)}{T} + \kappa^2 T \sum_{i > i_{\kappa,*}(T)} i^{-2\alpha}. \quad (23)$$

For the lower bound, from (22):

$$\sum_{i=1}^V u_i(T) \gtrsim \frac{i_{\kappa,*}(T)}{T} + \kappa^2 T \sum_{i \geq 2^{1/\alpha} i_{\kappa,*}(T)} i^{-2\alpha}.$$

The middle block between $i_{\kappa,*}$ and $2^{1/\alpha} i_{\kappa,*}$ satisfies

$$\sum_{i=i_{\kappa,*}+1}^{\lfloor 2^{1/\alpha} i_{\kappa,*} \rfloor} i^{-2\alpha} \leq (\lfloor 2^{1/\alpha} i_{\kappa,*} \rfloor - i_{\kappa,*}) \cdot i_{\kappa,*}^{-2\alpha} \lesssim i_{\kappa,*}^{1-2\alpha}.$$

Multiplying by $\kappa^2 T$ and using $i_{\kappa,*}^\alpha \asymp \kappa T$:

$$\kappa^2 T \cdot i_{\kappa,*}^{1-2\alpha} = \kappa^2 T \frac{i_{\kappa,*}}{i_{\kappa,*}^{2\alpha}} \asymp \kappa^2 T \frac{i_{\kappa,*}}{\kappa^2 T^2} = \frac{i_{\kappa,*}}{T}.$$

Now decompose the tail:

$$\kappa^2 T \sum_{i > i_{\kappa,*}} i^{-2\alpha} = \kappa^2 T \sum_{i=i_{\kappa,*}+1}^{\lfloor 2^{1/\alpha} i_{\kappa,*} \rfloor} i^{-2\alpha} + \kappa^2 T \sum_{i > \lfloor 2^{1/\alpha} i_{\kappa,*} \rfloor} i^{-2\alpha} \leq C \frac{i_{\kappa,*}}{T} + \kappa^2 T \sum_{i \geq 2^{1/\alpha} i_{\kappa,*}} i^{-2\alpha}.$$

Therefore

$$\sum_{i=1}^V u_i(T) \gtrsim \frac{i_{\kappa,*}}{T} + \kappa^2 T \sum_{i \geq 2^{1/\alpha} i_{\kappa,*}} i^{-2\alpha} \gtrsim \frac{i_{\kappa,*}}{T} + \kappa^2 T \sum_{i > i_{\kappa,*}} i^{-2\alpha}, \quad (24)$$

because the only extra part in $\sum_{i > i_{\kappa,*}}$ is exactly the middle block whose contribution is $\lesssim i_{\kappa,*}/T$, already present.

Combining (23) and (24) with (13) proves

$$\mathcal{L}_{\text{var}}(T) \asymp \frac{1}{b} \left(\frac{i_{\kappa,*}(T)}{T} + \kappa^2 T \sum_{i > i_{\kappa,*}(T)} i^{-2\alpha} \right).$$

Step 6: Regime-wise asymptotics. We convert the active-block plus inactive-tail expression into the closed-form scalings of Theorem 4 in the normalized case $\kappa = 1$.

Subcritical regime $T \lesssim T_{\text{crit}}(V)$. Here $i_*(T) = \lfloor T^{1/\alpha} \rfloor$, so the active block contributes $i_*(T)/T = T^{1/\alpha-1}$. For the tail, the integral test gives

$$T \sum_{i>i_*} i^{-2\alpha} \asymp \begin{cases} \frac{TV^{1-2\alpha}}{1-2\alpha} - \frac{T i_*^{1-2\alpha}}{1-2\alpha} = \frac{TV^{1-2\alpha}}{1-2\alpha} - \frac{T^{1/\alpha-1}}{1-2\alpha}, & 0 < \alpha < \frac{1}{2}, \\ T(\log V - \log i_*) = T(\log V - 2 \log T), & \alpha = \frac{1}{2}, \\ \frac{T i_*^{1-2\alpha}}{2\alpha-1} - \frac{TV^{1-2\alpha}}{2\alpha-1} = \frac{T^{1/\alpha-1}}{2\alpha-1} - \frac{TV^{1-2\alpha}}{2\alpha-1}, & \alpha > \frac{1}{2}, \end{cases} \quad (25)$$

Adding the active block $T^{1/\alpha-1}$ (which absorbs into the leading constant) gives

$$\frac{i_*}{T} + T \sum_{i>i_*} i^{-2\alpha} \asymp \begin{cases} TV^{1-2\alpha} - c_\alpha T^{1/\alpha-1}, & 0 < \alpha < \frac{1}{2}, \\ T(1 + \log V - 2 \log T), & \alpha = \frac{1}{2}, \\ T^{1/\alpha-1} - c_\alpha TV^{1-2\alpha}, & \alpha > \frac{1}{2}, \end{cases} \quad (26)$$

with $c_\alpha > 0$ an explicit α -dependent constant.

We then use ratio comparison to select the dominant term.

For the $\alpha \neq 1/2$ cases, (26) is a difference $A_\alpha(T, V) - B_\alpha(T, V)$ of two positive terms. A direct computation gives

$$\frac{B_\alpha}{A_\alpha} = \begin{cases} \frac{c_\alpha T^{1/\alpha-1}}{TV^{1-2\alpha}} = c_\alpha (T/V^\alpha)^{1/\alpha-2}, & 0 < \alpha < \frac{1}{2}, \\ \frac{c_\alpha TV^{1-2\alpha}}{T^{1/\alpha-1}} = c_\alpha (T/V^\alpha)^{2-1/\alpha}, & \alpha > \frac{1}{2}. \end{cases}$$

In both cases the exponent $p_\alpha > 0$, so under the subcritical assumption $T \lesssim V^\alpha$ we have $B_\alpha/A_\alpha \asymp (T/V^\alpha)^{p_\alpha} = O(1)$. Hence

$$A_\alpha - B_\alpha = A_\alpha(1 - B_\alpha/A_\alpha) \asymp A_\alpha,$$

which identifies the dominant term as A_α .

For $\alpha = 1/2$, the bracket $1 + \log V - 2 \log T$ is dominated by $\log V$ since $\log T \lesssim \log V$, so the same principle applies with $A = T \log V$.

Combining these with the prefactor $1/b$ from (3) yields the subcritical scalings

$$\mathcal{L}_{\text{var}}(T) \asymp \frac{1}{b} \cdot \begin{cases} TV^{1-2\alpha} & 0 < \alpha < 1/2, \\ T \log V & \alpha = 1/2, \\ T^{1/\alpha-1} & \alpha > 1/2, \end{cases} \quad T \lesssim T_{\text{crit}}(V). \quad (27)$$

Supercritical regime $T \gg T_{\text{crit}}(V)$. Here $i_*(T) = V$, so the inactive tail in the surrogate is empty and only the active block remains:

$$\mathcal{L}_{\text{var}}(T) \asymp \frac{1}{b} \cdot \frac{V}{T}.$$

Finally, dividing (27) and the supercritical bound by $\mathcal{E}(0)$ and applying Lemma 8 to substitute the regime-wise value of $\mathcal{E}(0)$ yields the relative-variance formulas stated in Theorem 4. \blacksquare

Appendix H. Proof of the Total Scaling and Critical Batch Size

H.1. Proof of Corollary 5

Proof Using the decomposition

$$\frac{\mathbb{E}[\mathcal{E}(\bar{W}_T)]}{\mathcal{E}(0)} = r_{\text{bias}}(T) + r_{\text{var}}(T),$$

we combine Theorem 4 and the relative variance $r_{\text{var}}(T) = \mathcal{L}_{\text{var}}(T)/\mathcal{E}(0)$ from Theorem 4. This gives the general representation

$$\frac{\mathbb{E}[\mathcal{E}(\bar{W}_T)]}{\mathcal{E}(0)} \asymp \frac{1}{z} \left(\sum_{i>i_*} i^{-\alpha} + \frac{1}{T^2} \sum_{i \leq i_*} i^\alpha \right) + \frac{1}{b\mathcal{E}(0)} \left(\frac{i_*}{T} + T \sum_{i>i_*} i^{-2\alpha} \right).$$

For $T \gg T_{\text{crit}}(V)$, we have $i_*(T) = V$ and all tail sums vanish, yielding

$$\frac{\mathbb{E}[\mathcal{E}(\bar{W}_T)]}{\mathcal{E}(0)} \asymp \frac{1}{z} \cdot \frac{1}{T^2} \sum_{i=1}^V i^\alpha + \frac{V}{b\mathcal{E}(0)T}.$$

The regime-wise formulas follow by substituting:

- $z \asymp V^{1-\alpha}/(1-\alpha)$ for $\alpha < 1$, $z \asymp \log V$ for $\alpha = 1$, $z \asymp \zeta(\alpha)$ for $\alpha > 1$ (Lemma 7);
- $\sum_{i=1}^V i^\alpha \asymp V^{\alpha+1}/(\alpha+1)$;
- $\mathcal{E}(0)$ from Lemma 8.

For example, when $\alpha = 1$:

$$\frac{1}{zT^2} \sum_{i=1}^V i \asymp \frac{V^2}{T^2 \log V}, \quad \frac{V}{b\mathcal{E}(0)T} \asymp \frac{V(\log V)^2}{bT}.$$

The subcritical formulas are obtained analogously by retaining both the tail and head sums before specializing. \blacksquare

H.2. Proof of Proposition 6 (Critical batch size)

Proof For $T \geq V^\alpha$, $i_*(T) = \min\{V, \lfloor T^{1/\alpha} \rfloor\} = V$ and all tail terms vanish. Hence the total relative excess risk from Corollary 5 reduces to

$$\frac{\mathbb{E}[\mathcal{E}(\bar{W}_T)]}{\mathcal{E}(0)} \asymp \underbrace{\frac{1}{z} \cdot \frac{1}{T^2} \sum_{i=1}^V i^\alpha}_{\text{bias}} + \underbrace{\frac{1}{b\mathcal{E}(0)} \cdot \frac{V}{T}}_{\text{variance}}.$$

The critical batch size $b_{\text{crit}}(T)$ is meant to balance these two contributions:

$$\frac{1}{zT^2} \sum_{i=1}^V i^\alpha \asymp \frac{V}{b_{\text{crit}}(T) \mathcal{E}(0) T}.$$

Rearranging gives:

$$b_{\text{crit}}(T) \asymp \frac{zVT}{\mathcal{E}(0) \sum_{i=1}^V i^\alpha}.$$

Since $\sum_{i=1}^V i^\alpha \asymp V^{\alpha+1}/(\alpha+1)$, this simplifies to

$$b_{\text{crit}}(T) \asymp \frac{z}{\mathcal{E}(0)} \cdot \frac{T}{V^\alpha}.$$

The regime-wise formulas follow by substituting the asymptotics of z (Lemma 7) and $\mathcal{E}(0)$ (Lemma 8). For example:

Case $0 < \alpha < 1/2$: $z \asymp V^{1-\alpha}/(1-\alpha)$ and $\mathcal{E}(0) \asymp (1-\alpha)^2 V^{-1}/(2(1-2\alpha))$, so

$$b_{\text{crit}}(T) \asymp \frac{V^{1-\alpha}}{V^{-1}} \cdot \frac{T}{V^\alpha} = TV^{2-2\alpha}.$$

Case $\alpha = 1$: $z \asymp \log V$ and $\mathcal{E}(0) \asymp \zeta(2)/(2(\log V)^2)$, so

$$b_{\text{crit}}(T) \asymp \frac{\log V}{(\log V)^{-2}} \cdot \frac{T}{V} = \frac{T(\log V)^3}{V}.$$

Case $\alpha > 1$: $z \asymp \zeta(\alpha)$ and $\mathcal{E}(0) \asymp \zeta(2\alpha)/(\zeta(\alpha)^2)$, so

$$b_{\text{crit}}(T) \asymp \frac{\zeta(\alpha)}{\zeta(2\alpha)/\zeta(\alpha)^2} \cdot \frac{T}{V^\alpha} \asymp \frac{T}{V^\alpha}.$$

For the compute-based formulas, substitute $T = N/b$ into $b_{\text{crit}} \asymp C \cdot T/V^\alpha$:

$$b_{\text{crit}} \asymp C \cdot \frac{N}{b_{\text{crit}} V^\alpha}, \quad b_{\text{crit}}^2 \asymp \frac{CN}{V^\alpha}, \quad b_{\text{crit}} \asymp \sqrt{\frac{CN}{V^\alpha}},$$

where $C = z/\mathcal{E}(0)$ depends on α as computed above. ■