

ESLM: RISK-AVERSE SELECTIVE LANGUAGE MODELING WITH HIERARCHICAL BATCH SELECTION

Anonymous authors

Paper under double-blind review

ABSTRACT

Large language model pretraining is compute-intensive, yet many tokens contribute marginally to learning, resulting in inefficiency. We introduce Efficient Selective Language Modeling (ESLM), an online, risk-aware batch selection algorithm that improves training efficiency and distributional robustness. ESLM operates in two phases: (i) instance-level selection via a shallow early-exit model pass that computes proxy per-instance statistics (e.g., loss or entropy) and retains data points using value-at-risk thresholding; and (ii) loss shaping with token-level selection via risk-aware thresholding on per-token scores. This data-centric mechanism reshapes the training objective, prioritizing high-risk tokens and eliminating redundant gradient computation. We frame ESLM as a bilevel game: the model competes with a masking adversary that selects worst-case token subsets under a constrained thresholding rule. In the loss-based setting, ESLM recovers conditional value-at-risk loss minimization, linking selective pretraining to distributionally robust optimization. We extend our approach to ADA-ESLM, which adaptively tunes the selection confidence during training. Experiments on GPT-2 pretraining show that ESLM significantly reduces training FLOPs while maintaining or improving perplexity and downstream performance compared to baselines. Our approach also scales across model sizes, pretraining corpora, and integrates naturally with knowledge distillation.

1 INTRODUCTION

The growing scale of large language models (LLMs) has brought substantial improvements in downstream performance at the expense of significantly higher pretraining costs (Brown et al., 2020). Training LLMs is notoriously compute-intensive, requiring massive GPU resources and often processing billions of tokens uniformly. Yet, many tokens, e.g., predictable or low-entropy ones, contribute little to model learning (Hüllermeier and Waegeman, 2021). Standard causal language modeling (CLM) treats all tokens equally in the loss, allocating compute uniformly to frequent or trivial tokens and more informative ones, leading to inefficient training and suboptimal use of resources (Lin et al., 2024).

Efforts to improve pretraining efficiency span architectural advances (Dao, 2023), token pruning (Hou et al., 2022), and increasingly, data-centric strategies (Xia et al., 2024; Wang et al., 2024). Among these, data-centric approaches show particular promise for sample efficiency through selective weighting or filtering of training examples (Katharopoulos and Fleuret, 2018). However, existing methods often rely on reference models or heuristics (Lin et al., 2024), operate at the sequence level (Yu et al., 2024), or require offline scoring (Xie et al., 2023b; Wettig et al., 2024). These choices limit adaptability and scaling of the model to massive web-scale corpora, and the ability to exploit token-level heterogeneity, despite being crucial for optimizing training dynamics and resource usage in LLM pretraining.

We address this gap with ESLM—*Efficient Selective Language Modeling*—a self-supervised data-centric framework that performs *online token-level batch selection* for efficient and robust pretraining. ESLM proceeds in two phases (see Figure 1): (i) a *proxy phase* that selects instances (sequences) via a shallow early-exit pass using only L model layers based on proxy statistics (e.g., per-instance loss or predictive entropy (Shannon, 1948)), keeping instances by (conditional-)value-at-risk (CVaR/VaR) thresholding; (ii) the training phase with risk-aware *token-level loss shaping*, which retains only the highest-risk tokens in the training objective via VaR thresholding over per-token statistics. This dynamic filtering shapes the training loss to emphasize uncertain or informative tokens, reducing redundant gradient updates and improving compute efficiency.

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

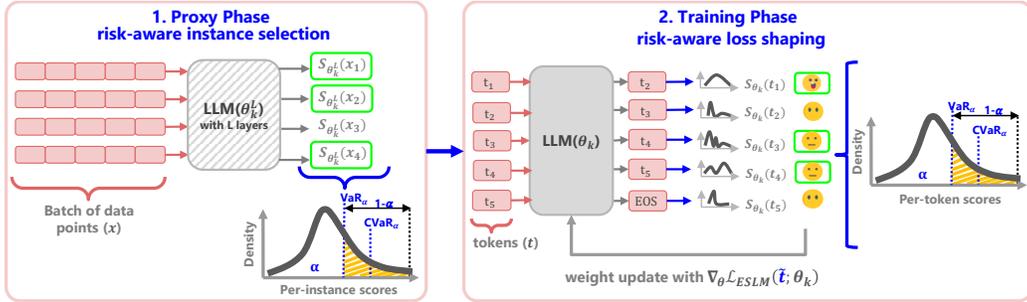


Figure 1: **ESLM illustration.** ESLM uses VaR thresholding on proxy per-instance scores to select high-risk data points, then applies token-level VaR to perform loss shaping. This reshapes the effective training distribution and loss by focusing computational resources on tokens with higher learning value.

Beyond its algorithmic simplicity, ESLM is grounded in a solid theoretical foundation. Its risk-aware selection mechanism can be viewed as a bilevel game in which the model competes with a constrained adversary that restricts learning to the most challenging tokens, directly linked to distributionally robust optimization through targeted reshaping of the training distribution. Building on this, we extend our approach to ADA-ESLM, an adaptive variant that dynamically calibrates the confidence level in response to the training dynamics, enabling a principled control over the compute-efficiency and generalization trade-off. Notably, ESLM requires no auxiliary supervision, reference models, or computationally expensive offline scoring. It also integrates naturally with knowledge distillation, allowing for risk-aware teacher supervision at the token level. Our key contributions are as follows:

- We propose ESLM, a risk-averse, two-phased selective language modeling that prioritizes high-risk inputs (i.e., informative or uncertain) for efficient LLM pretraining. ESLM first performs instance-level filtering using a lightweight proxy, then applies token-level loss shaping on selected instances. Selection in both phases uses VaR thresholding over loss or entropy statistics. We provide its two variants: **VaR-entropy** and **CVaR-loss** based on the risk score choice.
- We frame ESLM as a bilevel adversarial game between the model and a masker that perturbs the effective training distribution by selecting worst-case token subsets under a constrained thresholding rule. We further show that for loss-based selection, ESLM admits a distributionally robust optimization interpretation by recovering the CVaR objective (Rockafellar and Uryasev, 2002), a well-established risk-sensitive formulation from robust statistics (Ben-Tal et al., 2009).
- We propose ADA-ESLM, which adaptively adjusts the selection confidence level via a risk-aware controller guided by CVaR feedback to balance compute-efficiency and generalization.
- We demonstrate ESLM’s utility for knowledge distillation by enabling a sparsified risk-aware teacher supervision provided for selected high-risk tokens.
- Our experiments on GPT-2 (124M–1.5B) pretraining demonstrate that ESLM significantly reduces training FLOPs while maintaining or improving both validation perplexity and downstream task accuracy, with consistent gains across model sizes, dataset mixtures, and training settings.

2 BACKGROUND

This section introduces the key components underlying ESLM: CLM as the pretraining objective, token-level uncertainty estimation, and risk measures: VaR and CVaR.

Causal Language Modeling (CLM). CLM trains a language model (LM) θ to predict each token in a sequence given the previous context. Given a corpus \mathcal{C} of sequences $x = (x_1, \dots, x_M)$, each is of length T , drawn from a distribution \mathcal{D} over a vocabulary \mathcal{V} , the model factorizes the joint probability as: $P_{\theta}(x) = \prod_{j=1}^T P_{\theta}(x_j | x_{<j})$. CLM minimizes the average autoregressive loss:

$$\mathcal{L}_{\text{CLM}}(x; \theta) = \mathbb{E}_{x \sim \mathcal{D}}[\ell_{\theta}(x)] = \frac{1}{T} \sum_{j=1}^T -\log P_{\theta}(x_j | x_{<j}),$$

treating all tokens equally, despite many offering a limited learning signal (Lin et al., 2024).

Token-level risk. Let $S_{\theta}(x_j) \in \mathbb{R}$ denote the *risk score* of token x_j , under LM θ , computed via:

- (i) **Entropy** (Shannon, 1948): $H_{\theta}(x_j) = -\sum_{v \in \mathcal{V}} P_{\theta}(v | x_{<j}) \log P_{\theta}(v | x_{<j})$.
- (ii) **Loss**: $\ell_{\theta}(x_j) = -\log P_{\theta}(x_j | x_{<j})$.

Algorithm 1 ESLM

```

108 1: Input: LM  $\theta$ , dataset  $\mathcal{D}$ , learning rate  $\eta$ , confidence  $\alpha \in (0, 1)$ , batch size  $M$ , proxy  $L$  layers.
109 2: for each training iteration  $k = 1, \dots, K$  do
110 3:   Sample a batch of instances  $\mathcal{B} = \{x_1, \dots, x_M\} \sim \mathcal{D}$  of length  $T$   $\{x_j^t\}_{t=1}^T, \forall j \in \{1, \dots, M\}$ .
111 4:   Compute per-instance scores  $S_{\theta_k^L}(x)$  using early-exit  $\theta_k^L$  proxy. */ Entropy or loss
112 5:   Compute threshold  $S_{\theta_k^L, \alpha}^{\text{VaR}} \leftarrow \text{VaR}_\alpha \left( \{S_{\theta_k^L}(x_j)\}_{j=1}^M \right)$  using (1).
113 6:    $\bar{\mathcal{B}} \leftarrow \{x_j \in \mathcal{B} \mid S_{\theta_k^L}(x_j) \geq S_{\theta_k^L, \alpha}^{\text{VaR}}\}$ . */ High-risk instance selection
114 7:   Compute per-token statistics via:  $S_{\theta_k}(x_j^t) = \begin{cases} H_{\theta_k}(x_j^t) \text{ as in (i),} & \text{(VaR-entropy)} \\ \ell_{\theta_k}(x_j^t) \text{ as in (ii),} & \text{(CVaR-loss)} \end{cases}$ 
115 8:   Compute threshold  $S_{\theta_k, \alpha}^{\text{VaR}} \leftarrow \text{VaR}_\alpha \left( \{S_{\theta_k}(x_j^t)\}_{j \in \bar{\mathcal{B}}}\right)_{t=1}^T$  using (1).
116 9:    $\tilde{\mathcal{B}} \leftarrow \{x_j^t \in \bar{\mathcal{B}} \mid S_{\theta_k}(x_j^t) \geq S_{\theta_k, \alpha}^{\text{VaR}}\}$ . */ High-risk token selection
117 10:  Compute loss over selected tokens: */ Shaped loss
118 11:   $\mathcal{L}_{\tilde{\mathcal{B}}}(x; \theta_k) = \begin{cases} \mathbb{E}[\ell_{\theta_k}(x_j^t) \mid x_j^t \in \tilde{\mathcal{B}}], & \text{(VaR-entropy)} \\ \text{CVaR}_\alpha(\ell_{\theta_k}(x)) = \mathbb{E}[\ell_{\theta_k}(x_j^t) \mid x_j^t \in \tilde{\mathcal{B}}] \text{ using (2),} & \text{(CVaR-loss)} \end{cases}$ 
119 12:  Update model parameters using optimizer  $O$ :  $\theta_{k+1} \leftarrow O(\theta_k, \nabla_{\theta} \mathcal{L}_{\tilde{\mathcal{B}}}(x; \theta_k), \eta)$ .
120 13: end for
121 14: return  $\theta_K$ .

```

Both measures serve as proxies for token difficulty and informativeness—highlighting ambiguous, uncertain, or mispredicted tokens. **Instance-level risk scores** are obtained by aggregating token-level risk scores, i.e., the mean score over non-padded tokens within each sequence (line 4, Algorithm 1).

Risk measures. To prioritize high-impact tokens, we adopt risk-sensitive criteria from robust statistics (Gagne and Dayan, 2021). Let $S_{\theta}(x_j)$ denote per-token risk score computed via LM θ . The *value-at-risk* (VaR) (Rockafellar et al., 2000) at confidence level $\alpha \in (0, 1)$ is the minimum threshold such that only the top $(1 - \alpha)$ fraction of scores exceed the threshold; see Figure 1:

$$\text{VaR}_\alpha(S_\theta) := \inf\{\tau \in \mathbb{R} \mid P(S_\theta \geq \tau) \leq 1 - \alpha\}. \quad (1)$$

The corresponding CVaR is a coherent risk measure (Artzner et al., 1999) that computes the expected score among these highest-risk tokens (Rockafellar and Uryasev, 2002):

$$\text{CVaR}_\alpha(S_\theta) := \min_{\tau} \mathbb{E}_{x \sim \mathcal{D}} \left[\tau + \frac{1}{1 - \alpha} \max(0, S_\theta(x) - \tau) \right]. \quad (2)$$

These tail-risk measures allow us to reshape the training objective to emphasize tokens that are difficult or uncertain—an idea we exploit in the ESLM framework for efficient pretraining.

3 ESLM: RISK-AVERSE SELECTIVE LANGUAGE MODELING

We now introduce ESLM, a two-phase selective language modeling framework that improves pre-training efficiency by focusing optimization on high-risk inputs. We consider the standard causal language modeling setup presented in Section 2, where a language model with parameters θ is trained to minimize the expected token-level autoregressive loss. While effective, the expectation-based CLM objective assumes uniform importance across all tokens, leading to two key inefficiencies:

1. It wastes computation on trivially predictable tokens that dominate the loss landscape but offer little training signal.
2. It disregards token-level risk and overlooks rare, ambiguous, or out-of-distribution samples that are more informative for generalization and robustness.

To address these inefficiencies, we adopt the Selective Language Modeling (SLM) paradigm (Lin et al., 2024), which optimizes the model over a selected subset of tokens per training step. Formally, let $\pi_\phi(x)$ be a token selection policy that produces a binary mask $m = (m_1, \dots, m_T) \in \{0, 1\}^T$ for an input sequence x , the SLM objective becomes:

$$\mathcal{L}_{\text{SLM}}(\theta, \phi) = \mathbb{E}_{x \sim \mathcal{D}, m \sim \pi_\phi(x)} \left[\sum_{j=1}^T m_j \cdot \ell_\theta(x_j) \right],$$

where $\ell_\theta(x_j)$ is the per-token loss given in (ii). Existing approaches typically rely on learned or reference model (ϕ)-based policies for π_ϕ , that are expensive to train and may not generalize well (Lin et al., 2024), or design offline selectors (Xie et al., 2023b; Wettig et al., 2024). In contrast, we propose ESLM, a self-supervised online SLM framework rooted in statistical risk that eliminates the need for an auxiliary external selector to improve computational efficiency.

ESLM reshapes the loss towards *high-risk* inputs within each batch by dynamically filtering training signals both at the instance and token levels using risk-based thresholds derived from empirical batch distributions. Concretely, the risk is characterized by either (i) high predictive uncertainty (VaR-entropy selection) or (ii) high loss impact (CVaR-loss selection). At each step, a candidate batch $\mathcal{B} = \{x_1, \dots, x_M\} \sim \mathcal{D}$ is sampled. A shallow early-exit proxy model θ^L (first L layers of θ) is used to compute per-token risk scores (entropy (i) or loss (ii)), which are aggregated into per-instance scores: $S_{\theta^L}(x_j) = \frac{1}{T} \sum_{t=1}^T S_{\theta^L}(x_j^t)$. Given the empirical score distribution $\hat{\mathbb{P}}_{\mathcal{B}}$ over the batch, a VaR threshold at confidence level α is then applied to select high-risk instances:

$$\bar{\mathcal{B}} = \{x_j \in \mathcal{B} \mid S_{\theta^L}(x_j) \geq S_{\theta^L, \alpha}^{\text{VaR}}\}, \text{ where } S_{\theta^L, \alpha}^{\text{VaR}} = \inf \left\{ \tau \in \mathbb{R} \mid \hat{\mathbb{P}}_{\mathcal{B}}(S_{\theta^L}(x_j) \geq \tau) \leq 1 - \alpha \right\}.$$

On the reduced batch $\bar{\mathcal{B}}$, the model θ computes per-token risk scores $S_\theta(x_j^t)$ using either entropy (VaR-entropy) or loss (CVaR-loss). Given the empirical score distribution $\hat{\mathbb{P}}_{\bar{\mathcal{B}}}$ over the batch, a VaR threshold is applied at the token level, which defines a high-risk subset $\tilde{\mathcal{B}} = \{x_j^t \in \bar{\mathcal{B}} \mid S_\theta(x_j^t) \geq S_{\theta, \alpha}^{\text{VaR}}\}$, and an associated normalized training distribution: $Q_\tau \in \mathcal{P}_\alpha(\bar{\mathcal{B}}; \theta)$, $Q_\tau(x_j^t) \propto \mathbb{1}[x_j^t \in \tilde{\mathcal{B}}]$, where τ corresponds to the minimizer defined in (1). The training proceeds by minimizing:

$$\mathcal{L}_{\tilde{\mathcal{B}}}(\theta) = \mathbb{E}_{\tilde{\mathcal{B}} \sim \mathcal{D}} \left[\mathbb{E}_{x_j^t \sim Q_\tau} [\ell_\theta(x_j^t)] \right] = \mathbb{E}[\ell_\theta(x_j^t) \mid x_j^t \in \tilde{\mathcal{B}}],$$

which corresponds to CVaR (in (2)) when risk is based on token-level loss, and to an uncertainty-weighted loss when based on entropy. Our approach is detailed in Algorithm 1.

ESLM variations. While both VaR-entropy and CVaR-loss strategies select the upper tail of their respective score distributions; their inductive biases differ. The CVaR-loss selection emphasizes high-loss inputs, including both confidently incorrect predictions and uncertain correct ones. This helps the model correct overconfident mistakes and calibrate uncertainty. In contrast, the VaR-entropy selection focuses purely on predictive uncertainty, regardless of correctness, promoting learning in ambiguous or underexplored regions. We illustrate these differences through qualitative examples in Appendix F, showing that ESLM selects rare, or semantically rich tokens across domains.

The following formulations apply to token-level selection; however, instance-level selection is similar.

Bilevel game formulation. ESLM can be framed as a two-player adversarial game between the *model* and a *masker (adversary)*. This provides a bilevel optimization perspective where the masker perturbs the effective training distribution by choosing a threshold τ that determines which tokens are selected for training, under the VaR constraint, and the model minimizes its loss over the induced sub-distribution. Formally, the training process can be written as follows:

$$\min_{\theta} \mathbb{E}_{\tilde{\mathcal{B}} \sim \mathcal{D}} \left[\mathbb{E}_{x_j^t \sim Q_\tau} [\ell_\theta(x_j^t)] \right] \quad \text{subject to } \tau \in \arg \min_{\tilde{\tau} \in \mathbb{R}} \left\{ \tilde{\tau} \mid \hat{\mathbb{P}}_{\tilde{\mathcal{B}}} (S_\theta(x_j^t) \geq \tilde{\tau}) \leq 1 - \alpha \right\}, \quad (3)$$

where $\hat{\mathbb{P}}_{\tilde{\mathcal{B}}}$ is the empirical risk score distribution over the batch. This structure defines adversarial dynamics where the masker restricts the model to optimize over the most challenging subset of tokens, forcing it to improve performance on the tail distribution, and the model adapts to this shift. When the score function is $S_\theta(x_j^t) = \ell_\theta(x_j^t)$ (CVaR-loss), this procedure minimizes the CVaR at level α , thereby linking ESLM to classical risk-sensitive learning (Curi et al., 2020; Gagne and Dayan, 2021).

Distributionally robust optimization interpretation. ESLM admits a distributionally robust optimization (Duchi and Namkoong, 2021; Kuhn et al., 2025) interpretation. VaR thresholding restricts the training loss to a subset of tokens within the batch—those with scores in the top $(1 - \alpha)$ quantile. This induces an adversarial sub-distribution Q over the batch, supported only on the most challenging tokens. ESLM can then be seen as minimizing the worst-case expected loss over this ambiguity set:

$$\min_{\theta} \sup_{Q \in \mathcal{P}_\alpha(\bar{\mathcal{B}}; \theta)} \mathbb{E}_{x_j^t \sim Q} [\ell_\theta(x_j^t)]$$

$$\text{where } \mathcal{P}_\alpha(\bar{\mathcal{B}}; \theta) = \left\{ Q \ll \hat{\mathbb{P}}_{\bar{\mathcal{B}}} \mid \text{supp}(Q) \subseteq \{x_j^t \in \bar{\mathcal{B}} \mid S_\theta(x_j^t) \geq S_{\theta, \alpha}^{\text{VaR}}\} \right\}.$$

This robust optimization perspective explains why ESLM improves generalization: by optimizing performance under adversarial distributions, the model develops robustness to distributional shifts. The experimental results in Section 5.1 also confirm consistent generalization improvements by ESLM.

Risk-aware metric intuition. Unlike fixed or heuristic input selection thresholds, ESLM employs a quantile-based cutoff on the batch score distribution, selecting instances and tokens that fall within the high-risk tail. Aside from connecting ESLM to distributionally robust optimization, this design choice offers several advantages over the arbitrary thresholds: (i) the VaR/CVaR formulation provides a statistically principled mechanism to identify informative and difficult tokens, yielding both statistical guarantees on coverage and robustness to distribution shifts, (ii) the tail-risk formulation improves generalization (Section 5.1), a utility that extends beyond pure computational efficiency. Furthermore, by dynamically adjusting risk-awareness through the confidence level α , a token-level curriculum that balances computational efficiency and generalization can be derived, which we show with an adaptive extension of ESLM in Section 3.1.

Implementation and computational cost. We implement ESLM at the mini-batch level, compatible with distributed training. Each step has two passes: we first run an early-exit proxy forward with the first L transformer blocks of the same model in inference mode. From this shallow pass, we compute per-token statistics and aggregate to per-instance scores (mean over non-padded tokens). We then apply a VaR threshold over the M instance scores to select the subset. In the training pass, we run the full model on the selected instances and apply token-level loss shaping as given in Algorithm 1. The computational overhead from VaR thresholding is minimal, requiring $O(M \log M)$ time per batch (with batch size M). This cost is negligible compared to the dominant forward/backward FLOPs of the main model. We discuss the runtime overhead in Appendix D.3.

FLOPs accounting. Following Kaplan et al. (2020); Chowdhery et al. (2023), a full forward-backward pass costs $6N + 12LHQT$ FLOPs per token, for an LM with N parameters, L layers, H attention heads, head size Q , and sequence length T . In the proxy phase, we run a forward-only early-exit of depth $L' \leq L$, costing $2N_{\text{proxy}} + 6L'HQT$ FLOPs per token, where N_{proxy} is the parameters for embeddings, layer norms and the first L' transformer blocks. Since the risk scores need logits, we add $2DV$ FLOPs per token for the head projection, with embedding dimension D and vocabulary size V . During token-level loss shaping, masked tokens skip part of the backward in the final layer components (head + FFN matmuls + final/pre-FFN layer norms), yielding $\approx 4DV + 4D + 32D^2 + 4D$ FLOPs saving *per masked token* (Kaplan et al., 2020). ESLM’s optimized FLOPs savings come primarily from (i) dropping whole instances after the proxy pass, which avoids their full forward-backward in the training phase; and (ii) skipping last-block gradient computation for loss-masked tokens.

Downstream impact. Effective pretraining increasingly hinges on how data is selected (Tirumala et al., 2023; Mayilvahanan et al., 2025). Improvements in training loss does not guarantee better downstream generalization, particularly under distribution shift (Ramanujan et al., 2023; Isik et al., 2025). ESLM addresses this by providing hierarchical control over which parts of the input receive focus, concentrating optimization to high-risk tokens. In Section 5.1, we demonstrate that ESLM improves both loss-vs-compute efficiency and downstream performance than standard training.

Token vs instance-level selection. Unlike methods that filter or reweight entire sequences (Wang et al., 2024; Sow et al., 2025), ESLM uses a *hierarchical* scheme: a cheap proxy pass prunes low-value instances, then loss is shaped at the finer granularity of tokens. This combines compute-awareness with information-awareness, retaining useful tokens within kept sequences. ESLM is natively compatible with autoregressive training and avoids data pipeline changes. Under same compute budget, this hierarchical selection generalizes better than coarse instance-only selection methods (Section 5.1).

3.1 ADA-ESLM: ADAPTIVE CONFIDENCE THRESHOLDING

While a fixed confidence level α in ESLM yields strong efficiency gains (see Section 5.1), the optimal α level may vary throughout training. Early in training, broad token coverage may improve generalization, whereas later stages benefit from focusing on harder or more informative tokens. To accommodate this, we introduce ADA-ESLM, a dynamic variant that adjusts α during training using a *risk-sensitive controller* driven by CVaR feedback. In each evaluation step k , we compute CVaR_{α_k} of the per-token risk scores, using (2). We then track the changes in CVaR to detect shifts in training difficulty, estimated from model training dynamics, and update α using a multiplicative rule:

$$\alpha_{k+1} \leftarrow \alpha_k \cdot \exp(-\gamma \cdot \Delta_{\text{norm}}(\alpha_k)), \text{ where } \Delta_{\text{norm}}(\alpha_k) := \frac{\text{CVaR}_{\alpha_k} - \text{CVaR}_{\alpha_{k-1}}}{\text{CVaR}_{\alpha_{k-1}} + \varepsilon}.$$

Here, $\Delta_{\text{norm}}(\alpha_k)$ is a dimension and scale-independent signal capturing the relative change in CVaR, $\gamma > 0$ controls adaptation rate, and ε is a small constant for numerical stability. The core idea for this update rule is *stabilizing* CVaR: if $\Delta_{\text{norm}} > 0$ (i.e., CVaR increases), the model is encountering harder tokens, α is then decreased to include more tokens and expand the training signal. Conversely, if $\Delta_{\text{norm}} < 0$, the model is improving on difficult tokens. We increase α to focus learning on high-risk tokens. ADA-ESLM extends the adversarial game in (3) by equipping the masker with a CVaR-driven controller that adapts token sparsity in response to training dynamics, offering a form of curriculum learning. We provide the ADA-ESLM algorithm in Appendix B (see Algorithm 2).

3.2 ESLM-KD: RISK-AWARE KNOWLEDGE DISTILLATION WITH ESLM

Knowledge distillation transfers knowledge from a teacher model to a student by encouraging the student to match the teacher’s output distribution (Buciluă et al., 2006; Hinton, 2015). In language modeling, Rawat et al. (2024) showed that a small LM supervision improves the training of a much more capable LLM. While the standard framework operates over all tokens—typically using sequence- or word-level KL divergence (Kim and Rush, 2016)—we can utilize ESLM for risk-aware distillation.

To this end, we provide a use case, ESLM-KD, with the implementation details presented in Algorithm 3 in Appendix C. Specifically, we apply VaR_α thresholding using a student LM to select high-risk tokens, which are then used to compute the KL divergence between teacher and student logits. The student is trained only on these selected tokens, focusing its capacity on uncertain or error-prone regions. This strategy is teacher-agnostic, relying on the internal statistics of the student model for selection. ESLM-KD generates a sparse supervision signal based on selected tokens, resulting in improved compute and sample efficiency, as further empirically supported in Section 5.1.1.

4 RELATED WORK

Online data subset selection. Efficient data selection is essential for scaling LLM pretraining, where full-corpus training is often prohibitively expensive (Albalak et al., 2024). While early work focused on static or offline methods, such as filtering (Marion et al., 2023) or scoring examples before training (Coleman et al., 2020; Xie et al., 2023b; Wettig et al., 2024) or during fine-tuning (Xia et al., 2024), such methods lack adaptability and struggle to scale in large-batch or continual pretraining settings. Online data selection overcomes these limitations by adapting to the evolving state of the model. Early strategies on online example-level selection prioritized high-loss samples to accelerate convergence (Loshchilov and Hutter, 2015; Katharopoulos and Fleuret, 2018; Jiang et al., 2019) or leveraged gradients (Killamsetty et al., 2021). Recent works (Mindermann et al., 2022; Wang et al., 2024) apply gradient-based influence scoring (Sachdeva et al., 2024) to guide instance selection or leverage reference models for token selection (Fan and Jaggi, 2023; Lin et al., 2024); however, they often incur high memory due to expensive gradient computations or additional supervision costs from curated reference models and validation sets. In contrast, ESLM introduces a *lightweight*, self-supervised, hierarchical batch selection mechanism, eliminating offline preprocessing, external supervision, or costly gradient tracing. This yields an easily integrable approach that achieves a favorable trade-off between compute efficiency and robustness, while remaining agnostic to training configurations.

Risk-aversion in language modeling. Risk-sensitive optimization offers a principled mechanism to enhance robustness by focusing training on high-risk examples (Rockafellar et al., 2000). The CVaR objective has been previously studied in classification (Curi et al., 2020), submodular optimization (Maehara, 2015), and fair learning (Williamson and Menon, 2019), typically to mitigate the influence of tail-risk or worst-case samples. However, in the context of language modeling, CVaR-based approaches remain relatively underexplored. Notable exceptions include methods (Oren et al., 2019) that aggregate losses over topics to address distributional shift but these typically operate at the group level, or in fine-tuning LLMs with reinforcement learning (Chaudhary et al., 2024). On the contrary, ESLM brings risk-aware optimization to the token level for LLM pretraining. Each batch is shaped into a high-risk sub-distribution by the fine-grained risk control of ESLM, incorporating a distributionally robust view of token-level optimization. Unlike heuristic loss-based filtering, ESLM offers a theoretically grounded and practical approach for efficient and robust large-scale pretraining under uncertainty.

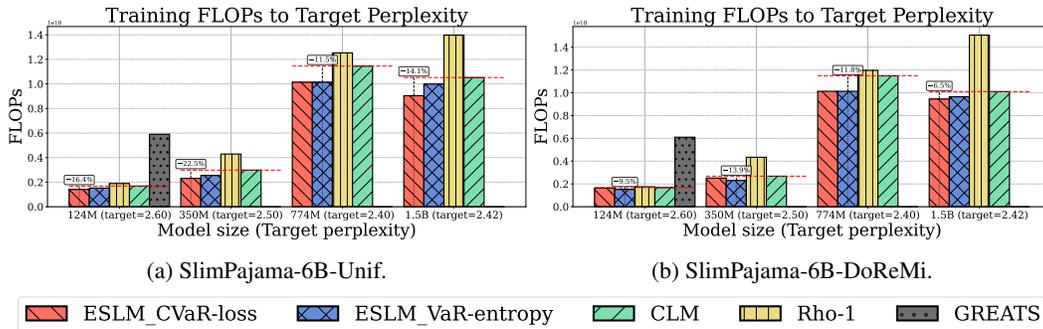


Figure 2: **Training FLOPs (\downarrow) required for convergence to target validation (log) perplexity.** We report training FLOPs required by the {124M, 350M, 774M, 1.5B} parameter models to converge to a target validation loss threshold across datasets. The % labels show the percentage FLOPs savings provided by ESLM relative to the best performing baseline. ESLM reduces training cost by focusing optimization on the high-risk instances and tokens and eliminating redundant gradient computation.

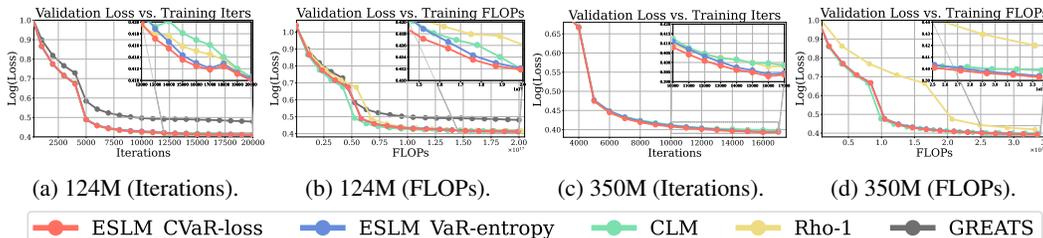


Figure 3: **Validation loss convergence.** We report convergence of validation loss versus training FLOPs/iterations of models trained on SlimPajama-6B-Unif mixture. ESLM variants provide faster convergence to lower loss values in terms of iterations, while requiring fewer FLOPs than baselines at the convergence. See Appendix E.1 for additional results on other pretraining corpora.

5 EXPERIMENTS

In this section, we evaluate ESLM with two variants (ESLM-CVaR-loss and ESLM-VaR-entropy) across diverse pretraining settings—varying model scales, data mixtures, and training budgets—to assess its impact on both efficiency and generalization.

In our experiments, we use the SlimPajama-6B (Soboleva et al., 2023) dataset: a 6B token mixture spanning seven domains {Arxiv, Book, CommonCrawl, C4, GitHub, StackExchange, Wikipedia}, used with both uniform and DoReMi (Xie et al., 2023a) domain weights (see Appendix D.2 for the exact weight values).

Experimental setup. We pretrain GPT-2 models with 124M, 350M, 774M, and 1.5B parameters using a BPE tokenizer (Sennrich et al., 2016) with vocabulary size 50,304. All models are trained with a sequence length of 1024, gradient accumulation over 40 steps, and mini-batch sizes of {4, 8, 12}. We use AdamW with cosine learning rate decay; full hyperparameters are provided in Appendix D.1. We apply ESLM with confidence level α tuned over the set {0.05, 0.1, 0.2}. Additional results with varying α levels are presented in Section 5.2.

Baselines. We compare ESLM variants against regular training and online batch selection methods: (1) CLM, introduced in Section 2, (2) Rho-1 (Lin et al., 2024), an online SLM using a reference model to score token loss differentials, (3) GREATS (Wang et al., 2024), a state-of-the-art online sample selection method based on high-quality validation data and per-sample gradients. GREATS’ high memory requirements, even with ghost inner product optimizations, limited our comparisons to the 124M setting. For distillation experiments (Section 5.1.1), we further compare against the dense distillation and SALT (Rawat et al., 2024) methods. We provide the baseline details in Appendix D.4.

Performance metrics. We assess our method concerning training efficiency and generalization ability by tracking the metrics: (i) training FLOPs required to converge to target validation perplexity, (ii) validation loss convergence versus training FLOPs/iterations, and (iii) zero-/few-shot accuracy

Table 1: **Generalization performance on downstream tasks.** All models (124M) are pretrained under a $\sim 3E17$ FLOPs budget on SlimPajama-6B-Unif mixture. We report the best observed accuracy (standard error) or exact match if provided, during training. **Highlighted** values indicate the best performance. See Appendix E.2 for the results under various model sizes and datasets.

Benchmark	# Shots	Method (124M)				
		ESLM-CVaR-loss	ESLM-VaR-entropy	CLM	Rho-1	GREATS
ARC-E (Clark et al., 2018)	0-shot	0.3712 _(0.0099)	0.3766 _(0.0099)	0.3644 _(0.0099)	0.3657 _(0.0099)	0.3236 _(0.0096)
LAMBADA (Paperno et al., 2016)	5-shot	0.1595 _(0.0051)	0.1721 _(0.0053)	0.1701 _(0.005)	0.1680 _(0.005)	0.0254 _(0.002)
SciQ (Welbl et al., 2017)	5-shot	0.714 _(0.0143)	0.71 _(0.0144)	0.6970 _(0.0145)	0.7000 _(0.0145)	0.4350 _(0.0157)
HellaSwag (Zellers et al., 2019)	5-shot	0.2930 _(0.0045)	0.2964 _(0.0046)	0.2901 _(0.0045)	0.2893 _(0.0045)	0.2621 _(0.0044)
TriviaQA (Joshi et al., 2017)	1-shot	0.0128 _(0.0001)	0.0066 _(0.0006)	0.0078 _(0.0007)	0.0090 _(0.0007)	0.0007 _(0.0002)
COPA (Wang et al., 2019)	5-shot	0.64 _(0.0482)	0.62 _(0.0488)	0.62 _(0.0488)	0.62 _(0.0488)	0.64 _(0.0482)
MultiRC (Wang et al., 2019)	5-shot	0.5633 _(0.0071)	0.5429 _(0.0072)	0.5338 _(0.0072)	0.5338 _(0.0072)	0.5497 _(0.0071)
OpenBookQA (Mihaylov et al., 2018)	5-shot	0.164 _(0.0166)	0.17 _(0.0168)	0.174 _(0.017)	0.164 _(0.0166)	0.148 _(0.0159)
PiQA (Bisk et al., 2020)	5-shot	0.6240 _(0.0113)	0.6191 _(0.0113)	0.6099 _(0.0114)	0.6180 _(0.0113)	0.5571 _(0.0116)
Average (\uparrow)		0.39353	0.39041	0.38523	0.38531	0.32684

(normalized, if provided) in downstream benchmark tasks from the lm-eval-harness (Gao et al., 2024) suite, spanning QA, reasoning, and generation. To track convergence behavior, we report the learning trajectory using a running average with a window size of 5 evaluation points. We further evaluate performance across model sizes and dataset mixtures. We estimate training FLOPs as explained in Section 3. The details on metrics and experimental setup are provided in Appendix D.

5.1 EXPERIMENTAL RESULTS

In this section, we report the performance of ESLM variants against the baseline methods, followed by the results of its implementation as a knowledge distillation mechanism and ablation analyses.

Validation loss vs training FLOPs/iterations. As presented in Figure 2, ESLM consistently requires fewer training FLOPs to reach target validation loss across model sizes and datasets. ESLM provides strong efficiency gains at convergence, reaching the level of, e.g., -22.5% reduced FLOPs relative to standard training in a 350M setting. Notably, these gains over CLM are essentially free, i.e., ESLM operates under the same optimizer, dataset, with no external supervision. Figure 3 further shows that ESLM accelerates learning convergence compared to the baselines, discovering lower validation loss with fewer iterations. While the proxy phase adds a small upfront cost, subsequent instance pruning and token-level loss shaping drive lower losses and net FLOPs savings at convergence across scales compared to CLM. Unlike Rho-1, which queries an external reference model—adding extra compute overhead and requiring high-quality data—ESLM instead leverages model-internal training dynamics for selection, avoiding such offline preprocessing. Similarly, while GREATS employs an efficient ghost inner product approximation (Wang et al., 2024), it still relies on curated validation data and per-sample gradient estimates, which are impractical at larger model scales due to high memory demands. In contrast, ESLM operates without gradient tracing and scales naturally. Moreover, GREATS performs selection at the instance level, often discarding informative tokens within partially useful sequences—leading to worse perplexity. This highlights the instance and token-level granularity of ESLM, which avoids this limitation by emphasizing valuable sub-sequence information.

Downstream performance. Table 1 summarizes the best zero-/few-shot accuracy achieved by 124M models trained on SlimPajama-6B-Unif under a fixed compute budget of $\sim 3E17$ FLOPs. ESLM variants significantly outperform baselines in average accuracy, with consistent gains over GREATS and Rho-1 across all tasks. Figure 4 further illustrates the accuracy norm convergence on the HellaSwag benchmark, where ESLM-CVaR-loss achieves faster early gains, while ESLM-VaR-entropy surpasses baselines in later training stages. These results, including additional evaluations in Appendix E.2, show that ESLM improves both training efficiency and generalization.

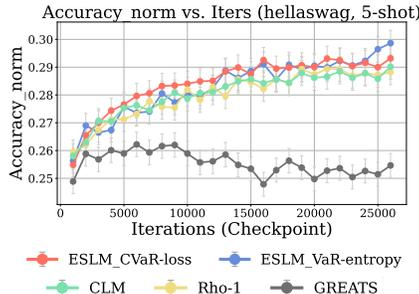


Figure 4: 5-shot accuracy (norm) (\uparrow) performance on HellaSwag throughout training. ESLM variants discover higher accuracy levels than baselines, with particular gains in the later training stages.

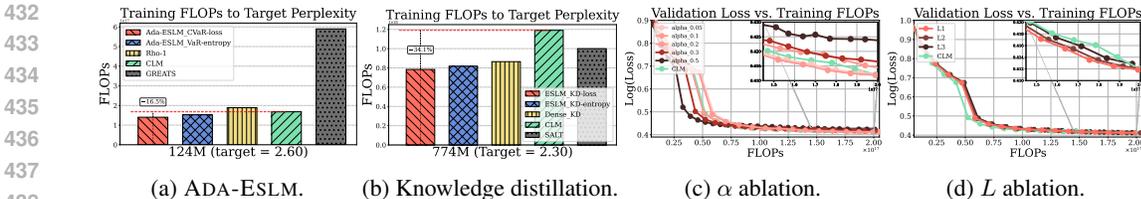


Figure 5: **Extended analyses for ESLM.** (a): ADA-ESLM reduces FLOPs for convergence to target validation loss (\downarrow) by adaptively tuning α based on training dynamics. (b): In knowledge distillation for 774M, ESLM-KD converges to target validation loss with substantially fewer FLOPs. (c): Varying the α level enables flexible control over the trade-off between training efficiency and model quality. (d): Proxy depth ablation shows that shallow proxies achieve similar convergence performance to deeper alternatives; hence, minimal proxy computation suffices for effective instance selection.

ADA-ESLM experiments. We train ADA-ESLM (124M) with $\alpha_0 = 0.1, \gamma = 0.5$ on SlimPajama-6B-Unif. Figure 5a reveals that ADA-ESLM achieves the target validation perplexity with -16.5% training FLOPs savings at convergence. As detailed in Appendix B (Figure 7), ADA-ESLM provides an implicit curriculum learning: the training process begins with broader input coverage and gradually shifts focus toward higher-risk tokens—without manual scheduling or external supervision. Adaptively adjusting α based on CVaR feedback stabilizes training while offering a principled trade-off between compute-efficiency and generalization. Downstream evaluations in Figure 8 and Table 2 (Appendix B) confirm that ADA-ESLM further achieves higher average downstream accuracy than baselines, improving generalization while maintaining high training efficiency.

5.1.1 EXPERIMENTS FOR KNOWLEDGE DISTILLATION WITH ESLM-KD

To utilize ESLM for risk-aware knowledge distillation (Section 3.2), we pretrain a 774M student LM using a 124M teacher on SlimPajama-6B-Unif. We set the distillation weight $\lambda = 0.5$ and teacher temperature $\rho = 1.0$. Training details are provided in Appendix C. We compare ESLM-KD against three baselines: CLM, dense distillation without token selection, and SALT (Rawat et al., 2024), a two-stage distillation-then-pretraining pipeline. As shown in Figure 5b, ESLM-KD models converge to the target validation perplexity with substantially fewer FLOPs; with -34% FLOPs savings compared to CLM. Furthermore, as reported in Appendix C (Table 3), it outperforms baseline models in downstream tasks, demonstrating the effectiveness of ESLM for efficient and generalizable distillation.

5.2 ABLATION AND ADDITIONAL ANALYSES

- **Confidence level (α):** In Figure 5c, we assess the sensitivity of ESLM to varying α values: $\{0.05, 0.1, 0.2, 0.3, 0.5\}$. Lower α values improve data coverage but increase computation, whereas higher α levels enhance efficiency at the cost of underutilization. We find $\alpha \in [0.1, 0.2]$ offers a favorable trade-off between compute savings and generalization.
- **Proxy depth (L):** In Figure 5d, we assess the effectiveness of ESLM under varying L levels: $\{1, 2, 3\}$ for 124M model. ESLM achieves similar validation loss convergence, hence, proxy computation with $L = 1$ suffices for effective instance selection while minimizing overhead costs.
- **Model size:** Across 124M, 350M, 774M, and 1.5B GPT-2 models, ESLM consistently improves efficiency and generalization (see Figures 2, 3, 4, and Appendix E).
- **Pretraining corpus:** We evaluate ESLM on SlimPajama-6B corpus with different domain mixture weights (uniform and DoReMi). The method generalizes well across corpora (see Figure 2) without requiring domain-specific tuning. Detailed results are reported in Appendix E.
- **Token selection analysis:** To better understand the behavior of ESLM variants, we analyze the selected tokens across different domains in Appendix F, which reveals that ESLM focuses on rare, or contextually ambiguous tokens—supporting its risk-aware design.
- **Instance vs token-level selection:** In Appendix-E.3, we report an ablation decomposing ESLM’s FLOPs savings into contributions from proxy instance selection and token-level loss shaping.
- **Robustness to label noise:** To show the broader utility of selective pretraining with high risk inputs, we further examine the robustness properties of ESLM grounded in the DRO framework using label noise injection into the pretraining setup. In practice, real-world pretraining corpora inherently contain noisy tokens from data quality heterogeneity, a challenge in large-scale web data pretraining. For this, we pretrain the models on SlimPajama-6B-Unif with 5% label noise, randomly permuting a portion of the next-token targets. The results in Figure 6 show that even

under noisy supervision, ESLM variants achieve the target validation loss with fewer FLOPs compared to standard training, demonstrating that its risk-aware token selection mitigates the effect of corrupted gradients. This experiment further supports the utility of ESLM’s hierarchical risk-aware filtering, not only in improving compute efficiency but also in enhancing robustness to label noise.

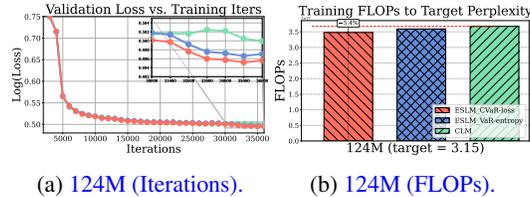


Figure 6: **Pretraining under 5% label noise.** (a): As shown with the learning trajectory using a running average with a window size of five evaluation points, ESLM variants converge faster to lower validation loss values under label noise compared to standard training. (b): Even under noisy supervision, ESLM achieves the target validation loss with 5.4% FLOPs savings compared to CLM.

6 LIMITATIONS

While ESLM offers an effective approach to selective pretraining, it inherently trades off completeness for efficiency by sparsifying backpropagation. Although this provides compute savings (FLOPs), in practice, it may underutilize the full training signal. As we discuss in Appendix D.3, integrating ESLM with sparsity-aware accelerators could further enhance resource utilization. Consistent with recent efficiency works (Lasby et al., 2024), achieving practical acceleration requires optimized hardware support for sparse backpropagation—an engineering direction beyond the paper’s scope.

7 CONCLUSION

We introduce ESLM, a selective language modeling with risk-aware hierarchical batch selection, which improves compute efficiency and generalization of LLM pretraining. Rather than training uniformly over all tokens, ESLM applies a risk-sensitive VaR threshold to prioritize high-utility tokens and skip redundant ones during backpropagation. This data-centric strategy effectively improves loss-per-FLOP efficiency, without modifying the model, optimizer, or dataset. By focusing optimization on the most informative inputs, ESLM improves generalization and enhances scalability in language modeling. As a future work, ESLM—along with its ADA-ESLM variant and integration with knowledge distillation—opens new directions in risk-aware token-level curriculum learning, adaptive compute allocation, and risk-aware data valuation for sustainable and efficient LLM scaling.

ETHICS & SOCIETAL IMPACT STATEMENT

Our work introduces ESLM, a selective language modeling framework which improves training efficiency and generalization performance in LLM pretraining via risk-aware instance and token selection. On the one hand, ESLM enables compute-efficient training by focusing optimization on the most informative parts of the input. This could reduce the energy footprint of large-scale training runs, make LLM development more accessible to institutions with limited compute budgets, and improve model robustness—particularly in out-of-distribution scenarios.

On the other hand, improved training efficiency may accelerate the development of powerful generative models, some of which could be misused for disinformation, synthetic media, or other harmful applications. In addition, token-level filtering methods—if miscalibrated—may reinforce spurious patterns or underrepresent minority language phenomena, inadvertently encoding or amplifying societal biases in the training data. Although ESLM is not tied to a specific application, its performance gains could boost the downstream impact of any application built upon the pretrained models. As our method provides purely training-time improvement, it does not increase model capacity or inference capability directly, which partially limits its risk surface.

Finally, our study does not involve intervention with or collection of data from human participants. Our experiments use widely adopted language-model pretraining corpora or mixtures assembled from publicly available sources under their respective licenses, with provided references.

REPRODUCIBILITY STATEMENT

For reproducibility, we provide the algorithmic description of ESLM in Algorithm 1 in Section 3. We further explain the implementation details in a separate section for reproducibility in Appendix D.6. In the same section, we provide a reference to the open source code repository on which our implementation is based, the evaluation benchmark suite used for evaluation, and specific implementation mechanisms. Furthermore, in Section 5, we detail the specific experimental setup along with references, specifically the type of model, pretraining corpus and specific dataset mixture weights (see also Appendix D.1-D.2) we use in our experiments.

REFERENCES

- 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. (2024). A survey on data selection for language models. *ArXiv*, 2402.16827.
- Artzner, P., Delbaen, F., Eber, J.-M., and Heath, D. (1999). Coherent measures of risk. *Mathematical Finance*, 9(3):203–228.
- Ben-Tal, A., Ghaoui, L., and Nemirovski, A. (2009). *Robust Optimization*. Princeton Series in Applied Mathematics. Princeton University Press.
- Bisk, Y., Zellers, R., Gao, J., Choi, Y., et al. (2020). PIQA: Reasoning about physical commonsense in natural language. *AAAI Conference on Artificial Intelligence*, pages 7432–7439.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. (2020). Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901.
- Buciluă, C., Caruana, R., and Niculescu-Mizil, A. (2006). Model compression. *International Conference on Knowledge Discovery and Data Mining*, pages 535–541.
- Chaudhary, S., Dinesha, U., Kalathil, D., and Shakkottai, S. (2024). Risk-averse fine-tuning of large language models. *Advances in Neural Information Processing Systems*, 37:107003–107038.
- Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., Barham, P., Chung, H. W., Sutton, C., Gehrmann, S., et al. (2023). PaLM: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113.

- 594 Clark, P., Cowhey, I., Etzioni, O., Khot, T., Sabharwal, A., Schoenick, C., and Tafjord, O. (2018).
595 Think you have Solved Question Answering? Try ARC, the AI2 Reasoning Challenge. *ArXiv*,
596 1803.05457.
- 597 Coleman, C., Yeh, C., Musmann, S., Mirzasoleiman, B., Bailis, P., Liang, P., Leskovec, J., and
598 Zaharia, M. (2020). Selection via proxy: Efficient data selection for deep learning. *International*
599 *Conference on Learning Representations*.
- 600 Curi, S., Levy, K. Y., Jegelka, S., and Krause, A. (2020). Adaptive sampling for stochastic risk-averse
601 learning. *Advances in Neural Information Processing Systems*, 33:1036–1047.
- 602 Dao, T. (2023). FlashAttention-2: Faster attention with better parallelism and work partitioning.
603 *ArXiv*, 2307.08691.
- 604 Duchi, J. and Namkoong, H. (2021). Learning models with uniform performance via distributionally
605 robust optimization. *The Annals of Statistics*, 49(3):1378–1406.
- 606 Fan, S. and Jaggi, M. (2023). Irreducible curriculum for language model pretraining. *ArXiv*,
607 2310.15389.
- 608 Fan, S., Pagliardini, M., and Jaggi, M. (2023). DoGE: Domain reweighting with generalization
609 estimation. *ArXiv*, 2310.15393.
- 610 Gagne, C. and Dayan, P. (2021). Two steps to risk sensitivity. *Advances in Neural Information*
611 *Processing Systems*, 34:22209–22220.
- 612 Gao, L., Tow, J., Abbasi, B., Biderman, S., Black, S., DiPofi, A., Foster, C., Golding, L., Hsu,
613 J., Le Noac’h, A., Li, H., McDonell, K., Muennighoff, N., Ociepa, C., Phang, J., Reynolds,
614 L., Schoelkopf, H., Skowron, A., Sutawika, L., Tang, E., Thite, A., Wang, B., Wang, K., and
615 Zou, A. (2024). The language model evaluation harness. [https://zenodo.org/records/
616 12608602](https://zenodo.org/records/12608602).
- 617 Hinton, G. (2015). Distilling the knowledge in a neural network. *ArXiv*, 1503.02531.
- 618 Hou, L., Pang, R. Y., Zhou, T., Wu, Y., Song, X., Song, X., and Zhou, D. (2022). Token dropping
619 for efficient BERT pretraining. *Annual Meeting of the Association for Computational Linguistics*,
620 pages 3774–3784.
- 621 Hüllermeier, E. and Waegeman, W. (2021). Aleatoric and epistemic uncertainty in machine learning:
622 An introduction to concepts and methods. *Machine learning*, 110(3):457–506.
- 623 Isik, B., Ponomareva, N., Hazimeh, H., Pappas, D., Vassilvitskii, S., and Koyejo, S. (2025). Scaling
624 laws for downstream task performance in machine translation. *International Conference on*
625 *Learning Representations*.
- 626 Jiang, A. H., Wong, D. L.-K., Zhou, G., Andersen, D. G., Dean, J., Ganger, G. R., Joshi, G., Kaminsky,
627 M., Kozuch, M., Lipton, Z. C., et al. (2019). Accelerating deep learning by focusing on the biggest
628 losers. *ArXiv*, 1910.00762.
- 629 Joshi, M., Choi, E., Weld, D., and Zettlemoyer, L. (2017). TriviaQA: A Large Scale Distantly
630 Supervised Challenge Dataset for Reading Comprehension. *ArXiv*, 1705.03551.
- 631 Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., Gray, S., Radford, A.,
632 Wu, J., and Amodei, D. (2020). Scaling laws for neural language models. *ArXiv*, 2001.08361.
- 633 Karpathy, A. (2022). NanoGPT [GitHub repository].
- 634 Katharopoulos, A. and Fleuret, F. (2018). Not all samples are created equal: Deep learning with
635 importance sampling. *International Conference on Machine Learning*, pages 2525–2534.
- 636 Killamsetty, K., Durga, S., Ramakrishnan, G., De, A., and Iyer, R. (2021). Grad-match: Gradient
637 matching based data subset selection for efficient deep model training. *International Conference*
638 *on Machine Learning*, pages 5464–5474.

- 648 Kim, Y. and Rush, A. M. (2016). Sequence-level knowledge distillation. *Conference on Empirical*
649 *Methods in Natural Language Processing*, pages 1317–1327.
- 650
- 651 Kuhn, D., Shafiee, S., and Wiesemann, W. (2025). Distributionally robust optimization. *ArXiv*,
652 2411.02549.
- 653 Lasby, M., Golubeva, A., Evci, U., Nica, M., and Ioannou, Y. (2024). Dynamic sparse training with
654 structured sparsity. *International Conference on Learning Representations*.
- 655
- 656 Lin, Z., Gou, Z., Gong, Y., Liu, X., Shen, Y., Xu, R., Lin, C., Yang, Y., Jiao, J., Duan, N., et al. (2024).
657 Rho-1: Not all tokens are what you need. *ArXiv*, 2404.07965.
- 658 Loshchilov, I. and Hutter, F. (2015). Online batch selection for faster training of neural networks.
659 *ArXiv*, 1511.06343.
- 660
- 661 Maehara, T. (2015). Risk averse submodular utility maximization. *Operations Research Letters*,
662 43(5):526–529.
- 663 Marion, M., Üstün, A., Pozzobon, L., Wang, A., Fadaee, M., and Hooker, S. (2023). When less is
664 more: Investigating data pruning for pretraining LLMs at scale. *ArXiv*, 2309.04564.
- 665
- 666 Mayilvahanan, P., Wiedemer, T., Mallick, S., Bethge, M., and Brendel, W. (2025). LLMs on the line:
667 Data determines loss-to-loss scaling laws. *ArXiv*, 2502.12120.
- 668 Mihaylov, T., Clark, P., Khot, T., and Sabharwal, A. (2018). Can a suit of armor conduct electricity?
669 a new dataset for open book question answering. *ArXiv*, 1809.02789.
- 670
- 671 Mindermann, S., Brauner, J. M., Razzak, M. T., Sharma, M., Kirsch, A., Xu, W., Hölting, B., Gomez,
672 A. N., Morisot, A., Farquhar, S., et al. (2022). Prioritized training on points that are learnable, worth
673 learning, and not yet learnt. *International Conference on Machine Learning*, pages 15630–15649.
- 674 Oren, Y., Sagawa, S., Hashimoto, T. B., and Liang, P. (2019). Distributionally robust language
675 modeling. *ArXiv*, 1909.02060.
- 676
- 677 Paperno, D., Kruszewski, G., Lazaridou, A., Pham, N. Q., Bernardi, R., Pezzelle, S., Baroni, M.,
678 Boleda, G., and Fernandez, R. (2016). The LAMBADA dataset: Word prediction requiring a
679 broad discourse context. *Annual Meeting of the Association for Computational Linguistics*, pages
680 1525–1534.
- 681 Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein,
682 N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy,
683 S., Steiner, B., Fang, L., Bai, J., and Chintala, S. (2019). PyTorch: An Imperative Style, High-
684 Performance Deep Learning Library. *Advances in Neural Information Processing Systems*, 32:8024–
685 8035.
- 686 Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., and Sutskever, I. (2019). Language models are
687 unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- 688
- 689 Ramanujan, V., Nguyen, T., Oh, S., Farhadi, A., and Schmidt, L. (2023). On the connection between
690 pre-training data diversity and fine-tuning robustness. *Advances in Neural Information Processing*
691 *Systems*, 36:66426–66437.
- 692 Rawat, A. S., Sadhanala, V., Rostamizadeh, A., Chakrabarti, A., Jitkrittum, W., Feinberg, V., Kim,
693 S., Harutyunyan, H., Saunshi, N., Nado, Z., et al. (2024). A little help goes a long way: Efficient
694 LLM training by leveraging small LMs. *ArXiv*, 2410.18779.
- 695
- 696 Rockafellar, R. T. and Uryasev, S. (2002). Conditional value-at-risk for general loss distributions.
697 *Journal of Banking & Finance*, 26(7):1443–1471.
- 698 Rockafellar, R. T., Uryasev, S., et al. (2000). Optimization of conditional value-at-risk. *Journal of*
699 *Risk*, 2:21–42.
- 700
- 701 Sachdeva, N., Coleman, B., Kang, W.-C., Ni, J., Hong, L., Chi, E. H., Caverlee, J., McAuley, J., and
Cheng, D. Z. (2024). How to train data-efficient LLMs. *ArXiv*, 2402.09668.

- 702 Sennrich, R., Haddow, B., and Birch, A. (2016). Neural machine translation of rare words with
703 subword units. *Annual Meeting of the Association for Computational Linguistics*, pages 1715–
704 1725.
- 705 Shannon, C. E. (1948). A mathematical theory of communication. *The Bell system technical journal*,
706 27(3):379–423.
- 707
- 708 Soboleva, D., Al-Khateeb, F., Myers, R., Steeves, J. R., Hestness, J., and Dey, N. (2023). SlimPajama:
709 A 627B token cleaned and deduplicated version of RedPajama.
- 710
- 711 Sow, D., Woiseschläger, H., Bulusu, S., Wang, S., Jacobsen, H.-A., and Liang, Y. (2025). Dynamic
712 loss-based sample reweighting for improved large language model pretraining. *ArXiv*, 2502.06733.
- 713 Tirumala, K., Simig, D., Aghajanyan, A., and Morcos, A. (2023). D4: Improving LLM pretraining
714 via document de-duplication and diversification. *Advances in Neural Information Processing*
715 *Systems*, 36:53983–53995.
- 716
- 717 Wang, A., Pruksachatkun, Y., Nangia, N., Singh, A., Michael, J., Hill, F., Levy, O., and Bowman, S.
718 (2019). SuperGLUE: A stickier benchmark for general-purpose language understanding systems.
719 *Advances in Neural Information Processing Systems*, 32.
- 720 Wang, J. T., Wu, T., Song, D., Mittal, P., and Jia, R. (2024). GREATS: Online selection of high-quality
721 data for LLM training in every iteration. *Advances in Neural Information Processing Systems*,
722 37:131197–131223.
- 723
- 724 Welbl, J., Liu, N. F., and Gardner, M. (2017). Crowdsourcing multiple choice science questions.
725 *ArXiv*, 1707.06209.
- 726 Wettig, A., Gupta, A., Malik, S., and Chen, D. (2024). QuRating: Selecting high-quality data for
727 training language models. *International Conference on Machine Learning*, pages 52915–52971.
- 728
- 729 Williamson, R. and Menon, A. (2019). Fairness risk measures. *International Conference on Machine*
730 *Learning*, pages 6786–6797.
- 731 Xia, M., Malladi, S., Gururangan, S., Arora, S., and Chen, D. (2024). LESS: Selecting influential
732 data for targeted instruction tuning. *International Conference on Machine Learning*, pages 54104–
733 54132.
- 734
- 735 Xie, S. M., Pham, H., Dong, X., Du, N., Liu, H., Lu, Y., Liang, P. S., Le, Q. V., Ma, T., and Yu, A. W.
736 (2023a). DoReMi: Optimizing data mixtures speeds up language model pretraining. *Advances in*
737 *Neural Information Processing Systems*, 36:69798–69818.
- 738 Xie, S. M., Santurkar, S., Ma, T., and Liang, P. S. (2023b). Data selection for language models via
739 importance resampling. *Advances in Neural Information Processing Systems*, 36:34201–34227.
- 740
- 741 Yu, Z., Das, S., and Xiong, C. (2024). MATES: Model-aware data selection for efficient pretraining
742 with data influence models. *Advances in Neural Information Processing Systems*, 37:108735–
743 108759.
- 744 Zellers, R., Holtzman, A., Bisk, Y., Farhadi, A., and Choi, Y. (2019). HellaSwag: Can a machine
745 really finish your sentence? *ArXiv*, 1905.07830.
- 746
- 747
- 748
- 749
- 750
- 751
- 752
- 753
- 754
- 755

Appendix

Table of Contents

A Statement for the Use of Large Language Models (LLMs)	16
B ADA-ESLM: adaptive confidence thresholding	16
C Risk-aware knowledge distillation with ESLM-KD	19
D Experiment details	19
D.1 Experimental setup	19
D.2 Pretraining corpus	19
D.3 Hardware & computational overhead	20
D.4 Baselines	21
D.5 Evaluation details	22
D.6 Reproducibility	22
E Additional experimental results	22
E.1 Validation loss versus training FLOPs/iterations results	22
E.2 Downstream performance evaluation results	22
E.3 Ablation for Instance selection vs Token-level Loss Shaping	26
F ESLM token selection analysis	26

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

A STATEMENT FOR THE USE OF LARGE LANGUAGE MODELS (LLMs)

We used LLMs solely as an assistive tool for polishing the writing and improving the readability of this paper. All research ideas, experimental design, implementation, and analysis were conducted by the authors. The paper itself directly studies an online, risk-aware batch selection method for improving the efficiency and distributional robustness of LLM pretraining.

B ADA-ESLM: ADAPTIVE CONFIDENCE THRESHOLDING

In Algorithm 2, we provide the algorithmic description of ADA-ESLM. Instead of using a fixed confidence level α throughout training, ADA-ESLM introduces a feedback-driven update mechanism that adjusts α based on the evolving difficulty of the training process. The underlying principle is to achieve a steady state through stabilizing the CVaR signal over time, allowing the model to gradually shift from broad token coverage to a more focused, high-risk subset. When CVaR increases over training intervals, it signals that the model is encountering more difficult (high-risk) examples, requiring a broader coverage. Conversely, a decrease in CVaR suggests that the model is improving on difficult tokens and can afford to focus more narrowly.

Concretely, at each evaluation step k (defined by the interval T_{eval}), ADA-ESLM measures the change in CVaR average tail token-level risk scores, which is a proxy for difficulty. Let CVaR_{α_k} denote the CVaR value at iteration k computed via (2). We define the normalized CVaR change, Δ_{norm} , via:

$$\Delta_{\text{norm}}(\alpha_k) := \frac{\text{CVaR}_{\alpha_k} - \text{CVaR}_{\alpha_{k-1}}}{\text{CVaR}_{\alpha_{k-1}} + \varepsilon},$$

which is a dimension and scale-independent feedback signal, and $\varepsilon > 0$ is a small constant for numerical stability. The controller updates the confidence level α multiplicatively using:

$$\alpha_{k+1} = \alpha_k \cdot \exp(-\gamma \cdot \Delta_{\text{norm}}(\alpha_k)),$$

where $\gamma > 0$ controls the update rate. The update rule captures the key intuition:

- If $\Delta_{\text{norm}} > 0$ (CVaR increases), then α is decreased to expand the input selection.
- If $\Delta_{\text{norm}} < 0$ (CVaR decreases), then α is increased to narrow focus to high-risk inputs.

This dynamic adjustment results in a form of *token-level curriculum learning* in which the model begins with broad exposure and progressively narrows focus to the most informative regions of the data. As we further show in Figure 7, ADA-ESLM gradually increases α over training and converges to a stable operating regime in the range $[0.1, 0.2]$ —a region empirically shown to yield a strong trade-off between training efficiency and data utility (Section 5.2, Figure 5c).

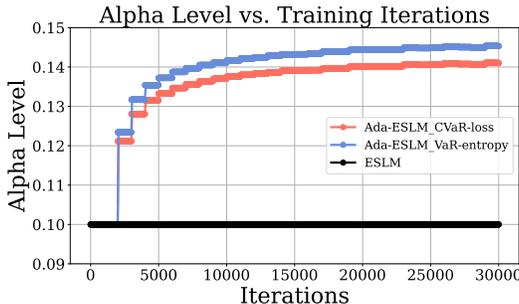


Figure 7: **ADA-ESLM confidence level (α) during training.** ADA-ESLM adjusts α dynamically using a CVaR-based controller to stabilize training. The learned α values converge to the $[0.1, 0.2]$ range—previously shown in Section 5.2 (Figure 5c) to balance training efficiency and data utilization.

In Table 2, we present the downstream performance of ADA-ESLM across standard benchmarks. Notably, the adaptive variant **ADA-ESLM-VaR-entropy** consistently outperforms both baseline methods (CLM, Rho-1, GREATS) and fixed- α ESLM variants. These results highlight the benefit of dynamically adjusting the token selection threshold during training, demonstrating that ADA-ESLM improves generalization while maintaining high training efficiency.

Algorithm 2 ADA-ESLM

1: **Input:** Language model θ , dataset \mathcal{D} , learning rate η , initial confidence level $\alpha_0 \in (0, 1)$, proxy L layers, sensitivity $\gamma > 0$, evaluation interval T_{eval} , batch size M , small constant $\varepsilon > 0$.

2: Initialize: $\text{CVaR}_0 \leftarrow 0$.

3: Initialize the list: $\text{CVaR_history} \leftarrow []$.

4: Append CVaR_0 to CVaR_history .

5: **for** each training iteration $k = 1, \dots, K$ **do**

6: Sample a batch of instances $\mathcal{B} = \{x_1, \dots, x_M\} \sim \mathcal{D}$ of length T $\{x_j^t\}_{t=1}^T, \forall j \in \{1, \dots, M\}$.

7: Compute per-instance scores $S_{\theta_k^L}(x)$ using early-exit θ_k^L proxy. */ Entropy or loss

8: Compute threshold $S_{\theta_k^L, \alpha}^{\text{VaR}} \leftarrow \text{VaR}_\alpha \left(\{S_{\theta_k^L}(x_j)\}_{j=1}^M \right)$ using (1).

9: $\bar{\mathcal{B}} \leftarrow \{x_j \in \mathcal{B} \mid S_{\theta_k^L}(x_j) \geq S_{\theta_k^L, \alpha}^{\text{VaR}}\}$. */ High-risk instance selection

10: Compute per-token statistics $S_{\theta_k}(x^t)$: */ Entropy or loss

$S_{\theta_k}(x_j^t) = \begin{cases} H_{\theta_k}(x_j^t) \text{ as in (i),} & \text{(VaR-entropy)} \\ \ell_{\theta_k}(x_j^t) \text{ as in (ii),} & \text{(CVaR-loss)} \end{cases}$

11: Compute threshold $S_{\theta_k, \alpha}^{\text{VaR}} \leftarrow \text{VaR}_\alpha \left(\{S_{\theta_k}(x_j^t)\}_{j \in \bar{\mathcal{B}}}\}_{t=1}^T \right)$ using (1).

12: $\tilde{\mathcal{B}} \leftarrow \{x_j^t \in \bar{\mathcal{B}} \mid S_{\theta_k}(x_j^t) \geq S_{\theta_k, \alpha}^{\text{VaR}}\}$. */ High-risk token selection

13: Compute loss over selected tokens: */ Shaped loss

$\mathcal{L}_{\tilde{\mathcal{B}}}(x; \theta_k) = \begin{cases} \mathbb{E}[\ell_{\theta_k}(x_j^t) \mid x_j^t \in \tilde{\mathcal{B}}], & \text{(VaR-entropy)} \\ \text{CVaR}_\alpha(\ell_{\theta_k}(x)) = \mathbb{E}[\ell_{\theta_k}(x_j^t) \mid x_j^t \in \tilde{\mathcal{B}}] \text{ using (2),} & \text{(CVaR-loss)} \end{cases}$

14: Update model parameters using optimizer O : $\theta_{k+1} \leftarrow O(\theta_k, \nabla_{\theta} \mathcal{L}_{\tilde{\mathcal{B}}}(x; \theta_k), \eta)$.

15: **if** $k \bmod T_{\text{eval}} = 0$ **then** */ Update α

16: Compute $\text{CVaR}_{\alpha_k} \leftarrow \text{CVaR}_{\alpha_k}(S_{\theta_k})$ using (2).

17: Retrieve $\text{CVaR}_{\alpha_{k-1}} \leftarrow \text{CVaR_history}[-1]$.

18: Compute normalized CVaR change:

$$\Delta_{\text{norm}}(\alpha_k) \leftarrow \frac{\text{CVaR}_{\alpha_k} - \text{CVaR}_{\alpha_{k-1}}}{|\text{CVaR}_{\alpha_{k-1}}| + \varepsilon}$$

19: Update confidence level: $\alpha_{k+1} \leftarrow \alpha_k \cdot \exp(-\gamma \cdot \Delta_{\text{norm}}(\alpha_k))$.

20: Append CVaR_{α_k} to CVaR_history .

21: **end if**

22: **end for**

23: **return** θ_K

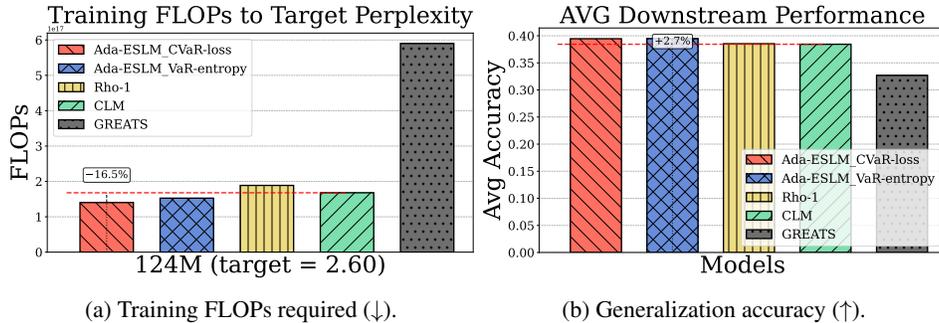


Figure 8: **ADA-ESLM efficiency and generalization performance.** (a): ADA-ESLM adaptively tunes the α level based on training dynamics, achieving the target validation (log) perplexity with fewer training FLOPs compared to baselines. (b): ADA-ESLM further improves generalization on downstream benchmarks reported in detail in Table 2, achieving higher average accuracy than baselines, trained under the same compute budget.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

Table 2: **Generalization performance of ADA-ESLM on downstream tasks.** All models (124M) are pretrained under a $\sim 3E17$ FLOPs budget on SlimPajama-6B-Unif mixture. We report the best observed accuracy_(standard error) or exact match if provided, during training. **Highlighted** values indicate the best performance. Results demonstrate that the dynamic CVaR-driven adjustment of α level leads to improved generalization over baselines, particularly with **ADA-ESLM-VaR-entropy** setting.

Benchmark	# Shots	Method (124M)						
		ADA-ESLM-CVaR-loss	ADA-ESLM-VaR-entropy	ESLM-CVaR-loss	ESLM-VaR-entropy	CLM	Rho-1	GREATS
ARC-E (Clark et al., 2018)	0-shot	0.3728 _(0.0099)	0.3884 _(0.01)	0.3712 _(0.0099)	0.3766 _(0.0099)	0.3644 _(0.0099)	0.3657 _(0.0099)	0.3236 _(0.0096)
LAMBADA (Paperno et al., 2016)	5-shot	0.1631 _(0.0051)	0.1628 _(0.0051)	0.1595 _(0.0051)	0.1721 _(0.0053)	0.1701 _(0.0051)	0.1680 _(0.005)	0.0254 _(0.002)
SciQ (Welbl et al., 2017)	5-shot	0.6960 _(0.0146)	0.7060 _(0.0144)	0.7140 _(0.0143)	0.7100 _(0.0144)	0.6970 _(0.0145)	0.7000 _(0.0145)	0.4350 _(0.0157)
HellaSwag (Zellers et al., 2019)	5-shot	0.2976 _(0.0046)	0.2952 _(0.0046)	0.2930 _(0.0045)	0.2964 _(0.0046)	0.2901 _(0.0045)	0.2893 _(0.0045)	0.2621 _(0.0044)
TriviaQA (Joshi et al., 2017)	1-shot	0.0082 _(0.0007)	0.0128 _(0.0008)	0.0128 _(0.0001)	0.0066 _(0.0006)	0.0078 _(0.0007)	0.0090 _(0.0007)	0.0007 _(0.0002)
COPA (Wang et al., 2019)	5-shot	0.6600 _(0.0476)	0.6300 _(0.0485)	0.6400 _(0.0482)	0.6200 _(0.0488)	0.6200 _(0.0488)	0.6200 _(0.0488)	0.6400 _(0.0482)
MultiRC (Wang et al., 2019)	5-shot	0.5600 _(0.0071)	0.5668 _(0.0071)	0.5633 _(0.0071)	0.5429 _(0.0071)	0.5338 _(0.0072)	0.5338 _(0.0072)	0.5497 _(0.0071)
OpenBookQA (Mihaylov et al., 2018)	5-shot	0.174 _(0.017)	0.178 _(0.0171)	0.164 _(0.0166)	0.170 _(0.0168)	0.166 _(0.0167)	0.164 _(0.0166)	0.148 _(0.0159)
PIQA (Bisk et al., 2020)	5-shot	0.6191 _(0.0113)	0.6115 _(0.0114)	0.6240 _(0.0113)	0.6191 _(0.0113)	0.6099 _(0.0114)	0.6180 _(0.0113)	0.5571 _(0.0116)
Average (\uparrow)		0.39453	0.39461	0.39353	0.39041	0.38434	0.38531	0.32684

C RISK-AWARE KNOWLEDGE DISTILLATION WITH ESLM-KD

We provide the implementation for our knowledge distillation setup, namely ESLM-KD, in Algorithm 3. The student model θ computes per-token risk scores over each batch, and high-risk tokens are selected via VaR_α thresholding. The student is then supervised only on these informative tokens using a combined loss: a weighted sum of KL divergence to the teacher (ϕ) and standard cross-entropy.

In our experiments (Section 5.1.1), we used a 124M GPT-2 model pretrained with the CLM objective (checkpoint 40,000) as the teacher to train a 774M student models on the SlimPajama-6B-Unif dataset. Based on hyperparameter tuning, we set the distillation weight to $\lambda = 0.5$ and the teacher temperature to $\rho = 1.0$. We compare ESLM-KD against three 774M baselines with the same teacher model: standard CLM training, Dense-KD (dense knowledge distillation without token selection), and SALT (Rawat et al., 2024), a two-staged distillation method which employs distillation in the first stage and then transitions to standard pretraining. For the SALT baseline, we set the distillation iterations to 12,000. We trained all distillation-based models with a compute budget of 1E18 FLOPs.

The experimental results in Figure 5b (Section 5.1) show that ESLM-KD achieves the target validation loss with significantly less training FLOPs, demonstrating its efficiency and effectiveness in large-scale distillation. Table 3 further compares the generalization performance of ESLM-KD against dense distillation using the same teacher model. The results show that integrating risk-aware token selection into distillation not only reduces compute cost but also improves downstream accuracy over full-token distillation.

Algorithm 3 ESLM-KD

1: **Input:** Teacher LM parameters ϕ , student LM parameters θ , dataset \mathcal{D} , learning rate η , confidence level $\alpha \in (0, 1)$, batch size M , teacher temperature $\rho > 0$, distillation loss weight $\lambda \in [0, 1]$.

2: **for** each training iteration $k = 1, \dots, K$ **do**

3: Sample a batch of tokens $\mathcal{B} = \{x_1, \dots, x_M\} \sim \mathcal{D}$.

4: Compute per-token statistics $S_\theta(x_j)$ using the student model: */ Entropy or loss

$$S_\theta(x_j) = \begin{cases} H_\theta(x_j) & \text{as in (i), (VaR-entropy)} \\ \ell_\theta(x_j) & \text{as in (ii), (CVaR-loss)} \end{cases}$$

5: Compute VaR threshold: $S_{\theta, \alpha}^{\text{VaR}} \leftarrow \text{VaR}_\alpha(\{S_\theta(x_j)\}_{j=1}^M)$ using (1).

6: Select high-risk tokens: $\tilde{\mathcal{B}} \leftarrow \{x_j \in \mathcal{B} \mid S_\theta(x_j) \geq S_{\theta, \alpha}^{\text{VaR}}\}$.

7: Compute combined student loss on selected tokens: */ Distillation + cross-entropy loss

$$\mathcal{L}_{\text{ESLM-KD}} = \frac{1}{|\tilde{\mathcal{B}}|} \sum_{x_j \in \tilde{\mathcal{B}}} [\lambda \cdot \text{KL}(P_\rho^\phi(x_j \mid x_{<j}) \parallel P_\rho^\theta(x_j \mid x_{<j})) + (1 - \lambda) \cdot \ell_{\theta_k}(x_j)]$$

8: Update student parameters using optimizer O : $\theta_{k+1} \leftarrow O(\theta_k, \nabla_\theta \mathcal{L}_{\text{ESLM-KD}}, \eta)$.

9: **end for**

10: **return** θ_K .

D EXPERIMENT DETAILS

D.1 EXPERIMENTAL SETUP

We set the training hyperparameters as in Table 4. We train GPT-2 models (Radford et al., 2019) of sizes 124M, 350M, 774M, and 1.5B parameters, with architecture details reported in Table 5.

D.2 PRETRAINING CORPUS

We utilize the SlimPajama-6B dataset for our experiments. SlimPajama-6B (Soboleva et al., 2023) mixture consists of seven data domains: {Arxiv, Book, CommonCrawl, C4, Github, Stackexchange, Wikipedia} with two weighted versions: uniform domain weights (SlimPajama-6B-Unif) and DoReMi (Xie et al., 2023a) domain weights (SlimPajama-6B-DoReMi).

Table 3: **Generalization performance of ESLM-KD on downstream tasks.** All models (774M) are pretrained under a $\sim 1\text{E}18$ FLOPs budget on SlimPajama-6B-Unif mixture, using the same teacher model. We report the best observed accuracy_(standard error) or exact match if provided, during training. **Highlighted** values indicate the best performance.

Benchmark	Method (774M)				
	# Shots	ESLM-KD-CVaR-loss	ESLM-KD-VaR-entropy	Dense-KD	SALT
ARC-E (Clark et al., 2018)	0-shot	0.3901 _(0.01)	0.3935 _(0.01)	0.3909 _(0.01)	0.3947 _(0.01)
LAMBADA (Paperno et al., 2016)	5-shot	0.2472 _(0.006)	0.2480 _(0.006)	0.2429 _(0.006)	0.2258 _(0.0058)
SciQ (Welbl et al., 2017)	5-shot	0.766 _(0.0134)	0.773 _(0.0133)	0.759 _(0.0135)	0.770 _(0.0133)
HellaSwag (Zellers et al., 2019)	5-shot	0.3200 _(0.0047)	0.3189 _(0.0047)	0.3179 _(0.0046)	0.3313 _(0.0047)
TriviaQA (Joshi et al., 2017)	1-shot	0.0273 _(0.0012)	0.0280 _(0.0012)	0.0231 _(0.0011)	0.0299 _(0.0013)
COPA (Wang et al., 2019)	5-shot	0.68 _(0.0469)	0.66 _(0.0476)	0.69 _(0.0465)	0.67 _(0.0473)
MultiRC (Wang et al., 2019)	5-shot	0.5408 _(0.0072)	0.5406 _(0.0072)	0.5360 _(0.0072)	0.5420 _(0.0072)
OpenBookQA (Mihaylov et al., 2018)	5-shot	0.270 _(0.0199)	0.290 _(0.0203)	0.278 _(0.0201)	0.278 _(0.02)
PiQA (Bisk et al., 2020)	5-shot	0.6300 _(0.0113)	0.6245 _(0.0113)	0.6327 _(0.0112)	0.6322 _(0.0113)
Average (\uparrow)		0.4301	0.4307	0.4299	0.4304

Table 4: Training and evaluation hyperparameters used in all experiments.

Hyperparameter	Value
ESLM-Specific	
Confidence level (α)	Optimized over {0.05, 0.1, 0.2}
Early-exit layers for proxy phase (L)	Optimized over {1, 2, 3}
General Setup	
Mini-batch size (M) in tokens	{ 12 (124M), 8 (350M/774M), 4 (1.5B) } \times 1024
Gradient accumulation steps	40
Effective batch size in tokens	{ 480 (124M), 320 (350M/774M), 160 (1.5B) } \times 1024
Sequence length (T)	1024
Vocabulary size ($ \mathcal{V} $)	50304
Dropout	0
Evaluation interval (T_{eval})	1000 iterations
Evaluation steps	200 iterations
Optimization	
Optimizer (O)	AdamW with $\beta_1 = 0.9, \beta_2 = 0.95$
Learning rate schedule	Cosine annealing with warmup
Max. learning rate (η)	0.0006 (124M/350M), 0.0001 (774M/1.5B)
Min. learning rate	0.00006
Warmup steps	2000
Decay iterations	200,000
Weight decay	0.1
Gradient clipping	1.0
Knowledge Distillation	
Teacher temperature (ρ)	1.0
Distillation loss weight (λ)	0.5

In Table 6, we report the domain weights for the experiments under SlimPajama-6B (Soboleva et al., 2023) mixture, using DoReMi (Xie et al., 2023a) and uniform weights.

D.3 HARDWARE & COMPUTATIONAL OVERHEAD

All experiments were conducted on the HTCondor-managed cluster equipped with NVIDIA A100 GPUs (80GB). Model pretraining and evaluation were parallelized using PyTorch’s Distributed Data Parallel (DDP) framework (Paszke et al., 2019) with the NCCL backend and mixed-precision (bfloat16) training. We used $4 \times$ A100 GPUs for experiments on the SlimPajama-6B mixtures.

Table 5: Architecture hyperparameters for GPT-2 model sizes.

Model size	Layers	Attention heads	Embed dimension
124M	12	12	768
350M	24	16	1024
774M	36	20	1280
1.5B	48	25	1600

Table 6: Domain weights used for experiments on the SlimPajama-6B mixture.

Domain	DoReMi	Unif
Arxiv	0.04235	0.1428
Book	0.08201	0.1428
CC	0.381	0.1428
C4	0.1141	0.1428
Github	0.0654	0.1428
Stackexchange	0.0847	0.1428
Wikipedia	0.2305	0.1428

Runtime analysis. In Table 7, we compare the wall-clock time of 124M models trained on the SlimPajama-6B-Unif, for convergence to a target validation loss value of 2.60. While ESLM achieves substantial reductions in training FLOPs, reaches lower validation loss, and stronger downstream performance, it incurs higher wall-clock time compared to standard training. However, it remains significantly more efficient than GREATS and Rho-1 baselines—nearly twice as fast as Rho-1 and over 8× faster than GREATS. We attribute this overhead to mismatches between sparse training operations and current hardware optimizations.

As also explained by the recent efficiency works (Lasby et al., 2024), achieving practical acceleration from FLOPs savings requires optimized hardware support for sparse backpropagation—an engineering direction which is beyond our paper’s scope. Although ESLM’s per-token risk scores are computed during the forward pass via sorting with $O(M \log M)$ complexity per batch—without requiring additional external inference or backpropagation—its proxy pass requires a shallow forward call, and VaR-based token filtering introduces sparsity into the training process. This sparsity, while beneficial for compute efficiency, leads to irregular and fragmented backpropagation paths that underutilize the dense compute capabilities of modern accelerators. Unlike the uniform operations of standard CLM, ESLM’s selective masking disrupts efficient tensor fusion, resulting in slower wall-clock runtime despite using fewer FLOPs. Nonetheless, ESLM provides a favorable trade-off: improved efficiency per FLOP and enhanced generalization. **We expect future work leveraging sparsity-aware hardware or sparse accelerators to further reduce this overhead (Lasby et al., 2024) and unlock the full potential of selective training.**

D.4 BASELINES

We identify the baseline methods against which we compare our ESLM approach, specifically from online batch selection methods for LLM pretraining and standard training as discussed in Section 5. We provide the baseline implementation details below:

- For the Rho-1 baseline (Lin et al., 2024), we used pretrained GPT-2 models trained via CLM objective as the reference model. Since training a high-quality reference model is the main bottleneck of the Rho-1 method, we used the last checkpoints of pretrained models as proxy

Table 7: Runtime comparison of 124M models trained on SlimPajama-6B-Unif, for convergence to a target validation loss value of 2.60. The overhead compared to the standard training is mainly due to the mismatch between sparsity introduced via instance & token selection and current hardware optimizations.

Method	Wall-clock time (hrs)
ESLM-VaR-entropy	11.73
ESLM-CVaR-loss	9.91
CLM	5.33
Rho-1	17.91
GREATS	99.89

1134 models. For pretraining on SlimPajama-6B mixtures (Unif and DoReMi), we used the last saved
1135 checkpoint of CLM GPT-2 models as the reference models. Specifically, we utilized 40000,
1136 30000, and 30000 checkpoints for 124M, 350M, 774M models, respectively. We set the loss
1137 threshold parameter to 0.1. For Rho-1’s total FLOPs calculation, we include additional FLOPs
1138 from the forward call on the reference model.

- 1139 • For the GREATS baseline, we follow the original setup by Wang et al. (2024), using a small
1140 validation set ($0.5 \times$ the batch size) and setting the batch selection budget to 0.9, aligning with
1141 ESLM’s $\alpha = 0.1$ level. To compute the training FLOPs for GREATS, we include the forward
1142 passes on both training and validation inputs, the backward pass through linear layers to obtain per-
1143 example gradients with respect to pre-activation outputs, and the additional FLOPs for computing
1144 ghost inner products. While GREATS is evaluated only on the GPT-2 124M model in the original
1145 paper, we also restrict our comparison to this setting. Despite adopting their ghost inner product
1146 optimization, we found the method to be highly memory-intensive when scaling to larger models,
1147 and it could not run stably beyond 124M size.

1148 D.5 EVALUATION DETAILS

1149 We evaluate pretrained models on a suite of standard language understanding benchmarks in the
1150 zero-shot and few-shot settings, using the `lm-evaluation-harness` evaluation suite (Gao et al.,
1151 2024), including HellaSwag (Zellers et al., 2019), LAMBADA (Paperno et al., 2016), ARC-Easy
1152 (Clark et al., 2018), TriviaQA (Joshi et al., 2017), SciQ (Welbl et al., 2017), COPA (Wang et al.,
1153 2019), MultiRC (Wang et al., 2019), OpenBookQA (Mihaylov et al., 2018), and PiQA (Bisk et al.,
1154 2020) tasks. We used the default settings provided by `lm-evaluation-harness`, which means
1155 all evaluations are performed on held-out validation splits, or test splits if provided, and standard
1156 errors are calculated using bootstrapping. Accuracy (norm if provided) or exact match is used as the
1157 primary metric.

1158 D.6 REPRODUCIBILITY

1159 For reproducibility, we provide the algorithmic description of the ESLM method in Algorithm 1.
1160 Our implementation builds on the open-source `NanoGPT` codebase (Karpathy, 2022). To handle
1161 training on the SlimPajama-6B dataset mixture, we adapted the open-source code of DoReMi (Xie
1162 et al., 2023a) and DoGE (Fan et al., 2023). As we detail in Section 3, we estimate the training
1163 FLOPs based on the theoretical estimate by Kaplan et al. (2020); Chowdhery et al. (2023). To skip
1164 redundant gradient computation as a result of token-level loss shaping, we use PyTorch’s backward
1165 hook mechanism (`Tensor.register_hook(hook)`) for gradient masking. For downstream
1166 evaluation, we utilize the publicly available `lm-evaluation-harness` suite (Gao et al., 2024).
1167 Our open-source implementation, along with references to the adapted codebases, will be released in
1168 the camera-ready version.

1169 E ADDITIONAL EXPERIMENTAL RESULTS

1170 In this section, we report additional experimental results, showing validation perplexity convergence
1171 in compute space (Appendix E.1) and generalization performance in downstream benchmark tasks
1172 (Appendix E.2) on different datasets across model sizes.

1173 E.1 VALIDATION LOSS VERSUS TRAINING FLOPS/ITERATIONS RESULTS

1174 As shown in Figures 9-10, ESLM variants consistently accelerate validation loss convergence in the
1175 compute space, requiring fewer training FLOPs to achieve comparable or superior validation loss
1176 relative to baseline models. This efficiency gain holds across diverse pretraining corpora and model
1177 scales, highlighting the robustness of ESLM across settings.

1178 E.2 DOWNSTREAM PERFORMANCE EVALUATION RESULTS

1179 Tables 8–12 present the generalization performance of ESLM models ranging from 124M to 774M
1180 parameters, trained on different mixtures and evaluated against baseline models on downstream

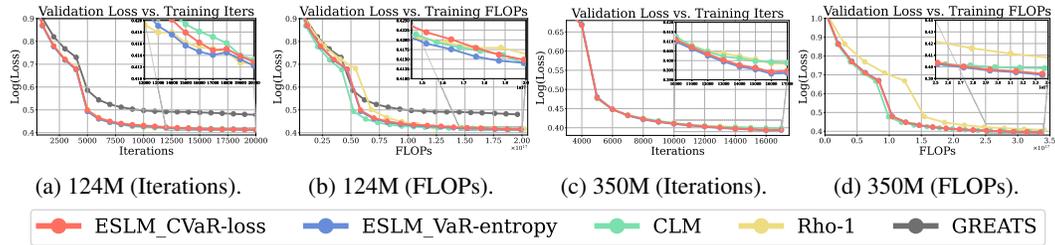


Figure 9: **Validation loss convergence on SlimPajama-6B-DoReMi.** We report convergence of validation loss versus training FLOPs/iterations of models trained on SlimPajama-6B-DoReMi mixture. ESLM variants provide faster convergence to lower loss values in terms of iterations, while requiring fewer FLOPs than baselines at the convergence.

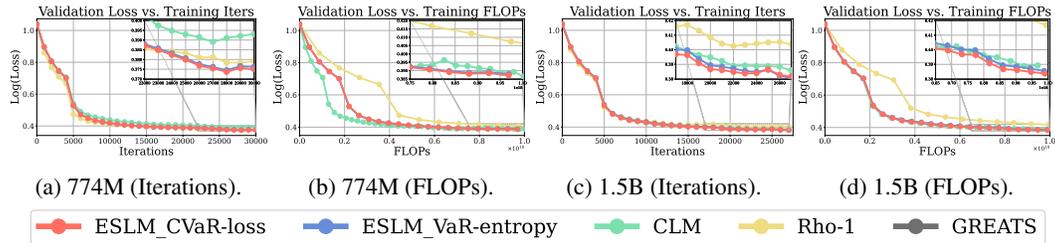


Figure 10: **Validation loss convergence on SlimPajama-6B-Unif.** We report convergence of validation loss versus training FLOPs/iterations of models trained on SlimPajama-6B-Unif mixture, for 774M and 1.5B models. ESLM variants provide faster convergence to lower loss values in terms of iterations, while requiring fewer FLOPs than baselines at the convergence.

benchmarks. All models are trained under a fixed compute budget measured in training FLOPs. In nearly all settings, ESLM variants consistently outperform baselines, achieving higher average downstream accuracy. While domain mixture weights influence absolute performance, ESLM maintains a consistent advantage, demonstrating that it is not only an efficient and simple approach but also yields better generalization quality.

Table 8: **Generalization performance of 124M models trained on SlimPajama-6B-DoReMi.** All models are pretrained under a $\sim 3E17$ FLOPs budget. We report the best observed accuracy (standard error) or exact match if provided, during training. **Highlighted** values indicate the best performance.

Benchmark	# Shots	Method (124M)				
		ESLM-CVaR-loss	ESLM-VaR-entropy	CLM	Rho-1	GREATS
ARC-E (Clark et al., 2018)	0-shot	0.3876 _(0.01)	0.3926 _(0.01)	0.3838 _(0.01)	0.3733 _(0.0099)	0.3190 _(0.0096)
LAMBADA (Paperno et al., 2016)	5-shot	0.1738 _(0.0053)	0.1727 _(0.0053)	0.1581 _(0.0051)	0.1672 _(0.0052)	0.0362 _(0.0026)
SciQ (Welbl et al., 2017)	5-shot	0.731 _(0.014)	0.718 _(0.0142)	0.735 _(0.014)	0.723 _(0.0142)	0.465 _(0.0158)
HellaSwag (Zellers et al., 2019)	5-shot	0.2936 _(0.0045)	0.2924 _(0.0045)	0.2945 _(0.0045)	0.2905 _(0.0045)	0.2639 _(0.0044)
TriviaQA (Joshi et al., 2017)	1-shot	0.0184 _(0.001)	0.0145 _(0.0009)	0.0100 _(0.0007)	0.0112 _(0.0008)	0.0007 _(0.0002)
COPA (Wang et al., 2019)	5-shot	0.65 _(0.0479)	0.66 _(0.0476)	0.66 _(0.0476)	0.64 _(0.0482)	0.65 _(0.0479)
MultiRC (Wang et al., 2019)	5-shot	0.5455 _(0.0072)	0.5367 _(0.0072)	0.5449 _(0.0072)	0.5486 _(0.0071)	0.5309 _(0.0072)
OpenBookQA (Mihaylov et al., 2018)	5-shot	0.176 _(0.017)	0.164 _(0.0166)	0.174 _(0.017)	0.164 _(0.0166)	0.148 _(0.0159)
PiQA (Bisk et al., 2020)	5-shot	0.6169 _(0.0113)	0.6175 _(0.0113)	0.6017 _(0.0114)	0.6033 _(0.0114)	0.5489 _(0.0116)
Average (\uparrow)		0.3992	0.3964	0.3957	0.3912	0.3291

Table 9: **Generalization performance of 350M models trained on SlimPajama-6B-Unif.** All models are pretrained under a $\sim 3.5E17$ FLOPs budget. We report the best observed accuracy (standard error) or exact match if provided, during training. **Highlighted** values indicate the best performance.

Benchmark	# Shots	Method (350M)			
		ESLM-CVaR-loss	ESLM-VaR-entropy	CLM	Rho-1
ARC-E (Clark et al., 2018)	0-shot	0.4078 _(0.0101)	0.3973 _(0.01)	0.4023 _(0.0101)	0.4006 _(0.0101)
LAMBADA (Paperno et al., 2016)	5-shot	0.2289 _(0.0059)	0.2070 _(0.0056)	0.2095 _(0.0057)	0.1993 _(0.0056)
SciQ (Welbl et al., 2017)	5-shot	0.749 _(0.0137)	0.768 _(0.0134)	0.744 _(0.0138)	0.760 _(0.0135)
HellaSwag (Zellers et al., 2019)	5-shot	0.3301 _(0.0047)	0.3298 _(0.0047)	0.3242 _(0.0047)	0.3207 _(0.0047)
TriviaQA (Joshi et al., 2017)	1-shot	0.0338 _(0.0013)	0.0271 _(0.0012)	0.0329 _(0.0013)	0.0263 _(0.0012)
COPA (Wang et al., 2019)	5-shot	0.68 _(0.0469)	0.68 _(0.0469)	0.68 _(0.0469)	0.67 _(0.0473)
MultiRC (Wang et al., 2019)	5-shot	0.5482 _(0.0071)	0.5556 _(0.0071)	0.5492 _(0.0071)	0.5676 _(0.0071)
OpenBookQA (Mihaylov et al., 2018)	5-shot	0.276 _(0.02)	0.288 _(0.0203)	0.280 _(0.0201)	0.284 _(0.0202)
PiQA (Bisk et al., 2020)	5-shot	0.6398 _(0.0112)	0.6420 _(0.0112)	0.6349 _(0.0112)	0.6322 _(0.0113)
Average (\uparrow)		0.4326	0.4327	0.4284	0.4289

Table 10: **Generalization performance of 350M models trained on SlimPajama-6B-DoReMi.** All models are pretrained under a $\sim 3.5E17$ FLOPs budget. We report the best observed accuracy (standard error) or exact match if provided, during training. **Highlighted** values indicate the best performance.

Benchmark	# Shots	Method (350M)			
		ESLM-CVaR-loss	ESLM-VaR-entropy	CLM	Rho-1
ARC-E (Clark et al., 2018)	0-shot	0.4170 _(0.0101)	0.4048 _(0.0101)	0.4196 _(0.0101)	0.4090 _(0.0101)
LAMBADA (Paperno et al., 2016)	5-shot	0.2400 _(0.006)	0.2344 _(0.0059)	0.2361 _(0.0059)	0.2144 _(0.0057)
SciQ (Welbl et al., 2017)	5-shot	0.775 _(0.0132)	0.794 _(0.0128)	0.782 _(0.0131)	0.773 _(0.0133)
HellaSwag (Zellers et al., 2019)	5-shot	0.3292 _(0.0047)	0.3271 _(0.0047)	0.3304 _(0.0047)	0.3213 _(0.0047)
TriviaQA (Joshi et al., 2017)	1-shot	0.0357 _(0.0014)	0.0431 _(0.0015)	0.0414 _(0.0015)	0.0395 _(0.0015)
COPA (Wang et al., 2019)	5-shot	0.68 _(0.0469)	0.70 _(0.0461)	0.69 _(0.0465)	0.68 _(0.0469)
MultiRC (Wang et al., 2019)	5-shot	0.5457 _(0.0072)	0.5645 _(0.0071)	0.5548 _(0.0071)	0.5558 _(0.0071)
OpenBookQA (Mihaylov et al., 2018)	5-shot	0.286 _(0.0202)	0.284 _(0.0202)	0.280 _(0.0201)	0.286 _(0.0202)
PiQA (Bisk et al., 2020)	5-shot	0.6289 _(0.0113)	0.6414 _(0.0112)	0.6354 _(0.0112)	0.6354 _(0.0112)
Average (\uparrow)		0.4375	0.4437	0.4410	0.4349

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

Table 11: **Generalization performance of 774M models trained on SlimPajama-6B-Unif.** All models are pretrained under a $\sim 1E18$ FLOPs budget. We report the best observed accuracy (standard error) or exact match if provided, during training. **Highlighted** values indicate the best performance.

Benchmark	Method (774M)				
	# Shots	ESLM-CVaR-loss	ESLM-VaR-entropy	CLM	Rho-1
ARC-E (Clark et al., 2018)	0-shot	0.4040 _(0.0101)	0.4061 _(0.0101)	0.4082 _(0.0101)	0.3985 _(0.01)
LAMBADA (Paperno et al., 2016)	5-shot	0.2408 _(0.006)	0.2410 _(0.006)	0.2336 _(0.0059)	0.2404 _(0.006)
SciQ (Welbl et al., 2017)	5-shot	0.763 _(0.0135)	0.760 _(0.0135)	0.762 _(0.0135)	0.780 _(0.0131)
HellaSwag (Zellers et al., 2019)	5-shot	0.3348 _(0.0047)	0.3399 _(0.0047)	0.3333 _(0.0047)	0.3332 _(0.0047)
TriviaQA (Joshi et al., 2017)	1-shot	0.0370 _(0.0014)	0.0315 _(0.0013)	0.0349 _(0.0014)	0.0298 _(0.0013)
COPA (Wang et al., 2019)	5-shot	0.67 _(0.0473)	0.68 _(0.0469)	0.69 _(0.0465)	0.69 _(0.0465)
MultiRC (Wang et al., 2019)	5-shot	0.5474 _(0.0071)	0.5622 _(0.0071)	0.5591 _(0.0071)	0.5680 _(0.0071)
OpenBookQA (Mihaylov et al., 2018)	5-shot	0.284 _(0.0202)	0.278 _(0.0201)	0.278 _(0.0201)	0.280 _(0.0201)
PiQA (Bisk et al., 2020)	5-shot	0.6458 _(0.0112)	0.6468 _(0.0112)	0.6436 _(0.0112)	0.6392 _(0.0112)
Average (\uparrow)		0.4363	0.4383	0.4380	0.4399

Table 12: **Generalization performance of 774M models trained on SlimPajama-6B-DoReMi.** All models are pretrained under a $\sim 1E18$ FLOPs budget. We report the best observed accuracy (standard error) or exact match if provided, during training. **Highlighted** values indicate the best performance.

Benchmark	Method (774M)				
	# Shots	ESLM-CVaR-loss	ESLM-VaR-entropy	CLM	Rho-1
ARC-E (Clark et al., 2018)	0-shot	0.4132 _(0.0101)	0.4158 _(0.0101)	0.4128 _(0.0101)	0.4141 _(0.0101)
LAMBADA (Paperno et al., 2016)	5-shot	0.2437 _(0.006)	0.2400 _(0.006)	0.2124 _(0.0057)	0.2229 _(0.0058)
SciQ (Welbl et al., 2017)	5-shot	0.799 _(0.0127)	0.801 _(0.0126)	0.78 _(0.0131)	0.8 _(0.0127)
HellaSwag (Zellers et al., 2019)	5-shot	0.3417 _(0.0047)	0.3383 _(0.0047)	0.3366 _(0.0047)	0.3382 _(0.0047)
TriviaQA (Joshi et al., 2017)	1-shot	0.0457 _(0.0016)	0.0470 _(0.0016)	0.0412 _(0.0015)	0.0388 _(0.0014)
COPA (Wang et al., 2019)	5-shot	0.71 _(0.0456)	0.69 _(0.0465)	0.68 _(0.0469)	0.67 _(0.0473)
MultiRC (Wang et al., 2019)	5-shot	0.5435 _(0.0072)	0.5602 _(0.0071)	0.5680 _(0.0071)	0.5470 _(0.0071)
OpenBookQA (Mihaylov et al., 2018)	5-shot	0.288 _(0.0203)	0.294 _(0.0204)	0.28 _(0.0201)	0.284 _(0.0202)
PiQA (Bisk et al., 2020)	5-shot	0.6360 _(0.0112)	0.6338 _(0.0112)	0.6430 _(0.0112)	0.6289 _(0.0113)
Average (\uparrow)		0.4467	0.4466	0.4393	0.4382

E.3 ABLATION FOR INSTANCE SELECTION VS TOKEN-LEVEL LOSS SHAPING

In Figure 11, we decompose the FLOPs savings of ESLM into contributions from the proxy instance-selection phase and the token-level loss-shaping phase on a 124M model trained on the SlimPajama-6B-Unif mixture. The instance selection phase provides the primary reduction by discarding low-value sequences early, while token-level shaping adds complementary savings by pruning redundant gradient updates within the retained sequences. Together, these two components improve pretraining efficiency over standard CLM training without degrading learning convergence.

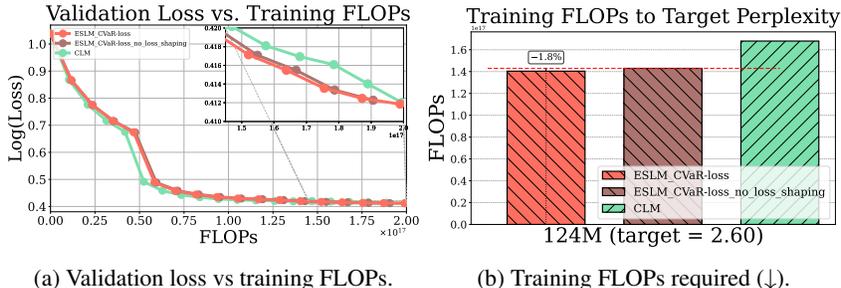


Figure 11: **Ablation on instance selection vs token-level loss shaping.** Applying token-level loss shaping further helps in compute-efficient convergence to the target validation loss with fewer FLOPs compared to the instance-only selection version.

F ESLM TOKEN SELECTION ANALYSIS

To better understand the behavior of ESLM, we conduct a qualitative analysis of token selection during pretraining. Specifically, we compare the selection patterns of two ESLM variants on the same SlimPajama-6B validation sequences. Figures 13-14 present examples where **highlighted tokens** represent those selected for backpropagation by 124M models using a fixed confidence level $\alpha = 0.1$ (i.e., top 90% high-risk tokens retained in the objective). In Figures 15-16, we further show selected tokens by 774M models under $\alpha = 0.2$ (i.e., top 80% high-risk tokens retained). We observe that both variants prioritize rare or informative tokens—such as named entities, foreign words, and domain-specific phrases—but differ in the nature of the signals they capture:

- **ESLM-VaR-entropy** emphasizes tokens associated with high predictive uncertainty, often selecting structurally or semantically transitional words, including common function words (e.g., “the”, “and”, “of”), punctuation, and formatting artifacts when they appear in unpredictable or shifting contexts. For instance, in the second passage (Figure 15), the model selects not only semantically meaningful words such as “anxiety-inducing”, “consumer-driven”, but also emphasizes “the”, “.”, and “with” in contexts where uncertainty over their grammatical role or continuation is high.
- **ESLM-CVaR-loss** instead selects tokens that incur high training loss—typically semantically complex or underfit tokens. In the first passage (Figure 16), we observe selection of technical terms such as “alkanes”, “aminated”; but tend to avoid repeating tokens such as “eq”. This variant tends to avoid punctuation and common syntactic tokens unless they contribute directly to high loss.

Figure 12 further illustrates the frequency of top-20 selected tokens by 774M ESLM models, from the validation examples given in Figures 15-16. The results reveal that as we allow for more tokens to be selected (α decreases), **ESLM-VaR-entropy** selects syntactically ambiguous tokens such as punctuations more than **ESLM-CVaR-loss**, reflecting its sensitivity to positional and contextual ambiguity, even in high-frequency tokens.

Crucially, this overall analysis also highlights the strength of token-level selection: ESLM captures the informativeness within sequences, in contrast to instance-level methods such as GREATS that filter entire examples. As a result, ESLM preserves valuable learning signals that would otherwise be discarded, offering a more fine-grained and efficient form of selective pretraining.

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

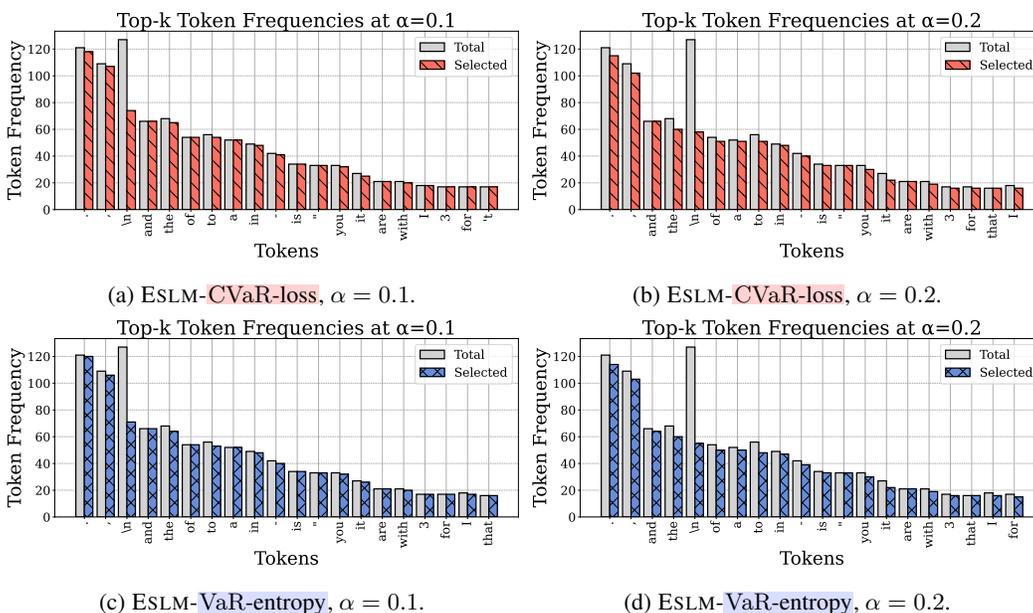


Figure 12: **Frequency of top-20 tokens selected by ESLM (774M) variants from validation sequences shown in Figures 15-16.** As α decreases (more tokens are selected from the batch), the **ESLM-VaR-entropy** emphasizes punctuation tokens more than **ESLM-CVaR-loss**, reflecting its sensitivity to positional and contextual ambiguity, even in high-frequency tokens.

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

Example 1 from domain: cc

ens, Lon Chaney Jr. Shiinomi Gakuen / The Shiinomi School (1955) Hiroshi Shimizu, Kyōko Kagawa, Yukiko Shimazaki, Jūjirō Kichi Uno, Drama Oliver Twist (1982) Clive Donner, George C. Scott, Tim Curry, Michael Hordern, Crime, Drama Cobra Woman (1944) Robert Siodmak, Maria Montez, Jon Hall, Sabu Imburnal (2008) Sherad Anthony Sanchez, Brian Monterola, Jelieta Mariveles-Ruca, Allen Lumanog Nemuri Kyoshiro 13: The Full Moon Swordsman (1969) Kazuo Mori, Hiroki Matsukata, Tomomi Satō, Sanae Nakahara, Action, Drama Shake Hands with the Devil (1959) Michael Anderson, James Cagney, Don Murray, Dana Wynter Sopyonje (1993) Kwon-taek Im, Myung-gon Kim, Jung-hae Oh, Kyu-chul Kim The American Soldier (1970) Rainer Werner Fassbinder, Karl Scheydtt, Elga Sorbas, Jan George, Drama Mikey and Nicky (1976) Elaine May, Peter Falk, John Cassavetes, Ned Beatty Małgorzata, morska / The Little Mermaid (1976) Karel Kachyna, Miroslava Safránková, Radovan Lukavský, Petr Svoboda, Family, Drama, Romance Strategic Air Command (1955) Anthony Mann, James Stewart, June Allyson, Frank Lovejoy, Action, Drama, War De Passagem / Passing By (2003) Ricardo Elias, Lohan Brandão, Thiago de Mello, Wilma de Souza, Drama The Citadel (1938) King Vidor, Robert Donat, Rosalind Russell, Ralph Richardson Ladies in Retirement (1941) Charles Vidor, Ida Lupino, Louis Hayward, Evelyn Keyes Madame Satō (2002) Karim Amrani, Lázaro Ramos, Marcelia Cartaxo, Flavio Bauraquí, Biography, Crime, Drama Youth in Fury (1960) Masahiro Shinoda, Shin'ichirō Mikami, Shima Iwashita, Kayoko Honoo Night Plane from Chungking (1943) Ralph Murphy, Robert Preston, Ellen Drew, Otto Kruger, Action, Drama, Romance, War The Gay Falcon (1941) Irving Reis, George Sanders, Wendy Barrie, Allen Jenkins, Crime, Drama, Mystery, Romance Hndzan / Sour Grape (1974) Bagrat Oganessian, A. Isahakyan, Sos Sargsyan, H. Azizyan, Drama Muri shinjū: Nihon no natsu / Japanese Summer: Double Suicide (1967) Nagisa Ōshima, Keiko Sakurai, Kei Satō, Tetsuo Ashida, Crime, Drama Lady Oscar (1979) Jacques Demy, Catriona MacColl, Barry Stokes, Christine Belfrage, Drama, History, Romance The Trail of 98 (1928) Clarence Brown, Dolores del Río, Ralph Forbes, Karl Dane The Saint Meets the Tiger (1943) Paul L. Stein, Hugh Sinclair, Jean Gillie, Gordon McLeod, Crime, Drama, Mystery The Quiet Duel (1949) Akira Kurosawa, Toshirō Mifune, Takashi Shimura, Miki Sanjō, Drama Beyond the Sea (2004) Kevin Spacey, Kate Bosworth, John Goodman<endoftext> "This Moment Tests the Character of the Nation": Rep. Barbara Lee Rejects Anti-Refugee Efforts StoryNovember 18, 2015 Watch Full Show Watch Full ShowNext Story Media Options Democratic congresswoman from California and former chair of the Congressional Black Caucus. House Speaker Paul Ryan and Senate Majority Leader Mitch McConnell have called for a "pause" in the U.S. program accepting Syrian refugees. Meanwhile, governors of at least 27 U.S. states have said they will not accept Syrian refugees. We speak to California Democratic Rep. Barbara Lee. StoryJan 19, 2018As Shutdown Looms over Immigration, Trump's Rejection of Refugees Could Have Global Domino Effect AMY GOODMAN: As we turn now to Washington, D.C., as House Speaker Paul Ryan and Senate Majority Leader Mitch McConnell are calling for a pause in the U.S. program accepting Syrian refugees, I want to bring into the conversation Congressman Barbara Lee of California. Your response to the crackdown? Now, 27 governors are saying they will not accept Syrian refugees. In fact, your theory, Peter Bouckaert, around

Example 2 from domain: book

out. Well, they weren't getting away with it. There had to be a confrontation. I got to my feet and marched toward the cabin. It was showdown time. Cards on the table, Gagliano! What's your game, atheist? Let's have the truth! I kicked open the cabin door and got a surprise. Sitting there drinking wine was my father. "You raised in a barn?" "I closed the door carefully. "Where you been?" "Looking for you," I said. "Where you been?" "Right here." "All the time?" "All the time." "Didn't you hear me calling?" "When?" It was useless to ask any more questions. I sat down and he poured me a little wine. "Eat something," he said, pushing the bread and cheese across the table. "What's hemorrhoids?" He told me, and I had to push the food away. "You're too young for hemorrhoids." "Not me. That woman." "She's got her troubles." He rolled some wine in his cheeks, staring thoughtfully. His eyes seemed dipped in blood. "Your mother's a wonderful woman," he said. I just looked at him. "Finest woman in the world." He stood up, lurching, and drifted heavily to the door and outside. I went to the door. He sat on a log a few feet away, talking to himself. "An angel," he said. Though the twilight was still warm, I put some logs in the stove and stretched out on the couch. Leaning on an elbow I watched my father through the open door. He was like a statue, chin in both hands. It was very quiet but beyond the silence you heard the uproar out there, bullfrogs croaking, birds and crickets singing, bugs buzzing, and the trees sighing in the wind. The crackling fire splashed the ceiling with wild shadows and filled the cabin with warmth. # ELEVEN It felt like midnight when I wakened. Someone had slipped off my jeans and shoes and laid a blanket over me. Shafts of moonlight poured through the windows. The fire was a mound of ashes in the stove. The other two beds were not occupied. I was alone. I put on my shoes and jeans and went outside. The moon was gigantic. From the direction of the mine I heard Frank Gagliano's drunken gravel laughter, then the voice of Rhoda Pruitt, then a roar from my father. I told myself not to go up there, to stay in the cabin, to leave them alone, but I would not listen to myself, and the presence of evil coming from there drew me up the trail, running eagerly on tiptoe enchanted by the sense of evil. They did not hear me, nor even the thunder of my heart, nor did they even see me in the frenzy of their cleaving together, grunting and sucking and squirming in the naked heavy slithering of arms and legs, caught up like a ball of squirming white snakes, bodywhite under the moon, grinding on a blanket all knotted together with them, clawing, gasping, groaning. Then I saw my father's face. It was the face of the devil on the door. I turned and ran. I ran to the cabin. I was cold, shivering. I threw wood into the fire. I shuddered, wrapped in a blanket by the fire, teeth clack, clack. Then I was thirsty, drink anything, the wine! I drank and drank. Shivering, hungry, famished. But not their cheese, their hemorrhoid cheese, their bread. I found the box with the sandwiches my mother had made for me, and I ate, and it was good in my mouth, sweet and good, but I shivered all the same, the blanket around my shoulders, their fire burning in my face. Then I discovered the bottle she had placed there, wrapped in a cloth, a pint of holy water. She had written upon it, written: "Holy water. Use as needed." Now I knew it, now I would do it. I went up there, running, with the bottle of holy water, a fool with holy water, I knew it, I knew I was a fool, but I didn't care. They had to know I was coming. It was only fair to let them know, they were entitled to that. I yelled, "Holy water!" I ran, yelling, "Holy water!" "Holy water on its way!" "Here comes the holy water!"

Figure 13: Example inputs from SlimPajama-6B-Unif mixture showing the **selected tokens** by **ESLM-VaR-entropy (124M)** with $\alpha = 0.1$. [Note: These examples are drawn from public datasets (Soboleva et al., 2023) and may contain intense language, political references, or mature content. These excerpts are included solely for the purpose of analyzing model behavior.]

1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565

Example 1 from domain: cc

ens, Lon Chaney Jr. Shiinomi Gakuen / The Shiinomi School (1955) Hiroshi Shimizu, Kyōko Kagawa, Yukiko Shimazaki, Jūjirō Kichiji Uno, Drama Oliver Twist (1982) Clive Donner, George C. Scott, Tim Curry, Michael Hordern, Crime, Drama Cobra Woman (1944) Robert Siodmak, Maria Montez, Jon Hall, Sabu Imburnal (2008) Sherad Anthony Sanchez, Brian Monterola, Jeliet Mariveles-Ruca, Allen Lumanog Nemuri Kyoshiro 13: The Full Moon Swordsman (1969) Kazuo Mori, Hiroki Matsukata, Tomomi Satō, Sanae Nakahara, Action, Drama Shake Hands with the Devil (1959) Michael Anderson, James Cagney, Don Murray, Dana Wynter Sopyonje (1993) Kwon-taek Im, Myung-gon Kim, Jung-hae Oh, Kyu-chul Kim The American Soldier (1970) Rainer Werner Fassbinder, Karl Scheydtt, Elga Sorbas, Jan George, Drama Mike and Nicky (1976) Elaine May, Peter Falk, John Cassavetes, Ned Beatty Małgorzata, morska / The Little Mermaid (1976) Karel Kachyna, Miroslava Safránková, Radovan Lukavský, Petr Svoboda, Family, Drama, Romance Strategic Air Command (1955) Anthony Mann, James Stewart, June Allyson, Frank Lovejoy, Action, Drama, War De Passagem / Passing By (2003) Ricardo Elias, Lohan Brandão, Thiago de Mello, Wilma de Souza, Drama The Citadel (1938) King Vidor, Robert Donat, Rosalind Russell, Ralph Richardson Ladies in Retirement (1941) Charles Vidor, Ida Lupino, Louis Hayward, Evelyn Keyes Madame Satō (2002) Karim Amrani, Lázaro Ramos, Marcelia Cartaxo, Flavio Bauraquí, Biography, Crime, Drama Youth in Fury (1960) Masahiro Shinoda, Shin'ichirō Mikami, Shima Iwashita, Kayoko Honoo Night Plane from Chungking (1943) Ralph Murphy, Robert Preston, Ellen Drew, Otto Kruger, Action, Drama, Romance, War The Gay Falcon (1941) Irving Reis, George Sanders, Wendy Barrie, Allen Jenkins, Crime, Drama, Mystery, Romance Hndzan / Sour Grape (1974) Bagrat Oganessian, A. Isahakyan, Sos Sargsyan, H. Azizyan, Drama Muri shinjū: Nihon no natsu / Japanese Summer: Double Suicide (1967) Nagisa Ōshima, Keiko Sakurai, Kei Satō, Crime, Drama Lady Oscar (1979) Jacques Demy, Catriona MacColl, Barry Stokes, Christine Belfrage, Drama, History, Romance The Trail of 98 (1928) Clarence Brown, Dolores del Río, Ralph Forbes, Karl Dane The Saint Meets the Tiger (1943) Paul L. Stein, Hugh Sinclair, Jean Gillie, Gordon McLeod, Crime, Drama, Mystery The Quiet Duel (1949) Akira Kurosawa, Toshirō Mifune, Takashi Shimura, Miki Sanjō, Drama Beyond the Sea (2004) Kevin Spacey, Kate Bosworth, John Goodman<endoftext> "This Moment Tests the Character of the Nation": Rep. Barbara Lee Rejects Anti-Refugee Efforts StoryNovember 18, 2015 Watch Full Show Watch Full ShowNext Story Media Options Democratic congresswoman from California and former chair of the Congressional Black Caucus. House Speaker Paul Ryan and Senate Majority Leader Mitch McConnell have called for a "pause" in the U.S. program accepting Syrian refugees. Meanwhile, governors of at least 27 U.S. states have said they will not accept Syrian refugees. We speak to California Democratic Rep. Barbara Lee. StoryJan 19, 2018As Shutdown Looms over Immigration, Trump's Rejection of Refugees Could Have Global Domino Effect AMY GOODMAN: As we turn now to Washington, D.C., as House Speaker Paul Ryan and Senate Majority Leader Mitch McConnell are calling for a pause in the U.S. program accepting Syrian refugees, I want to bring into the conversation Congressman Barbara Lee of California. Your response to the crackdown? Now, 27 governors are saying they will not accept Syrian refugees. In fact, your theory, Peter Bouckaert, around

Example 2 from domain: book

out. Well, they weren't getting away with it. There had to be a confrontation. I got to my feet and marched toward the cabin. It was showdown time. Cards on the table, Gagliano! What's your game, atheist? Let's have the truth! I kicked open the cabin door and got a surprise. Sitting there drinking wine was my father. "You raised in a barn?" "I closed the door carefully. "Where you been?" "Looking for you," I said. "Where you been?" "Right here." "All the time?" "All the time." "Didn't you hear me calling?" "When?" It was useless to ask any more questions. I sat down and he poured me a little wine. "Eat something," he said, pushing the bread and cheese across the table. "What's hemorrhoids?" He told me, and I had to push the food away. "You're too young for hemorrhoids." "Not me. That woman." "She's got her troubles." He rolled some wine in his cheeks, staring thoughtfully. His eyes seemed dipped in blood. "Your mother's a wonderful woman," he said. I just looked at him. "Finest woman in the world." He stood up, lurching, and drifted heavily to the door and outside. I went to the door. He sat on a log a few feet away, talking to himself. "An angel," he said. Though the twilight was still warm, I put some logs in the stove and stretched out on the couch. Leaning on an elbow I watched my father through the open door. He was like a statue, chin in both hands. It was very quiet but beyond the silence you heard the uproar out there, bullfrogs croaking, birds and crickets singing, bugs buzzing, and the trees sighing in the wind. The crackling fire splashed the ceiling with wild shadows and filled the cabin with warmth. # ELEVEN It felt like midnight when I wakened. Someone had slipped off my jeans and shoes and laid a blanket over me. Shafts of moonlight poured through the windows. The fire was a mound of ashes in the stove. The other two beds were not occupied. I was alone. I put on my shoes and jeans and went outside. The moon was gigantic. From the direction of the mine I heard Frank Gagliano's drunken gravel laughter, then the voice of Rhoda Pruitt, then a roar from my father. I told myself not to go up there, to stay in the cabin, to leave them alone, but I would not listen to myself, and the presence of evil coming from there drew me up the trail, running eagerly on tiptoe enchanted by the sense of evil. They did not hear me, nor even the thunder of my heart, nor did they even see me in the frenzy of their cleaving together, grunting and sucking and squirming in the naked heavy slithering of arms and legs, caught up like a ball of squirming white snakes, bodywhite under the moon, grinding on a blanket all knotted together with them, clawing, gasping, groaning. Then I saw my father's face. It was the face of the devil on the door. I turned and ran. I ran to the cabin. I was cold, shivering. I threw wood into the fire. I shuddered, wrapped in a blanket by the fire, teeth clack, clack. Then I was thirsty, drink anything, the wine! I drank and drank. Shivering, hungry, famished. But not their cheese, their hemorrhoid cheese, their bread. I found the box with the sandwiches my mother had made for me, and I ate, and it was good in my mouth, sweet and good, but I shivered all the same, the blanket around my shoulders, their fire burning in my face. Then I discovered the bottle she had placed there, wrapped in a cloth, a pint of holy water. She had written upon it, written: "Holy water. Use as needed." Now I knew it, now I would do it. I went up there, running, with the bottle of holy water, a fool with holy water, I knew it, I knew I was a fool, but I didn't care. They had to know I was coming. It was only fair to let them know, they were entitled to that. I yelled, "Holy water!" I ran, yelling, "Holy water!" "Holy water on its way!" "Here comes the holy water!"

Figure 14: Example inputs from SlimPajama-6B-Unif mixture showing the **selected tokens** by ESLM-CVaR-loss (124M) with $\alpha = 0.1$. [Note: These examples are drawn from public datasets (Soboleva et al., 2023) and may contain intense language, political references, or mature content. These excerpts are included solely for the purpose of analyzing model behavior.]

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

Example 1 from domain: cc

hydrochloric acid in methylene chloride-water, followed by separation of the organic phase, drying, and storage in solution at 0-5 A.C. or below.6.7 Analysis of Reagent Purity: determination of positive chlorine can be carried out iodometrically.6.7 Handling, Storage, and Precautions: toxic and may explode, especially on heating or when concentrated. Dilute, cold solutions of NCl3 in various organic solvents are stable for several days.6 Store under inert atmosphere. Use only behind a safety shield in an efficient fume hood. Amination of Aromatics. The reaction of benzene and derivatives with NCl3 and Aluminum Chloride in organic solvents can be a useful preparation for meta-substituted amines. However, yields are only moderate, and mixtures of isomers are often obtained. Arenes include mono-7-9 and dialkylbenzenes, 9,10 halobenzenes,11 biphenyl, and naphthalene.1 with trichloroamine/AlCl3, the conversions of toluene to *m*-toluidine (eq 1) and of 1,3-dimethylbenzene to 3,5-dimethylaniline (eq 2) in moderate yields have been observed. An addition-elimination mechanism involving a chloroarenium intermediate has been proposed for the amination reactions (eq 3).10 Amination of halobenzenes and halotoluenes with trichloroamine/AlCl3 proceeds by two competing processes in moderate yields.11 For example, fluorobenzene gives predominantly *m*-fluoroaniline and *p*-chloroaniline (eq 4). It has been proposed that the former is produced by a substitution (addition-elimination) mechanism, while the latter is formed by a pathway involving nucleophilic displacement of halide in a chloroarenium cation by a nitrogen containing nucleophile (eq 5). Amination of biphenyl gives 3-aminobiphenyl (eq 6) and amination of naphthalene gives a mixture of 1- and 2-amino derivatives in low yields.11 Amination of Alkanes. The trichloroamine/AlCl3 system has also been used for the amination of monocyclic,12,13 bicyclic,13,14 and tricyclic,15 alkanes. C5-C8 cycloalkanes and their mono- and dimethyl derivatives are aminated in good yields.13 Methylcyclohexane,12,16 and methylcyclopentane,13 are converted to 1-amino-1-methylcycloalkanes on treatment with trichloroamine/AlCl3 (eq 7). Treatment of decalin and hydrindane with the trichloroamine/AlCl3 system affords *cis*-9-aminodecalin (eq 8) and *cis*-8-aminohydrindane, respectively, in good yields.13 The trichloroamine/AlCl3 amination route provides a simple one-stop method of obtaining aminoadamantanes in high yield (eq 9).3,15 Diamantane,1 can also be efficiently aminated in this fashion. When hydrocarbons which do not contain a tertiary hydrogen are subjected to reaction with NCl3/AlCl3, cationic rearrangements and fragmentations are observed.18 Amination of Alkyl-Substituted Aromatics. Various monoalkyl substituted arenes have been aminated on the alkyl side chain to form *t*-benzylamines in the system trichloroamine/AlCl3/*t*-butyl bromide (an efficient additive).19-21 *p*-Alkyl and *p*-haloisopropylbenzenes give the corresponding aminated products in high yields (eq 10). Tertiary Amines from Chlorides. When simple tertiary alkyl chlorides are exposed to NCl3 and AlCl3 in methylene chloride at -10 A.C., varying yields of the corresponding amines can be obtained. *t*-Pentyl chloride under these conditions provides *t*-pentylamine in 82% yield (eq 11), while *n*-octylamine is obtained in 30% yield from the corresponding chloride.2 Primary and secondary halides give isomeric amines resulting from skeletal rearrangement, as well as aziridines (eqs 12-14). Mechanistic details of this reaction have been reported.22 Vicinal Dichlorides from Alkenes. The reaction of NCl3 with a variety of mono- and disubstituted alkenes, both cyclic and acyclic, aff

Example 2 from domain: book

on the city's pulse, who is doing more, who is using their time to their maximum. And it can be anxiety-inducing, even when you've actively chosen not to do something. Our obsession with the next best thing and the activities of others is a blight of our consumer-driven society, and it is felt most keenly in cities. It is up to us to quiet the voice inside that asks why we always feel late to the party. The truth is that there will always be so much more happening in a city than you can ever spread yourself across, in person or even in awareness. We will always be surrounded by more different things we can possibly do. It is the difficulty of choice when faced with such a glut of opportunity that feels paralyzing. Making decisions is scary, and yet being confident in the decisions we make is the key to so much happiness and fulfillment in life. The word "decision" originates from the Latin *de* and *caedere*, meaning to cut off. *de* literally slaying your options. It's learning when and what to opt in and out of that really matters, though. Have confidence in your choices: make sure that they reflect who you are, and what you enjoy. Don't succumb to peer pressure, or let yourself become a wingman in someone else's experience of city life. And don't end up doing nothing because you couldn't decide what to do. Planning ahead is a useful strategy in combating FOMO. Set dates to do things, book tickets for shows, concerts and tables at restaurants. Invite others to join you. This is a simple way of ensuring you will have things in your diary to look forward to. Engineer your own fun, and take others along for the ride; we all love the friends who are organized enough to book tickets in bulk and bring everyone together. Just be mindful of scheduling sufficient space for spontaneity too. Feeling Safe It is easy to believe that cities are dangerous. We are exposed to news reports and statistics that can terrify timid souls into thinking every stranger on the sidewalk is a criminal in waiting. Clearly crime is more prevalent in cities than in rural areas, but much of this is due to the greater concentration of people. It is vital not to be intimidated into living under the covers for fear of what might happen. Feeling safe is largely a matter of common sense and vigilance: the more vigilant we are as city inhabitants, the stronger we become together as a deterrent. In general, it makes sense to keep to places where there are other people. And just as being alert to your own safety is common sense, be aware of the safety of others too. If you happen to witness an incident, act with courage but caution. We've all heard the parable of the woman who was attacked on the street in broad daylight in front of many people, but no one intervened because they assumed someone else would. Should any of us find ourselves the unfortunate victim in such a situation, a good way to attract help is to shout out to someone individually, referring to them by what they are wearing, thereby giving them ownership of the situation and responsibility to act. It is important to foster your own feelings of safety. Don't put yourself in situations where you feel unsafe. Make connections with people in your neighborhood. Be active and alert, not passive or invisible. As a city dweller you have responsibility to be part of a community that looks out for its other members. We are all in it together. Feeling Clean Cities are dirty. Even the more clinical, manicured Mitteleuropean or Japanese cities have cars, and pollution, and inhabitants with germs who don't wash their hands and occasionally sneeze on the back of your neck. For anyone even moderately concerned with hygiene, urban living is a constant battle the moment you leave the sanctity of your own home. The grime of pollution is tough. Blowing your nose after a journey on any underground transport system is not a pretty sight, and imagining what's in your lungs after a day out on foot or bike is enough to induce panic. Unfortunately, dirty air is a trade-off we have to accept in return for the many pleasures of city living. Short of buying a respiratory mask, there's little you can do to shield yourself from pollution. Things are looking up, though. Fewer cars on the streets mean less pollution in the air, and thankfully most cities are on board with the idea that this is the way forward. Most of us can take small comfort from knowing that things are better today than they were for our forebears, who could almost chew what they inhaled. When it comes to germs, we all fall foul of the inconsistent and selfish habits of humankind. Despite all the advice to wash our hands, catch a sneeze in a tissue and so forth, all it takes is one rogue individual not playing along to ruin it.

Example 3 from domain: book

was the dog. "He went up the beach with Jamie," Rick said. Harriet dropped down beside me. There had been conversation and laughter as we approached, but now there was silence as they froze us out. I saw that Rick and Denny were smoking pot. Harriet noticed it too. "Be careful," she cautioned. "The Sheriff patrols this beach all the time." They smiled like wise old men. "You want a joint, Dad?" Denny said. "No, thanks." "How about you, mother?" It was ridiculous and he knew better. I said, "Your Mother isn't a pot smoker, so stop being a wise guy." "This stuff is pure gold, Dad. Sure you won't try it?" "No, thanks." "It won't hurt you, man." "Listen. I smoked pot before you were born, back when you could buy a full Prince Albert tin of it for four bits." "Ah, the good old days!" he heeded. "Tell us about it." "There isn't much to tell. Pot is a mind expander for people with shriveled brains. You need it because you're a moron." "Thanks a lot." He crushed his cigarette into the sand, pulled off his shoes and socks, and trudged toward the water. Harriet looked after him with soft eyes. "That wasn't very nice," she said. I got up and went after him. He turned as I came splashing up to the creeping tide, then continued on down the beach. I caught up with him and put my arm around his shoulder. He slapped it away. "Leave me alone." "I'm sorry." "There you go, sorry again. You're always sorry after you insult somebody. You make sure you insult them first, and then you're sorry." "I try to be honest." "Honest! You're as devious as a snake, twisting and talking until you have it your way. You're the most two-faced bastard I ever saw." I was about to say I was sorry again, but I caught myself just in time. We splashed along for another fifty yards, our white feet in the thin embroidery of foam whisking across the dark sand, until we came to a skiff beached above the water line, seaweed and debris cluttered around it. He didn't want me with him, but I hung in there stubbornly as he leaned against the old boat and lit a cigarette. I didn't know what to say to him and he didn't know what to say to me. "Let's start back," I said. "I'm fed up with you, Dad." "Oh?" "I want you to stop calling me a moron. Ever since I can remember, all the way back to kindergarten, you've called me a moron. Why don't you cut it out?" "Okay." Maybe the pot did it. Maybe it was a break-through of his anger, the hot night and the curious circumstance that had brought us together at that moment. Maybe he had wanted to say it for years, but the right mood and moment had eluded him, but now he said it, and it sounded like a carefully prepared statement he had tucked away for a propitious time. "Dad, you're a lousy writer." That couldn't be my son Denny. It had to be the marijuana, just as it had been the wine with my father when I was twenty. He had bullied me for years and on Christmas Eve, hostile with wine, I had challenged him. We had fought it out in our front yard in North Sacramento, rolling in the dirt, kicking and gouging and cursing until the neighbors separated us. So it was Christmas Eve again. "I think Mother writes better than you do. I've read your novels. They're corny, sentimental cop-outs, and I'm not even talking about your screenplays." "The screenplays aren't much," I admitted. "Why did you ever become a writer, Dad? How the hell did you ever get published?" "Oh, shit, I'm not that bad! H. L. Mencken thought I was pretty good. He published me first." "You stink, Dad, you really do." "The Tyrant isn't a bad book. It got great reviews." "How many times did it sell?" "Not many, but it made a pretty good movie." "Have you seen it on TV lately?" I passed that one. "Anything else?"

Figure 15: Example inputs from SlimPajama-6B-Unif mixture showing the selected tokens by ES LM-VaR-entropy (774M) with $\alpha = 0.1$. [Note: These examples are drawn from public datasets (Soboleva et al., 2023) and may contain intense language, political references, or mature content. These excerpts are included solely for the purpose of analyzing model behavior.]

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

Example 1 from domain: cc

hydrochloric acid in methylene chloride-water, followed by separation of the organic phase, drying, and storage in solution at 0-5 A.C or below.6.7 Analysis of Reagent Purity: determination of positive chlorine can be carried out iodometrically. 6.7 Handling, Storage, and Precautions: toxic and may explode, especially on heating or when concentrated. Dilute, cold solutions of NCl3 in various organic solvents are stable for several days.6 Store under inert atmosphere. Use only behind a safety shield in an efficient fume hood. Amination of Aromatics. The reaction of benzene and derivatives with NCl3 and Aluminum Chloride in organic solvents can be a useful preparation for meta-substituted amines. However, yields are only moderate, and mixtures of isomers are often obtained. Arenes include mono-7-9 and dialkylbenzenes, 9,10 halobenzenes,11 biphenyl, and naphthalene.1 With trichloroamine/AlCl3, the conversions of toluene to *m*-toluidine (eq 1) and of 1,3-dimethylbenzene to 3,5-dimethylaniline (eq 2) in moderate yields have been observed. An addition-elimination mechanism involving a chloroareonium intermediate has been proposed for the amination reactions (eq 3).10 Amination of halobenzenes and halotoluenes with trichloroamine/AlCl3 proceeds by two competing processes in moderate yields.11 For example, fluorobenzene gives predominantly *m*-fluoroaniline and *p*-chloroaniline (eq 4). It has been proposed that the former is produced by a substitution (addition-elimination) mechanism, while the latter is formed by a pathway involving nucleophilic displacement of halide in a chloroareonium cation by a nitrogen containing nucleophile (eq 5). Amination of biphenyl gives 3-aminobiphenyl (eq 6) and amination of naphthalene gives a mixture of 1- and 2-amino derivatives in low yields.1 Amination of Alkanes. The trichloroamine/AlCl3 system has also been used for the amination of monocyclic,12,13 bicyclic,13,14 and tricyclic,15 alkanes. C5-C8 cycloalkanes and their mono- and dimethyl derivatives are aminated in good yields.13 Methylcyclohexane,12,16 and methylcyclopentane,13 are converted to 1-amino-1-methylcycloalkanes on treatment with trichloroamine/AlCl3 (eq 7). Treatment of decalin and hydriindane with the trichloroamine/AlCl3 system affords *cis*-9-aminodecalin (eq 8) and *cis*-8-aminohydriindane, respectively, in good yields.13 The trichloroamine/AlCl3 amination route provides a simple one-stop method of obtaining aminoadamantanes in high yield (eq 9).3,15 Diamantane,17 can also be efficiently aminated in this fashion. When hydrocarbons which do not contain a tertiary hydrogen are subjected to reaction with NCl3/AlCl3, cationic rearrangements and fragmentations are observed.18 Amination of Alkyl-Substituted Aromatics. Various monoalkyl substituted arenes have been aminated on the alkyl side chain to form *t*-benzylamines in the system trichloroamine/AlCl3/*t*-butyl bromide (an efficient additive).19-21 *p*-Alkyl and *p*-haloisopropylbenzenes give the corresponding aminated products in high yields (eq 10). Tertiary Amines from Chlorides. When simple tertiary alkyl chlorides are exposed to NCl3 and AlCl3 in methylene chloride at -10 A.C, varying yields of the corresponding amines can be obtained. *t*-Pentyl chloride under these conditions provides *t*-pentylamine in 82% yield (eq 11), while *n*-octylamine is obtained in 30% yield from the corresponding chloride.2 Primary and secondary halides give isomeric amines resulting from skeletal rearrangement, as well as aziridines (eqs 12-14). Mechanistic details of this reaction have been reported.22 Vicinal Dichlorides from Alkenes. The reaction of NCl3 with a variety of mono- and disubstituted alkenes, both cyclic and acyclic, aff

Example 2 from domain: book

on the city's pulse, who is doing more, who is using their time to their maximum. And it can be anxiety-inducing, even when you've actively chosen not to do something. Our obsession with the next best thing and the activities of others is a blight of our consumer-driven society, and it is felt most keenly in cities. It is up to us to quiet the voice inside that asks why we always feel late to the party. The truth is that there will always be so much more happening in a city than you can ever spread yourself across, in person or even in awareness. We will always be surrounded by more different things we can possibly do. It is the difficulty of choice when faced with such a glut of opportunity that feels paralyzing. Making decisions is scary, and yet being confident in the decisions we make is the key to so much happiness and fulfillment in life. The word "decision" originates from the Latin *de* and *caedere*, meaning to cut off. Literally slaying your options. It's learning when and what to opt in and out of that really matters, though. Have confidence in your choices: make sure that they reflect who you are, and what you enjoy. Don't succumb to peer pressure, or let yourself become a wingman in someone else's experience of city life. And don't end up doing nothing because you couldn't decide what to do. Planning ahead is a useful strategy in combating FOMO. Set dates to do things, book tickets for shows, concerts and tables at restaurants. Invite others to join you. This is a simple way of ensuring you will have things in your diary to look forward to. Engineer your own fun, and take others along for the ride; we all love the friends who are organized enough to book tickets in bulk and bring everyone together. Just be mindful of scheduling sufficient space for spontaneity too. Feeling Safe It is easy to believe that cities are dangerous. We are exposed to news reports and statistics that can terrify timid souls into thinking every stranger on the sidewalk is a criminal in waiting. Clearly crime is more prevalent in cities than in rural areas, but much of this is due to the greater concentration of people. It is vital not to be intimidated into living under the covers for fear of what might happen. Feeling safe is largely a matter of common sense and vigilance: the more vigilant we are as city inhabitants, the stronger we become together as a deterrent. In general, it makes sense to keep to places where there are other people. And just as being alert to your own safety is common sense, be aware of the safety of others too. If you happen to witness an incident, act with courage but caution. We've all heard the parable of the woman who was attacked on the street in broad daylight in front of many people, but no one intervened because they assumed someone else would. Should any of us find ourselves the unfortunate victim in such a situation, a good way to attract help is to shout out to someone individually, referring to them by what they are wearing, thereby giving them ownership of the situation and responsibility to act. It is important to foster your own feelings of safety. Don't put yourself in situations where you feel unsafe. Make connections with people in your neighborhood. Be active and alert, not passive or invisible. As a city dweller you have responsibility to be part of a community that looks out for its other members. We are all in it together. Feeling Clean Cities are dirty. Even the more clinical, more European or Japanese cities have cars, and pollution, and inhabitants with germs who don't wash their hands and occasionally sneeze on the back of your neck. For anyone even moderately concerned with hygiene, urban living is a constant battle the moment you leave the sanctity of your own home. The grime of pollution is tough. Blowing your nose after a journey on any underground transport system is not a pretty sight, and imagining what's in your lungs after a day out on foot or bike is enough to induce panic. Unfortunately, dirty air is a trade-off we have to accept in return for the many pleasures of city living. Short of buying a respiratory mask, there's little you can do to shield yourself from pollution. Things are looking up, though. Fewer cars on the streets mean less pollution in the air, and thankfully most cities are on board with the idea that this is the way forward. Most of us can take small comfort from knowing that things are better today than they were for our forebears, who could almost chew what they inhaled. When it comes to germs, we all fall foul of the inconsistent and selfish habits of humankind. Despite all the advice to wash our hands, catch a sneeze in a tissue and so forth, all it takes is one rogue individual not playing along to ruin it.

Example 3 from domain: book

was the dog. "He went up the beach with Jamie," Rick said. Harriet dropped down beside me. There had been conversation and laughter as we approached, but now there was silence as they froze us out. I saw that Rick and Denny were smoking pot. Harriet noticed it too. "Be careful," she cautioned. "The Sheriff patrols this beach all the time." They smiled like wise old men. "You want a joint, Dad?" Denny said. "No, thanks." "How about you, mother?" It was ridiculous and he knew better. I said, "Your Mother isn't a pot smoker, so stop being a wise guy." "This stuff is pure gold, Dad. Sure you won't try it?" "No, thanks." "It won't hurt you, man." "Listen. I smoked pot before you were born, back when you could buy a full Prince Albert tin of it for four bits." "Ah, the good old days!" he needed. "Tell us about it." "There isn't much to tell. Pot is a mind expander for people with shriveled brains. You need it because you're a moron." "Thanks a lot." He crushed his cigarette into the sand, pulled off his shoes and socks, and trudged toward the water. Harriet looked after him with soft eyes. "That wasn't very nice," she said. I got up and went after him. He turned as I came splashing up to the creeping tide, then continued on down the beach. I caught up with him and put my arm around his shoulder. He slapped it away. "Leave me alone." "I'm sorry." "There you go, sorry again. You're always sorry after you insult somebody. You make sure you insult them first, and then you're sorry." "I try to be honest." "Honest! You're as devious as a snake, twisting and talking until you have it your way. You're the most two-faced bastard I ever saw." I was about to say I was sorry again, but I caught myself just in time. We splashed along for another fifty yards, our white feet in the thin embroidery of foam whisking across the dark sand, until we came to a skiff beached above the water line, seaweed and debris cluttered around it. He didn't want me with him, but I hung in there stubbornly as he leaned against the old boat and lit a cigarette. I didn't know what to say to him and he didn't know what to say to me. "Let's start back," I said. "I'm fed up with you, Dad." "Oh?" "I want you to stop calling me a moron. Ever since I can remember, all the way back to kindergarten, you've called me a moron. Why don't you cut it out?" "Okay." Maybe the pot did it. Maybe it was a break-through of his anger, the hot night and the curious circumstance that had brought us together at that moment. Maybe he had wanted to say it for years, but the right mood and moment had eluded him, but now he said it, and it sounded like a carefully prepared statement he had tucked away for a propitious time. "Dad, you're a lousy writer." That couldn't be my son Denny. It had to be the marijuana, just as it had been the wine with my father when I was twenty. He had bullied me for years and on Christmas Eve, hostile with wine, I had challenged him. We had fought it out in our front yard in North Sacramento, rolling in the grass, kicking and gouging and cursing until the neighbors separated us. So it was Christmas Eve again. "I think Mother writes better than you do. I've read your novels. They're corny, sentimental cop-outs, and I'm not even talking about your screenplays." "The screenplays aren't much." I admitted. "Why did you ever become a writer, Dad? How the hell did you ever get published?" "Oh, shit, I'm not that bad! H. L. Mencken thought I was pretty good. He published me first." "You stink, Dad, you really do." "The Tyrant isn't a bad book. It got great reviews." "How many times did it sell?" "Not many, but it made a pretty good movie." "Have you seen it on TV lately?" I passed that one. "Anything else?"

Figure 16: Example inputs from SlimPajama-6B-Unif mixture showing the selected tokens by ESLM-CVaR-loss (774M) with $\alpha = 0.1$. [Note: These examples are drawn from public datasets (Soboleva et al., 2023) and may contain intense language, political references, or mature content. These excerpts are included solely for the purpose of analyzing model behavior.]