# TS<sup>2</sup>: Training with Sparsemax+, Testing with Softmax for Accurate and Diverse LLM Fine-Tuning

# **Anonymous authors**

Paper under double-blind review

# **ABSTRACT**

Large Language Models typically rely on Supervised Fine-Tuning (SFT) with Cross-Entropy (CE) loss to specialize in downstream tasks. However, CE forces the distribution toward one-hot targets and ignores alternative continuations, thereby limiting output diversity—a key drawback for generative applications that rely on sampling-based exploration. In this paper, we propose "Training with Sparsemax+, Testing with Softmax (TS<sup>2</sup>)". Intuitively, sparsemax and its tailored loss mask the gradients of probabilities outside the support set, leaving excessive probability mass on irrelevant tail classes when evaluating with softmax. To address this issue, we propose an improved variant, Sparsemax+, for training, which augments the sparsemax loss with a suppression term that penalizes the out-ofsupport probabilities. At testing, we decode with softmax, yielding calibrated, non-degenerate probabilities where plausible near-ties survive. We fine-tuned Llama-3.1-8B and Qwen-2.5-7B with TS<sup>2</sup>, achieving consistent improvements in accuracy and output diversity across chat, code, and open-domain benchmarks. Together, these results demonstrate that TS<sup>2</sup> provides a practical, drop-in solution for fine-tuning LLMs that are both more accurate and more creative.

#### 1 Introduction

Supervised fine-tuning (SFT) is one of the major steps in the Large Language Models (LLMs) post-training stage: with a small amount of high-quality annotated data, it teaches models to organize language better and produce instruction-following responses. The default loss function is cross-entropy loss, mainly because it coincides with maximum likelihood and is a strictly proper scoring rule, so minimizing it recovers the data generating conditional under well-specification (Gneiting & Raftery, 2007). However, the same geometry drives the probability mass toward the one-hot target and away from plausible alternatives, yielding overconfident posteriors and reduced useful diversity. A large body of work seeks to counteract this overconfidence and recover useful diversity. One branch changes only the decoding, e.g., nucleus sampling and best-of-N, leaving training dynamics and calibration untouched (Holtzman et al., 2020). Another branch alters the training signal itself. The recent GEM framework reframes SFT as reverse-KL minimization with an entropy regularizer, improving variety and mitigating overfitting (Li et al., 2025). These approaches highlight a fundamental issue: promoting diversity can conflict with keeping probabilities calibrated and tails disciplined.

We argue that the field lacks a precise operational notion of useful "diversity" for instruction following. In many tasks, we do not want to "spread probability" indiscriminately over the entire vocabulary. Instead, we want probability mass concentrated among a handful of semantically plausible next tokens, those with a real chance of leading to a high quality continuation, while aggressively deflating the long tail of obviously incorrect tokens toward (near) zero. The right diversity is within the plausible set, not across the whole simplex. The forward KL  $\mathrm{KL}(p \parallel q)$  is mean-seeking, incentivizing probability wherever the data has support; the reverse  $\mathrm{KL}\,\mathrm{KL}(q \parallel p)$  is mode-seeking, concentrating mass on promising regions (Minka, 2005). This lens helps explain why CE with entropy maximization (a forward-KL-flavored objective under softmax) can inflate low-probability tokens, while reverse-KL flavored objectives like GEM avoid gratuitous tail mass. Yet even reverse-KL does not guarantee that clearly implausible tokens go to zero.

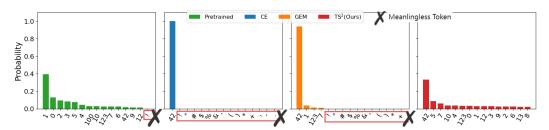


Figure 1: Token Distribution for single digit generation (detailed in Appendix C.4).

Our approach takes a geometric route by decoupling the mapping from logits to probabilities between training and testing. Specifically, we train with sparsemax and optimize a modified Fenchel-Young loss tailored to this mapping (Martins & Astudillo, 2016; Blondel et al., 2019), while at inference we revert to softmax, which restores calibrated and smooth probabilities on the same logits. Our tailored loss contains a tail penalty that drives non-support tokens to zero while ensuring the gold token is never penalized, even if it lies outside the instantaneous sparsemax support. Notably, CE with softmax collapses diversity: all non-gold logits, even plausible ones, are pushed toward zero. In contrast, sparsemax maintains a sparse support set by zeroing gradients of non-support tokens, preserving plausible candidates. However, if sparsemax were also used at inference, a converged model would still produce one-hot outputs—similar to CE with softmax decoding—thus limiting diversity (Martins & Astudillo, 2016; Blondel et al., 2019). Figure 1 shows that the pretrained model naturally exhibits diversity, but such diversity is lost during SFT: CE training drives the distribution into a one-hot solution, while GEM manages to retain only a few candidates, with most remaining probability mass assigned to irrelevant symbols. In contrast, our method delivers stronger and more stable diversity: the candidate set is both meaningful and varied, striking a balance between plausibility and coverage.

This decoupled recipe, Train with Sparsemax+, Test with Softmax (TS<sup>2</sup>), has two key effects. During training, sparse gradients act as a principled early-stopping mechanism by avoiding wasted updates on already separated tail candidates. At inference, reverting to softmax restores smooth, calibrated probabilities so that plausible near-ties survive and sampling can explore them without aggressive temperature tuning. By construction, our method achieves local diversity among plausible tokens while assigning near-zero confidence to implausible ones. We position  $TS^2$  among complementary strategies. Inference-only methods (e.g., nucleus, top-k, best-of-N) improve sample variety but leave training untouched; our approach reshapes training dynamics while remaining fully compatible with such decoders (Holtzman et al., 2020). Entropy targeting methods (e.g., GEM) promote spread but do not enforce exact zeros on implausible tokens; our penalty term supplies this "hard" suppression, while sparsemax ensures spread occurs where it matters (Li et al., 2025). Finally, because  $TS^2$  decouples mappings rather than altering model architecture, it integrates seamlessly into existing SFT pipelines. The contribution of this paper are summarized in the following:

- We frame the problem as achieving Tail-Suppressed Plausible Diversity (TSPD) and propose TS<sup>2</sup>, which decouples training and inference by using a Sparsemax+ loss with tail penalty for training and standard softmax for decoding.
- We provide a theoretical analysis showing how TS<sup>2</sup> avoids the distributional collapse common to CE training via a gradient-masking mechanism, thereby preserving diversity at inference.
- We demonstrate in practice that our TS<sup>2</sup> significantly improves winrates, sample efficiency in code generation, and output diversity across multiple benchmarks compared to existing methods.

# 2 DISTRIBUTION COLLAPSE AND OUR INSIGHT

Recent studies have observed an "alignment tax" in large language models (LLMs): while supervised fine-tuning (SFT) improves faithfulness and task adherence (Brown et al., 2020), it often comes at the cost of reduced output diversity and partial forgetting of pre-trained knowledge (O'Mahony et al., 2024; Kim et al., 2025). Pre-trained LLMs naturally exhibit a broad generative repertoire, producing multiple semantically valid outputs for the same prompt (Wang et al., 2025). However, after SFT, models tend to respond with highly deterministic and homogeneous output (Li et al., 2025), weakening their utility in downstream applications such as planning (Song et al., 2023), writing (Lee

111

112

113

114

115

116 117

118 119

120

121

122 123

124

125 126

127

128

133

134 135

136 137

138

139

140

141

142

143

144

145

146

147 148 149

150 151

152 153

154

155

156

157

158 159 160

161

et al., 2022), or code generation (Liu et al., 2023), all of which fundamentally rely on the ability to explore diverse candidate responses.

A central obstacle in supervised fine-tuning is that cross-entropy (CE), driving the predictive distribution towards a one-hot distribution, causing all probability mass to collapse onto the gold token. This distribution collapse ensures convergence, but comes at a severe cost: the model suppresses all alternatives to nearly zero, leading to deterministic outputs, both in the choice of tokens and in the semantic content of the whole responses. The mechanism destroys diversity, erasing helpful variations preserved in the pre-trained distribution and thereby yielding monotonous generations.

#### 2.1 Our guiding insight: Tail-Suppressed Plausible Diversity

In generative modeling, output diversity is essential. A target distribution should retain a compact set of plausible candidates with non-negligible probability while suppressing irrelevant long-tail tokens toward zero. We formalize this as Tail-Suppressed Plausible Diversity (TSPD), which remedies the distribution collapse commonly observed in existing SFT.

**Notation.** We consider prompt–response pairs  $(x,y) \in \mathcal{D}$  from a supervised dataset  $\mathcal{D}$ . Let  $f_{\theta}$ denote a pre-trained LLM parameterized by  $\theta$ . For a prompt x, let  $z = f_{\theta}(x) \in \mathbb{R}^{K}$  denote the corresponding logit vector<sup>1</sup>. We define the probability simplex as  $\Delta^{K-1} = \{ p \in \mathbb{R}^{K} \mid p_{i} \geq 1 \}$  $0, \sum_{i=1}^{K} p_i = 1$ , where p = g(z) denotes a probability distribution obtained from the logits z via a probability mapping function  $g(\cdot)$ .

**Definition 1** (Tail-Suppressed Plausible Diversity  $(m, \varepsilon_{\text{head}}, \varepsilon_{\text{tail}})$ ). Given a prompt-response pair (x, y), let  $\mathbf{p} = g(f_{\theta}(x)) \in \Delta^{K-1}$  be a distribution over a vocabulary  $\mathcal{V}$ . Fix an integer  $m \geq 2$  and thresholds  $0 < \varepsilon_{\text{head}} \leq \frac{1}{m}$  and  $0 \leq \varepsilon_{\text{tail}} \leq 1 - m \varepsilon_{\text{head}}$ . Let  $\text{Top}_m(\mathbf{p})$  denote the indices of the m largest coordinates of  $\mathbf{p}$ . If  $y \in \text{Top}_m(\mathbf{p})$ , let  $\mathcal{S} := \text{Top}_m(\mathbf{p})$ ; otherwise, let  $\mathcal{S} := \text{Top}_{m-1}(\mathbf{p}) \cup \{y\}$ . We say that p satisfies TSPD of order m if

(Head Preservation) 
$$\min_{j \in S} p_j \ge \varepsilon_{\text{head}},$$
 (1a)

which ensures that candidates in S retain non-negligible probability, whereas tokens outside S receive essentially zero probability, thereby preserving uncertainty and transferable knowledge at inference. If one chooses  $\varepsilon_{\rm head} = 1/m$  exactly, then the strict requirement  $\varepsilon_{\rm tail} \geq 0$  forces  $m\varepsilon_{\rm head} = 1$ and  $p_j = 0 \ \forall j \notin S$ ; therefore, in practice one can take  $\varepsilon_{\rm head} < 1/m$  and relax  $\varepsilon_{\rm tail} > 0$ .

**Corollary 1.** If Definition 1 holds and  $\varepsilon_{\text{tail}} < \varepsilon_{\text{head}}$ , then  $\max_{j \notin S} p_j \le \varepsilon_{\text{tail}} < \varepsilon_{\text{head}} \le \min_{i \in S} p_i$ , so each plausible sample has strictly higher probability than any tail sample.

**Corollary 2.** If all probability mass collapses onto the ground-truth token, i.e.,  $p_{ij} = 1$  and  $p_{ij'} = 1$  $0 \ \forall y' \neq y$ , then **p** fails to qualify the TSPD  $(m(\geq 2), \varepsilon_{\text{head}}, \varepsilon_{\text{tail}})$ .

In the next section, we motivate our method that operationalizes this principle, directly countering the diversity-reducing bias of CE loss while retaining the benefits of supervised fine-tuning.

#### 3 ACHIEVING TAIL-SUPPRESSED PLAUSIBLE DIVERSITY

A natural way to realize TSPD in Equation (1) is to exploit the sparsity of the sparsemax mapping sparsemax(z), which projects logits z onto the probability simplex  $\Delta^{K-1}$ , yielding exact zeros outside a data-dependent support. Formally,

$$p^{\mathsf{sp}}(\boldsymbol{z}) = [\boldsymbol{z} - \tau(\boldsymbol{z})\mathbf{1}]_{+} := \operatorname{sparsemax}(\boldsymbol{z}),$$

where  $[z]_+ := \max\{z,0\}$  is applied elementwise, and the threshold is defined as  $\tau(z) =$  $\frac{\sum_{j \in S^{\text{sp}}(\boldsymbol{z})} z_j - 1}{|S^{\text{sp}}(\boldsymbol{z})|}$ , with  $S^{\text{sp}}(\boldsymbol{z}) = \{j : z_j - \tau(\boldsymbol{z}) > 0\}$  denoting the support set. In effect, sparsemax automatically identifies a compact support set of plausible candidates  $S^{\sf sp}(z)$  and prunes away the long tail. Compared to softmax probability mapping  $p^{\sf sf}(z) = \frac{\exp(z)}{\sum_{i=1}^K \exp(z_i)} := \operatorname{softmax}(z)$ , its Jacobian is sparse. See Lemma 3 for more details.

 $<sup>^{1}</sup>y$  and x can be sequential, where an auto-regressive formulation is used.

**Lemma 3** (Gradients vanish outside the sparsemax support). Let  $p = \operatorname{sparsemax}(z)$  and  $S^{\operatorname{sp}}(z)$  be its support. Consider the sparsemax loss  $\mathcal{L}_{\operatorname{sp}}(z,y)$  with target y. If  $y \in S^{\operatorname{sp}}(z)$ , then  $\forall i \notin S^{\operatorname{sp}}(z)$ ,  $\frac{\partial \mathcal{L}_{\operatorname{sp}}(z,y)}{\partial z_i} = 0$ .

While sparsemax provides margin-induced sparsity, it nonetheless tends to collapse into a nearly one-hot distribution once the leading logit surpasses the margin threshold. Such collapse inevitably reduces sampling diversity, making sparsemax undesirable for inference.

This motivates us to instead carry out decoding with softmax. Under this choice, the gradient-vanishing property established in Lemma 3 remains advantageous during training: by nullifying gradients outside the active support whenever the target is included, it mitigates the cross-entropy—style erosion of plausible near-optimal alternatives, thereby inducing an implicit early-stopping effect.

**Theorem 4** (Sparsemax expands pairwise gaps faster than softmax). Let  $z \in \mathbb{R}^K$ ,  $p^{sf} = \operatorname{softmax}(z)$ , and  $p^{sp} = \operatorname{sparsemax}(z)$ . For any indices  $i \neq j$ , let  $u := z_i - z_j$  and we have

$$\frac{\partial}{\partial u}(p_i^{\rm sp}-p_j^{\rm sp})=1 \quad \forall \, i,j \in \mathcal{S}^{\rm sp}, \qquad \text{sparsemax}$$

$$\frac{\partial}{\partial u}\left(p_i^{\rm sf}-p_j^{\rm sf}\right)<1, \qquad \text{softmax}$$

Given the same logits, Theorem 4 shows that sparsemax linearly preserves pairwise probability gaps within its active support and collapses to a one-hot prediction once a finite margin is attained, whereas softmax strictly contracts such gaps. Consequently, sparsemax induces sharp discrimination and faster label collapse during training, while applying softmax to the same logits at inference preserves non-degenerate mass on plausible candidates—maintaining output diversity that is desirable for generative tasks.

**Corollary 5** (Softmax remains TSPD-valid when sparsemax is one-hot). Let  $z \in \mathbb{R}^K$  with  $y = \arg\max_j z_j$ , and  $\delta_j := z_y - z_j$ . Assume sparsemax is one-hot at y, i.e.,  $\delta_{\min} := \min_{j \neq y} \delta_j \geq \gamma > 0$  (e.g.,  $\gamma = 1$ ), and the top-m head is bounded:  $\delta_{(k)} := z_c - z_{(k)} \leq B \ \forall k = 2, \ldots, m$ . Set  $A_m = m + (K - m)e^{-\gamma}$ . Then for  $p^{\text{sf}} = \operatorname{softmax}(z)$  we have

$$p_y^{\mathrm{sf}} \geq \frac{1}{A_m}, \quad p_{(k)}^{\mathrm{sf}} \geq \frac{e^{-B}}{A_m} \; (\forall k=2,\ldots,m), \quad \sum_{k>m} p_{(k)}^{\mathrm{sf}} \leq \frac{(K-m)e^{-\gamma}}{A_m}.$$

Consequently,  $p^{\text{sf}}$  satisfies TSPD of order m with any thresholds  $0 < \varepsilon_{\text{head}} \le \frac{e^{-B}}{A_m}$ ,  $\frac{(K-m)e^{-\gamma}}{A_m} \le \varepsilon_{\text{tail}} \le 1 - m \varepsilon_{\text{head}}$ .

**Remark 1.** Without the head bound  $\delta_{(k)} \leq B$  ( $\forall k \leq m$ ),  $p_{(m)}^{sf}$  can be made arbitrarily small even when  $\delta_{\min} \geq 1$ , so only a vanishingly small head floor  $\varepsilon_{\text{head}}$  can be guaranteed for general m.

According to Corollary 5, the cumulated tail mass of softmax outside the top-m satisfies  $\sum_{k>m} p_{(k)}^{\rm sf} \leq \frac{(K-m)e^{-\gamma}}{A_m}$ . This upper bound is strictly increasing in K (for fixed  $m,\gamma$ ) and approaches 1 as  $K\to\infty$ . Thus, with large vocabularies, the admissible tail under softmax at inference becomes nearly 1, indicating that sparsemax training has not sufficiently penalized tail tokens, contradicting the goal of suppressing irrelevant tail mass.

To address these issues, we propose a fine-tuning strategy of **Training with Sparsemax+**, **Testing with Softmax**. Sparsemax+ builds on Sparsemax, inheriting margin-induced sparsity to introduce gradient masking during training, thereby implicitly enforcing an early-stopping effect once the top-1 candidate is clearly separated. It further incorporates a lightweight *Tail-suppressing Loss* to explicitly penalize residual probability on tail tokens, ensuring that tail mass is sharply suppressed. At inference, we revert to softmax over the same logits, which restores smooth, calibrated probabilities across the plausible candidates within the support set, while keeping the irrelevant tail mass negligible due to the additional suppressing effect. In this way, the model learns to *separate and prune* the logits during training, yet *preserve and diversify* the output distribution during inference, achieving the desired support-aware diversity.

# 4 TS<sup>2</sup>: TRAINING WITH SPARSEMAX+, TESTING WITH SOFTMAX

In the following, we present supervised fine-tuning based on the Fenchel-Young loss, which encompasses both the softmax and sparsemax mappings. It then motivates our Sparsemax+ loss.

#### 4.1 DIFFERENT PREDICTION MAPPINGS WITH THE UNIFIED FENCHEL-YOUNG LOSS

For any strictly convex regularization function  $\Omega: \Delta^{K-1} \to \mathbb{R}$ , the corresponding regularized prediction function is  $p_*(z) = \arg\max_{p \in \Delta^{K-1}} \langle p, z \rangle - \Omega(p)$ . The associated Fenchel-Young loss can be represented as

$$L_{\Omega}(z;y) = \Omega(e_y) - \Omega(p_*) + \langle z, p_* - e_y \rangle, \tag{3}$$

where y is the gold label and  $e_y$  is the corresponding one-hot vector. Different choices of  $\Omega$  yield different prediction mappings and losses.

Softmax Softmax corresponds to using the negative Shannon entropy as regularizer  $\Omega(\boldsymbol{p}) = \sum_{i=1}^K p_i \log p_i$ , which gives  $\boldsymbol{p}_*(\boldsymbol{z}) = \frac{\exp(\boldsymbol{z})}{\sum_{i=1}^K \exp(z_i)} := \operatorname{softmax}(\boldsymbol{z})$ . The Fenchel-Young loss reduces to the standard CE loss  $L_{\operatorname{softmax}}(\boldsymbol{z};y) = \log \sum_{i=1}^K \exp(z_i) - z_y = -\log \frac{\exp(z_y)}{\sum_{i=1}^K \exp(z_i)}$ .

**Sparsemax** Sparsemax corresponds to using the negative Gini entropy as regularizer  $\Omega(\boldsymbol{p})=\frac{1}{2}\sum_{i=1}^K p_i(1-p_i)$ , which gives  $\boldsymbol{p}_*(\boldsymbol{z})=[\boldsymbol{z}-\tau(\boldsymbol{z})\mathbf{1}]_+:=\operatorname{sparsemax}(\boldsymbol{z})$ . The corresponding Fenchel-Young loss, called the sparsemax loss, is  $L_{\operatorname{sparsemax}}(\boldsymbol{z};y)=-z_y+\frac{1}{2}\sum_{j\in S^{\operatorname{sp}}(\boldsymbol{z})}\left(z_j^2-\tau^2(\boldsymbol{z})\right)+\frac{1}{2}$ .

In conclusion, when training with sparsemax but performing inference with softmax, although  $\operatorname{softmax}(z)$  does not yield a one-hot output like  $\operatorname{sparsemax}(z)$ , it still assigns the highest probability to the correct class. Importantly, it naturally enables early stopping and preserves distributional diversity across all classes, which is consistent with the goal of diversifying plausible candidates.

Given prompt-response pairs (x,y) from a supervised dataset, let  $z \in \mathbb{R}^K$  be a logit vector and p be probability mapping either via softmax or sparsemax. If the gradient of the sparsemax loss vanishes, i.e.,  $\nabla_z \mathcal{L}_{\text{sparsemax}}(z;y) = 0$ , then it follows that  $\text{sparsemax}(z)_y = 1$ . For any index  $\forall j \neq y$ ,  $\text{sparsemax}(z)_j = 0$ , it holds that  $\text{softmax}(z)_j > 0$ . That is, softmax assigns non-zero probability to all entries, including those which sparsemax maps to zero. According to Corollary 5, the cumulated tail mass of softmax outside the top-m satisfies  $\sum_{k>m} p_{(k)}^{\text{sf}} \leq \frac{(K-m)e^{-\gamma}}{A_m}$ . With large vocabularies, the admissible tail under softmax at inference becomes nearly 1. This behavior is undesirable, as assigning non-negligible probabilities to clearly incorrect classes may lead the model to produce semantically meaningless outputs.

**Sparsemax**+ To address this issue, we introduce a lightweight tail-suppressing loss that explicitly suppresses probabilities assigned to the non-plausible candidates. Given logits  $z \in \mathbb{R}^K$ , let  $p^{sf} = \operatorname{softmax}(z) \in \Delta^{K-1}$ . The tail-suppressing loss is defined as

$$\mathcal{L}_{\sup}(\boldsymbol{p}; y) = -\log(1 - \sum_{i \notin \mathcal{S}} p_i^{\mathsf{sf}}),$$

where S is defined in Definition 1. This penalty drives the probabilities of tail tokens toward zero, thereby avoiding residual mass on clearly implausible candidates.

**Remark 2.** The tail suppressing loss can be interpreted as a direct generalization of the standard softmax CE to the group-label setting. Specifically, given logits z and softmax distribution  $p^{sf} = \operatorname{softmax}(z)$ , the suppressing term can be written as

$$L_{\sup}(\boldsymbol{z}) = -\log(1 - \sum_{i \notin S} p_i^{\mathsf{sf}}) = -\log \sum_{i \in S} p_i^{\mathsf{sf}},$$

which is exactly the softmax cross-entropy with the target label being the merged "super-class" S. In the special case where  $S = \{y\}$  is a singleton, this reduces to the usual CE loss  $-\log p_y^{\sf sf}$ . Thus, the suppressing loss can be viewed as encouraging the softmax probability mass to concentrate on a set of plausible candidates while retaining the probabilistic interpretation of cross-entropy.

271

272

273

274

275

276

277

278

279

281

282

283

284285286

287

288

289 290

291

292

293

295

296

297

298299300

301

302

303

304

305

306

307

308

310

311

312

313

314

315

316

317

318 319

320

321

322

323

# Algorithm 1 TS<sup>2</sup>: Training with Sparsemax+, Testing with Softmax

```
Input: pre-trained model f_{\theta}; training dataset \mathcal{D}_{tr} = \{(x,y)\}; test dataset \mathcal{D}_{te} = \{x\}.
     Hyperparameters: epochs T; batch size B; learning rate \eta > 0; suppression weight \alpha > 0.
 1: for t = 1 to T do
                                                                                                                      ▶ Training Phase
          for mini-batch \{(\boldsymbol{x}_b, \boldsymbol{y}_b)\}_{b=1}^B \subset \mathcal{D}_{tr} do
                Compute logits z_b \leftarrow f_{\theta}(x_b), \forall b = 1, 2, \dots, B
 3:
                Compute loss L_b \leftarrow L_{\text{spm}} + (\boldsymbol{z}_b; \boldsymbol{y}_b), \forall b = 1, 2, \dots, B
 4:

    ▷ Sparsemax + loss

 5:
           Update \theta \leftarrow \theta - \eta \nabla_{\theta} \frac{1}{B} \sum_{b=1}^{B} L_b
 6:
 7: end for
 8: for test input x \in \mathcal{D}_{te} do

    ▶ Testing Phase

           Compute logits z \leftarrow f_{\theta}(x)
 9:
10:
           Predict probability p \leftarrow \operatorname{softmax}(z)
           Evaluation on p
                                                                                                                   11:
12: end for
```

Combining sparsemax with the tail-suppressing loss yields our proposed *Sparsemax+ loss*:

$$L_{\text{spm}^+}(\boldsymbol{z}; y) = -z_y + \frac{1}{2} \sum_{j \in S^{\text{sp}}(\boldsymbol{z})} \left( z_j^2 - \tau^2(\boldsymbol{z}) \right) + \alpha \left( -\log \left( 1 - \sum_{i \notin S^{\text{sp}}(\boldsymbol{z}), i \neq y} p_i^{\text{sf}} \right) \right), \tag{4}$$

where  $\tau(z)$  is the sparsemax threshold and  $\alpha>0$  controls the strength of the suppression. For simplicity, we find that directly implementing the candidate set  $\mathcal S$  from Definition 1 using the sparsemax support  $S^{\mathsf{sp}}(z)$  achieves superior performance.

We summarize our fine-tuning strategy of **Training with Sparsemax+**, **Testing with Softmax** in Algorithm 1. From  $L_{\mathrm{spm}^+}(\boldsymbol{z};\boldsymbol{y})$  in equation 4, we see that it prevents CE-style erosion of plausible near-ties by amplifying relative ratios among top logits while nulling the rest, thereby achieving two goals: sparsemax selects a stable support set with early stopping of gradient flow, and the suppressing term explicitly drives unreasonable tokens toward zero to prevent spurious mass at inference.

#### 5 Experiments

To situate our work within the current state-of-the-art, we build upon the experimental foundation of GEM (Li et al., 2025), adopting a similar training setup. Our primary methodological difference is the substitution of the GEM objective with our proposed TS² loss. Furthermore, while GEM evaluates OpenLLM Leaderboard tasks using a standard one shot setting, we employ a multi-response, best-of-N protocol. We argue this is a more faithful and informative evaluation for diversity aware models, as it measures model's latent ability to find the correct answer rather than penalizing it for plausible "hesitation" in a single attempt.

**Setup.** We conduct experiments on two powerful, open source base models: Llama-3.1-8B and Qwen-2-7B . For supervised finetuning, we use the high quality UltraFeedback dataset (Cui et al., 2024), a large-scale corpus of preference aligned responses generated by a diverse set of models. All models are finetuned for 3 epochs using the AdamW optimizer with an effective batch size of 128. We employ a cosine learning rate schedule with an initial rate of  $2 \times 10^{-5}$  and a warm-up ratio of 0.03, a standard practice for fine-tuning modern LLMs (Yu et al., 2024; Liu et al., 2024). The maximum sequence length is capped at 2,048 tokens. For our proposed TS<sup>2</sup> method, the suppression weight  $\alpha$  (see Equation 4) is empirically determined for each model architecture, with optimal values reported alongside results. Further implementation details are provided in the Appendix B.

We compare TS<sup>2</sup> against a suite of strong and relevant baselines to provide a comprehensive evaluation: **Cross-entropy** (**CE**): The standard SFT objective, which serves as our primary baseline. **CE** with Weight Decay (**CE+WD**): A common regularization technique shown to help preserve diversity in instruction tuning (Ouyang et al., 2022; Bai et al., 2022). We use a weight decay coefficient of 0.1. **NEFTune** (**NEFT**): A regularization method that adds noise to word embeddings during training to mitigate overfitting and improve generalization (Jain et al., 2023). **GEM**: The current

state-of-the-art method for diversity preserving SFT, which we use as our main point of comparison (Li et al., 2025).

# 5.1 IMPROVING ACCURACY AND DIVERSITY IN OPEN-ENDED GENERATION

We first evaluate TS<sup>2</sup> in open ended domains to assess its ability to navigate the critical trade-off between response quality and diversity. While standard fine-tuning often improves quality at the cost of collapsing the output distribution, we hypothesize that TS<sup>2</sup> can break this trade-off by simultaneously enhancing generation quality and fostering a rich, useful diversity beneficial for sampling-based decoding. To test this, we evaluate on two distinct benchmarks. For conversational chat, we use the AlpacaEval dataset (Dubois et al., 2024) with a best-of-32 (BoN@32) protocol; a state-of-the-art reward model, FsfairX-LLaMA3-RM-v0.1 (Lambert et al., 2024), selects the best response, which is then compared against GPT-4 to determine a win rate. For code generation, we measure the pass@k metric on the HumanEvalbenchmark (Chen et al., 2021), which assesses the model's ability to generate functionally correct Python code via execution.

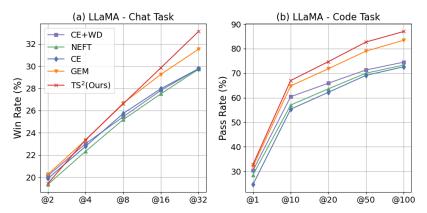


Figure 2: Performance of Llama-3.1-8B on open-ended tasks. Left: Win rate on AlpacaEval vs. sampling budget (N). Right: Pass rate on HumanEval vs. sampling budget (k).  $TS^2$  consistently outperforms baselines.

Model	Method	Win Rate (%) $\uparrow$	N-gram ↑	100 - Self-BLEU $\uparrow$	Sent-BERT $\uparrow$
	CE	29.77	17.78	47.04	9.97
	CE+WD	29.72	17.78	47.14	10.03
LLaMA-3.1-8B	NEFT	29.77	17.74	47.41	10.07
	GEM	31.53	20.32	49.82	11.16
	TS <sup>2</sup> (Ours)	33.12	23.78	53.87	12.80
	CE	31.41	17.23	16.77	7.95
	CE+WD	31.05	17.43	17.08	8.06
Qwen-2-7B	NEFT	30.36	16.59	24.59	8.06
	GEM	33.89	24.35	31.19	9.25
	TS <sup>2</sup> (Ours)	37.48	30.15	39.04	9.81

Table 1: Win rate (Best of N@32) and diversity metrics for Llama-3.1-8B and Qwen-2-7B on AlpacaEval.  $TS^2$  achieves the best results across both quality and diversity on both architectures.

**Performance on Chat and Code Generation.** As shown in Figure 2, TS<sup>2</sup> demonstrates a clear performance advantage on Llama-3.1-8B. In chat generation, its win rate at a budget of N=32 responses, reaches 33.12%, which is an improvement of 11.2% relative over the baseline cross entropy loss and a 5.0% relative improvement over the strong GEM baseline. This advantage extends to structured problem solving, on HumanEval, TS<sup>2</sup> achieves a pass@100 of 87.00%, which is 4.3% increase relative to GEM and 19.8% to that of CE. Notably, the diversity fostered by our method translates to superior sample efficiency: the pass@50 rate for TS<sup>2</sup> (82.70%) nearly matches GEM's pass@100 performance (83.40%), indicating that correct solutions can be found with fewer samples. Similar results are also observed for Qwen-2-7B model. Detailed breakdown of results for both Llama-3.1-8B and Qwen-2-7B are detailed in the Table 3.

Crucially, these performance gains do not come at the cost of diversity. As detailed in Table 1,  $TS^2$  not only achieves the highest win rate but also scores best across all three diversity metrics. It improves N-gram diversity by 17.0% ,BLEU diversity by 8.1% and sentence-bert diversity by 10.7% over GEM for LLama-3.1-8B. Similarly for Qwen-2-7B, the same metrics are improved by 23.8%, 25.1% and 6% respectively over GEM. This result confirms that  $TS^2$  successfully breaks the quality-diversity trade-off, producing responses that are simultaneously judged as higher quality by a reward model while being measurably more varied.

#### 5.1.1 DIVERSITY ON CREATIVE WRITING TASKS

To further probe the nature of the diversity generated by  $TS^2$ , we evaluate it on purely creative tasks: generating poems from 573 titles in the poetry8 dataset and stories from 500 prompts from ROCStories (Mostafazadeh et al., 2016). As shown in Table 2,  $TS^2$  once again achieves the highest scores across all three diversity metrics on both tasks, confirming its ability to produce a wider range of high-quality, creative outputs compared to all baselines.

Method		Poem		Story				
Memou	N-gram ↑	100 - Self-BLEU↑	Sent-BERT ↑	N-gram ↑	100 - Self-BLEU↑	Sent-BERT ↑		
CE	38.87	55.38	14.83	44.47	67.20	22.15		
CE+WD	38.92	55.69	14.17	44.43	67.26	22.22		
NEFT	38.80	55.68	14.13	44.31	67.21	22.04		
GEM	46.59	57.50	14.70	50.05	69.15	24.02		
TS <sup>2</sup> (Ours)	49.70	59.41	16.52	52.10	70.36	24.98		

Table 2: Diversity evaluation on creative writing tasks for Llama-3.1-8B. Higher is better.

#### 5.2 Preserving Pre-trained Capabilities on Standard Benchmarks

To assess generalization and knowledge retention, we evaluate models on six tasks from the Open-LLM Leaderboard: ARC, GSM8K, HellaSwag, MMLU, TruthfulQA, and WinoGrande. Instead of the standard greedy one-shot decoding that penalizes models preserving multiple reasoning paths, we propose a best-of-n (BoN) strategy on the OpenLLM leaderboard, which is better aligned with evaluating the capabilities of diversity-preserving models.

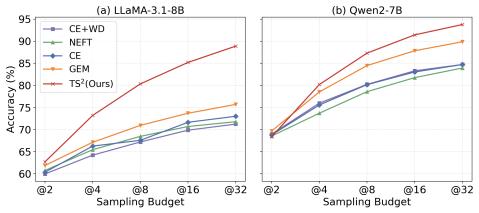


Figure 3: Average Best-of-N accuracy across six OpenLLM Leaderboard tasks. While competitive in few-shot settings (@2), TS<sup>2</sup>'s performance scales far more effectively with the sampling budget, revealing its superior knowledge retention.

We argue that a more faithful metric is Best-of-N (BoN) accuracy. This protocol measures the model's latent ability to identify the correct answer within a small sampling budget, which better reflects the true underlying capabilities of a well-calibrated, diverse model. For fair comparision, all methods are evaluated under the same BoN protocol and we report the average accuracy across all tasks.

Figure 3 validates this hypothesis. While all methods are competitive at a small sampling budget, TS<sup>2</sup>'s performance scales significantly better as 'N'(responses) increases. On Llama-3.1-8B, the

average accuracy of TS<sup>2</sup> at N=32 reaches 88.88%, a massive 13.2-point absolute (+17.4% relative) improvement over GEM (75.69%). The trend is consistent on Qwen-2-7B, where TS<sup>2</sup> again achieves the highest accuracy, demonstrating the robustness of our TS<sup>2</sup> across different model architectures.

This shows that TS<sup>2</sup> effectively preserves the model's pre-trained knowledge. Unlike CE, which collapses the distribution and discards valid alternatives, TS<sup>2</sup> maintains a clean, calibrated set of high-quality reasoning paths. With sampled responses, the model consistently finds the correct solution. A detailed breakdown of performance on each of the six tasks is provided in the Table 5.

#### 5.3 ABLATION STUDY

To assess the contribution of each component, we run an ablation study on AlpacaEval, comparing win rate against GPT-4 and BLEU diversity. TS<sup>2</sup> integrates three elements: (1) sparsemax-based training, (2) softmax decoding, and (3) a tail-suppression penalty. We evaluate three variants: **Decoupling Only** (sparsemax training, softmax inference, no penalty), **Unified Sparsemax** (sparsemax for both training and inference), and **Suppression Only** (CE loss with suppression term).

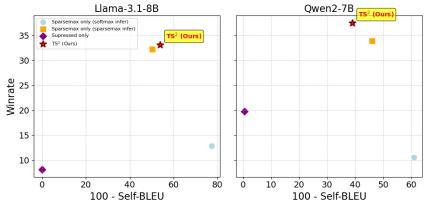


Figure 4: Ablation study on Llama-3.1-8B and Qwen-2-7B.

Figure 4 demonstrates that all components of TS<sup>2</sup> are essential. First, using the Decoupling Only strategy results in a massive increase in diversity, high BLEU diversity score, but a catastrophic drop in win rate. This shows that while decoupling unlocks variety, the suppression penalty is crucial for ensuring that this diversity is high-quality and not just uncalibrated noise.

Conversely, the Unified Sparsemax approach achieves a competitive win rate but offers lesser diversity than our full method. This confirms that the switch to softmax at inference is key to translating the learned logit geometry into a rich, sample-able probability distribution. Finally, applying the Suppression Only penalty to a standard CE baseline fails on both metrics, proving it is not a standalone improvement but works in synergy with the sparsemax-defined support set.

Meanwhile, the TS<sup>2</sup> method successfully integrates these components, achieving the best balance of high win rate and high diversity across both model architectures. This analysis confirms that the sparsemax objective, the decoupled inference, and the suppression penalty are all necessary and synergistic elements of our approach.

#### 6 CONCLUSION

In this work, we make the first step toward decoupling training and inference by adopting different prediction mappings in supervised finetuning. By combining **Sparsemax+ loss**; a tailored design that leverages margin induced sparsity with an additional suppression term for non plausible tokens; with softmax decoding at inference, our approach achieves significant improvements over existing SFT paradigms. It preserves support-aware diversity while maintaining high accuracy, thereby alleviating the alignment tax. Despite its simplicity, our method consistently outperforms CE and GEM across both chat and code tasks, achieving the highest win rates and more diverse generations. Unlike prior methods that inevitably trade off diversity against accuracy, our paradigm improves both, providing a natural remedy to distribution collapse and open up new directions for advancing alignment with broad and long term impact.

# ETHICS STATEMENT

This work investigates new algorithms for supervised fine-tuning of large language models. Our objective is to improve training stability and output diversity, thereby broadening the range of downstream applications. The methods introduced in this paper are purely algorithmic and evaluated on public datasets.

#### REPRODUCIBILITY STATEMENT

Experiment details for reproducing our numerical results can be found in Appendix B and Appendix C. To ensure anonymity and prevent potential information leakage during the review process, our source code will be released publicly after the blind review phase.

# LLM USAGE STATEMENT

We used large language model to correct grammar errors, polish the writing, and adjust the formatting of the paper.

#### REFERENCES

Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. Training a helpful and harmless assistant with reinforcement learning from human feedback, 2022. URL https://arxiv.org/abs/2204.05862.

Mathieu Blondel, André F. T. Martins, and Vlad Niculae. Learning classifiers with fenchel-young losses: Generalized entropies, margins, and algorithms. In *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pp. 606–615. PMLR, 2019. URL https://proceedings.mlr.press/v89/blondel19a.html.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code, 2021. URL https://arxiv.org/abs/2107.03374.

Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Bingxiang He, Wei Zhu, Yuan Ni, Guotong Xie, Ruobing Xie, Yankai Lin, Zhiyuan Liu, and Maosong Sun. Ultrafeedback: Boosting language models with scaled ai feedback, 2024. URL https://arxiv.org/abs/2310.01377.

Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Alpacafarm: A simulation framework for methods that learn from human feedback, 2024. URL https://arxiv.org/abs/2305.14387.

- Tilmann Gneiting and Adrian E. Raftery. Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102(477):359–378, 2007. doi: 10.1198/016214506000001437.
  - Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=rygGQyrFvH.
  - Neel Jain, Ping yeh Chiang, Yuxin Wen, John Kirchenbauer, Hong-Min Chu, Gowthami Somepalli, Brian R. Bartoldson, Bhavya Kailkhura, Avi Schwarzschild, Aniruddha Saha, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Neftune: Noisy embeddings improve instruction finetuning, 2023.
  - Jiyeon Kim, Hyunji Lee, Hyowon Cho, Joel Jang, Hyeonbin Hwang, Seungpil Won, Youbin Ahn, Dohaeng Lee, and Minjoon Seo. Knowledge entropy decay during language model pretraining hinders new knowledge acquisition. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=eHehzSDUFp.
  - Robert Kirk, Ishita Mediratta, Christoforos Nalmpantis, Jelena Luketina, Eric Hambro, Edward Grefenstette, and Roberta Raileanu. Understanding the effects of rlhf on llm generalisation and diversity, 2024. URL https://arxiv.org/abs/2310.06452.
  - Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, Noah A. Smith, and Hannaneh Hajishirzi. Rewardbench: Evaluating reward models for language modeling, 2024. URL https://arxiv.org/abs/2403.13787.
  - Mina Lee, Percy Liang, and Qian Yang. Coauthor: Designing a human-ai collaborative writing dataset for exploring language model capabilities. In *Proceedings of the 2022 CHI conference on human factors in computing systems*, pp. 1–19, 2022.
  - Ziniu Li, Congliang Chen, Tian Xu, Zeyu Qin, Jiancong Xiao, Ruoyu Sun, and Zhi-Quan Luo. Preserving diversity in supervised fine-tuning of large language models. In *International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=940YQccSM6.
  - Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation. *Advances in Neural Information Processing Systems*, 36:21558–21572, 2023.
  - Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. What makes good data for alignment? a comprehensive study of automatic data selection in instruction tuning, 2024. URL https://arxiv.org/abs/2312.15685.
  - André F. T. Martins and Ramón Fernandez Astudillo. From softmax to sparsemax: A sparse model of attention and multi-label classification. In *Proceedings of the 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pp. 1614–1623. PMLR, 2016. URL https://proceedings.mlr.press/v48/martins16.html.
  - Thomas Minka. Divergence measures and message passing. Technical Report MSR-TR-2005-173, Microsoft Research, 2005. URL https://www.microsoft.com/en-us/research/publication/divergence-measures-and-message-passing/.
  - Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong He, Devi Parikh, Dhruv Batra, Lucy Vanderwende, Pushmeet Kohli, and James Allen. A corpus and evaluation framework for deeper understanding of commonsense stories, 2016. URL https://arxiv.org/abs/1604.01696.
  - Laura O'Mahony, Leo Grinsztajn, Hailey Schoelkopf, and Stella Biderman. Attributing mode collapse in the fine-tuning of large language models. In *ICLR 2024 Workshop on Mathematical and Empirical Understanding of Foundation Models*, 2024. URL https://openreview.net/forum?id=3pDMYjpOxk.

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022.

- Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M Sadler, Wei-Lun Chao, and Yu Su. Llm-planner: Few-shot grounded planning for embodied agents with large language models. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 2998–3009, 2023.
- Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2818–2826, 2016.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Evan Z Wang, Federico Cassano, Catherine Wu, Yunfeng Bai, William Song, Vaskar Nath, Ziwen Han, Sean M. Hendryx, Summer Yue, and Hugh Zhang. Planning in natural language improves LLM search for code generation. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=48WAZhwHHw.
- Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. Magicoder: Empowering code generation with oss-instruct, 2024. URL https://arxiv.org/abs/2312.02120.
- Sean Welleck, Ilia Kulikov, Stephen Roller, Emily Dinan, Kyunghyun Cho, and Jason Weston. Neural text generation with unlikelihood training. In *International Conference on Learning Representations*, 2019.
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T. Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for large language models, 2024. URL https://arxiv.org/abs/2309.12284.

#### A RELATED WORK

Our work, TS<sup>2</sup>, intersects with three primary areas of research: supervised finetuning (SFT) and its inherent limitations, methods for enhancing generative diversity in large language models (LLMs), and the use of sparse activation functions in neural networks.

# A.1 SUPERVISED FINETUNING AND THE ALIGNMENT TAX

Supervised finetuning is a major landmark in adapting pre-trained LLMs to downstream applications, enabling them to follow instructions and adhere to specific conversational styles (Ouyang et al., 2022; Touvron et al., 2023). The standard practice involves minimizing a cross-entropy (CE) loss on a dataset of high quality datasets. While effective, this approach is known to induce an "alignment tax" (O'Mahony et al., 2024), where models become overly specialized to the finetuning distribution. This often leads to a reduction in creative capacity, a phenomenon sometimes termed "knowledge forgetting" or a collapse in output diversity (Kim et al., 2025; Li et al., 2025). The CE objective, by driving the model's posterior towards a one-hot representation of the target token, aggressively penalizes all alternative continuations, including those that are semantically plausible. This results in overconfident and deterministic models. Our work directly addresses this limitation by replacing the CE objective with a loss that preserves a set of plausible next-tokens, thereby mitigating the distributional collapse and retaining more of the pre-trained model's capabilities.

#### A.2 ENHANCING GENERATIVE DIVERSITY

Efforts to counteract the loss of diversity in finetuned LLMs can be broadly categorized into **decoding-time** and **training-time** strategies.

**Decoding-Time Strategies:** A popular line of work focuses on modifying the sampling process at inference. Techniques such as **temperature scaling**, **top-k sampling**, and **nucleus** (**top-p**) **sampling** (Holtzman et al., 2020) manipulate the output probability distribution to encourage variety. Similarly, **best-of-N sampling**, where multiple candidate responses are generated and ranked by a reward model (Bai et al., 2022), can improve output quality by exploring a wider search space. While widely used and effective, these methods are applied post-hoc and do not address the underlying overconfidence of the model's learned distribution. TS<sup>2</sup> is complementary to these techniques but fundamentally different, as it reshapes the logit geometry during training to produce a more inherently diverse and well-calibrated posterior.

**Training-Time Strategies:** Another branch of research modifies the training objective itself. **Label smoothing** (Szegedy et al., 2016) is a regularization technique that discourages overconfidence by training on soft targets. More recently, unlikelihood training was proposed to explicitly penalize undesirable tokens or repetitive patterns (Welleck et al., 2019). Closest to our work is the recent **GEM framework** (Li et al., 2025), which recasts SFT as a reverse-KL minimization problem with an entropy regularizer. GEM successfully improves diversity by preventing the model's posterior from collapsing. However, it does not enforce a hard separation between plausible and implausible tokens, potentially leaving residual mass on the long tail of the distribution. TS<sup>2</sup> offers a more direct approach: the sparsemax function provides a principled mechanism for identifying a compact support set of plausible tokens, while our proposed suppression penalty explicitly drives the probability of out-of-support tokens to zero, achieving a cleaner and more decisive separation.

#### A.3 SPARSE ACTIVATIONS IN NEURAL NETWORKS

The sparsemax function, which we leverage for our training objective, is a projection onto the probability simplex that can produce exact zeros (Martins & Astudillo, 2016). It was originally introduced as a sparse alternative to softmax for attention mechanisms and structured prediction tasks, valued for its ability to select a small subset of relevant inputs. The sparsemax loss is a specific instance of a Fenchel-Young loss, a broader class of losses that provides a unified framework for various structured prediction mappings (Blondel et al., 2019). While sparsemax has been explored for classification and attention, its application to generative LLM fine-tuning for diversity preservation is novel. Critically, our work is the first to propose a decoupled paradigm: we use the desirable properties of sparsemax (e.g., gradient masking for non-support tokens) during training but revert to the smooth, fully-supported softmax for inference. This decoupling is the key to unlocking calibrated diversity, a concept not explored in prior work that typically uses the same mapping for both training and testing.

# B EXPERIMENT DETAILS

We conduct all training on H200-141GB GPUs, employing the DeepSpeed framework with ZeRO-2 optimization and gradient checkpointing enabled. Offloading is disabled. For efficient and reproducible training, we adopt flash-attention-2 with deterministic backward passes. Our base models are Llama-3.1-8B and Qwen-2-7B, optimized using AdamW with a total batch size of 128. The learning rate is initialized at  $2\times10^{-5}$  with a warm-up ratio of 0.03 and follows a cosine decay schedule, as suggested by prior work (Yu et al., 2024; Liu et al., 2024; Cui et al., 2024),. Training is run for 3 epochs. All supervised datasets are reformatted into the chat style with the Llama-3.1-8B-Instruct and Qwen2.5-7B-Instruct tokenizer. For inference, we employ vLLM to accelerate response generation.

The supervised finetuning is done on the binarized UltraFeedback dataset curated by the HuggingfaceH4 team<sup>2</sup>, which contains 61,135 training examples and 1,000 held-out test prompts. Inputs longer than 2,048 tokens are truncated, while shorter ones are padded. To achieve a global batch size of 128, we use 4 GPUs, each with a per-device batch size of 8 and gradient accumulation of 4. A single training run requires roughly 12 GPU hours. For CE+WD baselines, the weight decay coefficients is 0.1. For NEFT, we set the noise scale to 5, consistent with Jain et al. (2023).

**Evaluation Protocol** For chatting, we use 805 prompts from AlpacaEval and score outputs with the FsfairX-LLaMA3-RM-v0.1 reward model. The maximum decoding length is 2,048, and

<sup>&</sup>lt;sup>2</sup>https://huggingface.co/datasets/HuggingFaceH4/ultrafeedback\_binarized

each prompt yields 32 samples using temperature=0.6, top-k=50, and top-p=0.9. Win rate is computed against GPT- $4^3$  responses via the Bradley–Terry model:

$$P(y \succ y'|x) = \frac{\exp(r(x,y))}{\exp(r(x,y)) + \exp(r(x,y'))}.$$

For code generation, we adopt the HumanEval benchmark (164 Python problems). Prompts follow the template of (Wei et al., 2024).

```
You are an exceptionally intelligent coding assistant that consistently delivers accurate and reliable responses to user instructions. @@ Instruction {instruction}
```

For each task, we sample 200 outputs with the same decoding configuration. The evaluation metric is pass rates, which are computed using execution-based evaluation scripts from Magicoder<sup>4</sup>.

# C ADDITIONAL RESULTS

#### C.1 CHAT AND CODE GENERATION

			T 3.54	2.4.0D								
Chat	LLaMA-3.1-8B					Qwen-2-7B						
	CE+WD	NEFT	CE	GEM	$\mathrm{TS}^2(\alpha=0.25)$	CE+WD	NEFT	CE	GEM	$\mathrm{TS}^2(\alpha=0.5)$		
@2	20.14	19.35	19.88	20.26	19.43	18.35	18.13	18.4	18.06	18.72		
@4	23.02	22.33	22.78	23.34	23.37	21.59	21.49	21.78	21.93	22.54		
@8	25.44	25.19	25.74	26.67	26.61	24.58	24.39	27.9	26.26	27.66		
@16	27.82	27.51	27.97	29.26	29.85	27.76	27.02	27.9	30.02	32.77		
@32	29.77	29.72	29.77	31.53	33.12	31.05	30.36	31.41	33.89	37.48		
Code			LLaMA	-3.1-8B		Qwen-2-7B						
0040	CE+WD	NEFT	CE	GEM	$TS^2(\alpha=0.25)$	CE+WD	NEFT	CE	GEM	$TS^2(\alpha=0.5)$		
@1	30.30	28.50	24.60	31.90	32.80	45.10	45.30	44.90	41.80	42.20		
@10	60.40	57.00	55.30	64.80	67.00	76.80	76.50	76.00	78.50	78.20		
	65.90	63.60	62.20	71.80	74.70	81.30	81.00	81.10	84.50	84.60		
@20	05.90	05.00	02.20	, 1.00								
@20 @50	71.30	70.00	69.10	79.00	82.70	84.10	83.20	83.50	87.20	87.80		

Table 3: Performance comparison of different methods on LLaMA-3.1-8B and Qwen-2-7B models for the chat and code geneartion tasks

Table 3 details the performance of both models on the open-ended generation tasks. For chat generation, it presents the win rate against GPT-4 across various best-of-N sampling budgets (N = 2, 4, 8, 16, 32). For code generation, it shows the corresponding pass@k rates for k = 1, 10, 20, 50, and 100.

#### C.2 CREATIVE WRITING

We further investigate output diversity on two creative writing tasks: poetry and short stories. For poetry, we use 573 titles drawn from the Huggingface poetry8 dataset, which covers themes such as love, nature, and mythology. For stories, we construct 500 prompts from the ROCStories dataset (Mostafazadeh et al., 2016). In both settings, the instruction is to write a piece titled "[X]" in under 200 words, where [X] is sampled from the corresponding dataset.

Diversity is measured along three dimensions following Kirk et al. (2024): **N-gram**, the fraction of distinct n-grams within a single response (intra-diversity); **Self-BLEU**, computed by treating each sample as the reference for the others (inter-diversity); **Sentence-BERT dissimilarity**, the mean

 $<sup>^3</sup>$ https://github.com/tatsu-lab/alpaca\_eval/tree/main/results/gpt4\_1106\_preview

 $<sup>^4</sup>$ https://github.com/ise-uiuc/magicoder/blob/main/experiments/text2code.py

cosine distance between generated responses in the embedding space. All scores are scaled to the range [0, 100], with higher values indicating greater diversity.

For evaluation, each model generates 16 completions per prompt using temperature=0.6, top-k=50, and top-p=0.9. Results are summarized in Table 4. It is evident that methods such as CE+WD and NEFT bring only marginal improvements in diversity. GEM consistently improves intra- and inter-diversity, while  $\mathrm{TS}^2$  achieves the highest scores.

Method (Llama-3.1-8B)		Poem		Story			
meniou (Zimim ziri oZ)	N-gram ↑	100 - Self-BLEU↑	Sent-BERT ↑	N-gram ↑	100 - Self-BLEU↑	Sent-BERT ↑	
CE+WD	38.92	55.69	14.17	44.43	67.26	22.22	
NEFT	38.80	55.68	14.13	44.31	67.21	22.04	
CE	38.87	55.38	14.83	44.47	67.20	22.15	
GEM	46.59	57.50	14.70	50.05	69.15	24.02	
$TS^2 (\alpha = 0.25)$	49.70	59.41	16.52	52.10	70.36	24.98	
Method (Owen-2-7B)		Poem			Story		
memou (Qwen 2 72)	N-gram ↑	100 - Self-BLEU↑	Sent-BERT ↑	N-gram ↑	100 - Self-BLEU↑	Sent-BERT ↑	

Method (Owen-2-7B)				2.225				
memou (Quen 2 72)	N-gram ↑	100 - Self-BLEU↑	Sent-BERT ↑	N-gram ↑	100 - Self-BLEU↑	Sent-BERT ↑		
CE+WD	44.29	44.9	8.66	56.62	50.06	19.01		
NEFT	44.37	45.09	8.55	59.66	52.2	18.94		
CE	43.94	44.92	8.56	56.44	49.83	18.86		
GEM	50.29	48.62	9.54	60.91	56.05	20.98		
$TS^2 (\alpha = 0.5)$	53.46	51.10	10.26	62.13	57.17	20.95		

Table 4: Diversity evaluation on creative writing tasks (poem and story). Higher values indicate greater diversity (N-gram, 100 - Self-BLEU, and Sentence-BERT.

#### C.3 OPENLLM LEADERBOARD TASKS

Table 5 reports results on six representative OpenLLM leaderboard tasks under varying sampling budgets. These benchmarks collectively reflect a broad spectrum of model capabilities: ARC focuses on grade-school science questions, reflecting *commonsense reasoning*; GSM8K requires multi-step solutions, capturing *mathematical reasoning*; HellaSwag emphasizes physical commonsense and narrative continuation, probing *contextual understanding*; MMLU spans 57 subjects, testing *broad factual knowledge*; TruthfulQA challenges models with common misconceptions, measuring *robustness*; and WinoGrande is a coreference benchmark, assessing *pronoun disambiguation and fine-grained language understanding*.

Building on this setup, we observe that other methods exhibit only limited gains as the sampling budget increases. In contrast, TS<sup>2</sup> consistently improves performance across tasks, achieving the largest boosts under larger budgets, often surpassing all baselines by a substantial margin. The improvements are especially pronounced for LLaMA-3.1-8B, where diversity-oriented training translates into 10–15 point gains under BoN sampling. For Qwen-2-7B, whose baseline win rates already exceed 90%, the relative gains appear smaller but still confirm the benefits of preserving diversity during training.

# C.4 MACRO- AND MICRO-LEVEL ANALYSIS OF TOKEN DISTRIBUTIONS

To understand why our method simultaneously improves accuracy and diversity, we analyze token probability distributions from two complementary perspectives: (i) a *macro-level* analysis of model outputs on a real benchmark, and (ii) a *micro-level* controlled probing task.

**Macro-level distribution.** We evaluate the models(Llama) on the AlpacaEval dataset. For each generated response, we record the probability of every selected token and compute the average probability of that response. We then plot these values across all responses to obtain a global view of the distribution. As shown in Figure 5, CE exhibits the highest mean probability ( $\approx 0.90$ ) with the smallest variance, indicating collapsed and overly uniform predictions. GEM lowers the mean probability to about 0.86 with a larger variance, consistent with its entropy-regularized updates that discourage overconfidence. Moving along the sequence  $CE \rightarrow GEM \rightarrow Sparsemax$  (sparse inference)  $\rightarrow TS^2$ , we observe a systematic trend: mean probability decreases (remaining above

				(a) I	.lama-3.1-	γR				
Method	ARC@2	ARC@4	ARC@8	ARC@16	ARC@32	Hella@2	Hella@4	Hella@8	Hella@16	Hell
CE+WD	75.59	80.27	80.27	83.28	83.61	66.95	70.55	72.95		Tien
NEFT	75.67	79.26	81.27	81.61	81.61	65.45	69.53	72.93		
CE	76.59	79.50	71.27	82.60	83.60	67.03	61.96	63.06		
GEM $TS^2(\alpha = 0.25)$	78.60 <b>78.93</b>	82.27 <b>85.95</b>	83.94 <b>88.96</b>	85.28 <b>90.30</b>	85.61 <b>91.63</b>	66.51 65.47	71.71 <b>76.76</b>	74.84 <b>84.43</b>		
$13^{-}(\alpha = 0.23)$	76.93	65.95	00.90	90.30	91.03	03.47	70.70	04.43	00.55	
Method	Wino@2	Wino@4	Wino@8 V	Vino@16 W	ino@32 MM	ILU@2 M	IMLU@4 I	MMLU@8	MMLU@16	MML
CE+WD	59.65	61.48	63.30	64.01	64.96	60.75	63.12	65.42	66.93	
NEFT CE	61.24 59.75	63.14 <b>80.27</b>	64.64 80.27	66.46 83.28	66.77 83.61	61.20 60.86	63.85 63.98	66.24 66.23	68.03 67.90	
GEM	61.80	64.64	66.69	68.43	69.85	62.04	66.12	69.32	71.85	
$TS^2(\alpha = 0.25)$	66.14	75.77	80.51	83.98	87.21	64.19	73.64	81.24	85.82	
Method	Truth@2	Truth@4	Truth@8	Truth@16	Truth@32	GSM@2	2 GSM@4	GSM@8	GSM@16	GSN
CE+WD	43.02	45.29	46.88	48.59	49.20	53.84	64.59	74.37	82.03	
NEFT	46.74	50.06	51.41			54.21				
CE GEM	43.21	45.04	47.86			54.66				
									02 (4	
	47.86 <b>51.16</b>	51.29 <b>62.42</b>	55.32 <b>71.85</b>			54.44 50.27				
$TS^2(\alpha = 0.25)$				80.05	87.39	50.27				
$TS^2(\alpha = 0.25)$	51.16	62.42	71.85	(b)	87.39 Qwen2-7	50.27 B	64.67	75.06	82.49	Hell
$\mathrm{TS}^2(\alpha=0.25)$ Method	51.16 ARC@2	62.42 ARC@4	71.85 ARC@8	(b) ARC@16	Qwen2-7] ARC@32	50.27 B Hella@2	Hella@4	75.06 Hella@8	Hella@16	
$\mathrm{TS}^2(\alpha=0.25)$ Method CE+WD	ARC@2 84.61	ARC@4 87.62	71.85 ARC@8 88.96	(b) ARC@16 89.96	Qwen2-71 ARC@32 89.96	50.27 B Hella@2 80.20	Hella@4 85.50	75.06 Hella@8 88.76	Hella@16 90.67	
$\mathrm{TS}^2(\alpha=0.25)$ Method	51.16 ARC@2	62.42 ARC@4	71.85 ARC@8	(b) ARC@16	Qwen2-7] ARC@32	50.27 B Hella@2	Hella@4	75.06 Hella@8	Hella@16	
$\frac{\mathrm{TS}^2(\alpha=0.25)}{\mathrm{Method}}$ $\frac{\mathrm{CE+WD}}{\mathrm{NEFT}}$ $\frac{\mathrm{CE}}{\mathrm{GEM}}$	51.16 ARC@2 84.61 84.94 84.61 84.94	ARC@4 87.62 87.29 87.95 88.96	71.85 ARC@8 88.96 88.62 89.62 92.30	(b) ARC@16 89.96 89.96 90.30 92.60	Qwen2-71 ARC@32 89.96 89.96 90.30 93.64	50.27  B  Hella@2  80.20  80.20  80.45  80.01	Hella@4 85.50 85.39 85.60 <b>87.54</b>	Hella@8 88.76 88.67 88.79 91.94	Hella@16 90.67 90.59 90.62 94.52	
$\frac{\mathrm{TS}^2(\alpha=0.25)}{\mathrm{Method}}$ $\frac{\mathrm{CE+WD}}{\mathrm{NEFT}}$ $\frac{\mathrm{CE}}{\mathrm{CE}}$	51.16 ARC@2 84.61 84.94 84.61	ARC@4 87.62 87.29 87.95	71.85 ARC@8 88.96 88.62 89.62	(b) ARC@16 89.96 89.96 90.30	Qwen2-71 ARC@32 89.96 89.96 90.30	50.27 B Hella@2 80.20 80.20 <b>80.45</b>	Hella@4 85.50 85.39 85.60	Hella@8 88.76 88.67 88.79	Hella@16 90.67 90.59 90.62	
$\begin{aligned} & \text{Method} \\ & \text{CE+WD} \\ & \text{NEFT} \\ & \text{CE} \\ & \text{GEM} \\ & \text{TS}^2(\alpha=0.5) \end{aligned}$	51.16  ARC@2  84.61  84.94  84.61  84.94  83.27	ARC@4 87.62 87.29 87.95 88.96 89.63	71.85 ARC@8 88.96 88.62 89.62 92.30 92.97	(b) ARC@16 89.96 89.96 90.30 92.60 93.97	Qwen2-71 ARC@32 89.96 89.96 90.30 93.64 <b>94.31</b>	50.27  B  Hella@2  80.20 80.20 80.45 80.01 75.39	Hella@4 85.50 85.39 85.60 <b>87.54</b> 87.36	Hella@8 88.76 88.67 88.79 91.94	Hella@16 90.67 90.59 90.62 94.52	
$\begin{tabular}{ll} TS^2(\alpha=0.25) \\ \hline Method \\ CE+WD \\ NEFT \\ CE \\ GEM \\ TS^2(\alpha=0.5) \\ \hline Method \\ CE+WD \\ \hline \end{tabular}$	51.16  ARC@2  84.61  84.94  84.61  84.94  83.27	ARC@4 87.62 87.29 87.95 88.96 89.63	71.85 ARC@8 88.96 88.62 89.62 92.30 92.97	(b) ARC@16 89.96 89.96 90.30 92.60 93.97	Qwen2-71 ARC@32 89.96 89.96 90.30 93.64 <b>94.31</b>	50.27  B  Hella@2  80.20 80.20 80.45 80.01 75.39	Hella@4 85.50 85.39 85.60 <b>87.54</b> 87.36	Hella@8 88.76 88.67 88.79 91.94 <b>94.18</b>	Hella@16 90.67 90.59 90.62 94.52 96.97	
$\begin{aligned} & \text{Method} \\ & \text{CE+WD} \\ & \text{NEFT} \\ & \text{CE} \\ & \text{GEM} \\ & \text{TS}^2(\alpha=0.5) \end{aligned}$ $& \text{Method}$ $& \text{CE+WD} \\ & \text{NEFT}$	ARC@2  84.61  84.94  84.61  84.94  83.27  Wino@2  70.79  71.19	ARC@4  87.62 87.29 87.95 88.96 89.63  Wino@4 V 77.34 77.26	ARC@8  88.96 88.62 89.62 92.30 92.97  Vino@8 W 81.84 82.00	(b)  ARC@16  89.96 89.96 90.30 92.60 93.97  Fino@16 Wi 83.34 83.58	Qwen2-77 ARC@32  89.96 89.96 90.30 93.64 94.31  no@32 MM 83.89 83.97	50.27  B  Hella@2  80.20 80.45 80.01 75.39  LU@2 M  69.42 69.28	Hella@4 85.50 85.39 85.60 <b>87.54</b> 87.36 MLU@4 N 77.86 71.49	Hella@8 88.76 88.67 88.79 91.94 <b>94.18</b> MMLU@8 84.53 79.83	Hella@16 90.67 90.59 90.62 94.52 96.97  MMLU@16 90.37 85.97	
$\begin{aligned} & \text{Method} \\ & \text{CE+WD} \\ & \text{NEFT} \\ & \text{CE} \\ & \text{GEM} \\ & \text{TS}^2(\alpha=0.5) \end{aligned}$ $& \text{Method}$ $& \text{CE+WD} \\ & \text{NEFT} \\ & \text{CE} \end{aligned}$	ARC@2  84.61  84.94  84.63  84.94  83.27  Wino@2  70.79  71.19  70.63	ARC@4  87.62 87.29 87.95 88.96 89.63  Wino@4 V 77.34 77.26	ARC@8  88.96 88.62 89.62 92.30 92.97  Vino@8 W 81.84 82.00 82.16	(b)  ARC@16  89.96 89.96 90.30 92.60 93.97  Fino@16 Wi 83.34 83.58 83.66	Qwen2-71 ARC@32 89.96 89.96 90.30 93.64 <b>94.31</b> no@32 MM 83.89 83.97 84.13	B Hella@2 80.20 80.20 80.45 80.01 75.39 LU@2 M 69.42 69.28 69.28	Hella@4 85.50 85.39 85.60 87.54 87.36 MLU@4 M 77.86 71.49 76.65	Hella@8 88.76 88.67 88.79 91.94 94.18 MMLU@8 84.53 79.83 83.89	Hella@16 90.67 90.59 90.62 94.52 96.97  MMLU@16 90.37 85.97 89.46	
$\begin{aligned} & \text{Method} \\ & \text{CE+WD} \\ & \text{NEFT} \\ & \text{CE} \\ & \text{GEM} \\ & \text{TS}^2(\alpha=0.5) \end{aligned}$ $& \text{Method}$ $& \text{CE+WD} \\ & \text{NEFT}$	ARC@2  84.61  84.94  84.61  84.94  83.27  Wino@2  70.79  71.19	ARC@4  87.62 87.29 87.95 88.96 89.63  Wino@4 V 77.34 77.26	ARC@8  88.96 88.62 89.62 92.30 92.97  Vino@8 W 81.84 82.00	(b)  ARC@16  89.96 89.96 90.30 92.60 93.97  Fino@16 Wi 83.34 83.58	Qwen2-77 ARC@32  89.96 89.96 90.30 93.64 94.31  no@32 MM 83.89 83.97	50.27  B  Hella@2  80.20 80.45 80.01 75.39  LU@2 M  69.42 69.28	Hella@4 85.50 85.39 85.60 <b>87.54</b> 87.36 MLU@4 N 77.86 71.49	Hella@8 88.76 88.67 88.79 91.94 <b>94.18</b> MMLU@8 84.53 79.83	Hella@16 90.67 90.59 90.62 94.52 96.97  MMLU@16 90.37 85.97	
$\begin{aligned} & \text{Method} \\ & \text{CE+WD} \\ & \text{NEFT} \\ & \text{CE} \\ & \text{GEM} \\ & \text{TS}^2(\alpha=0.5) \\ & \text{Method} \end{aligned}$	ARC@2  84.61  84.94  84.61  84.94  87.70  84.94  83.27  84.94  83.27  84.94  83.27	ARC@4  87.62  87.95  88.96  89.63  Wino@4 V  77.34  77.26  81.76  84.37	71.85  ARC@8  88.96  88.62  89.62  92.30  Vino@8 W  81.84  82.00  82.16  87.05  91.31	(b)  ARC@16  89.96 89.96 90.30 92.60  93.97  Fino@16 Wi 83.34 83.58 83.66 89.50 95.26	Qwen2-71 ARC@32 89.96 89.96 90.30 93.64 94.31 no@32 MM 83.89 83.97 84.13 90.37 95.97	B Hella@2 80.20 80.45 80.01 75.39  LU@2 M 69.42 69.28 69.37 69.01	Hella@4 85.50 85.39 85.60 <b>87.54</b> 87.36 MLU@4 M 77.86 71.49 76.65 80.36 <b>82.11</b>	Hella@8 88.76 88.67 88.79 91.94 94.18 MMLU@8 84.53 79.83 83.89 88.86 90.67	Hella@16 90.67 90.59 90.62 94.52 96.97  MMLU@16  90.37 85.97 89.46 93.86 94.84	MML
$\begin{aligned} & \text{Method} \\ & \text{CE+WD} \\ & \text{NEFT} \\ & \text{CE} \\ & \text{GEM} \\ & \text{TS}^2(\alpha=0.5) \\ & \text{Method} \\ & \text{CE+WD} \\ & \text{NEFT} \\ & \text{CE} \\ & \text{GEM} \\ & \text{TS}^2(\alpha=0.5) \\ & \text{Method} \\ & Meth$	ARC@2  84.61  84.94  84.61  84.94  83.27  Wino@2  70.79  71.19  70.63  73.63  71.11  Truth@2	ARC@4  87.62  87.99  87.95  88.96  89.63  Wino@4 V  77.34  77.26  81.76  84.37	ARC@8  88.96 88.62 89.62 92.30 92.97  Vino@8 W 81.84 82.00 82.16 87.05 91.31	(b)  ARC@16  89.96 89.96 90.30 92.60 93.97  7ino@16 Wi 83.34 83.58 83.66 89.50 95.26	Qwen2-71 ARC@32 89.96 89.96 90.30 93.64 94.31  no@32 MM 83.89 83.97 84.13 90.37 95.97	B Hella@2 80.20 80.45 80.01 75.39  LU@2 M 69.42 69.28 69.37 69.91 69.01	Hella@4  85.50 85.39 85.60 87.54 87.36  MLU@4 N  77.86 71.49 76.65 80.36 82.11	Hella@8 88.76 88.67 88.79 91.94 94.18 4MLU@8 84.53 79.83 83.89 88.86 90.67	Hella@16 90.67 90.59 90.62 94.52 96.97  MMLU@16 90.37 85.97 89.46 93.86 94.84  GSM@16	MML
$\begin{tabular}{ll} Method & CE+WD \\ NEFT & CE \\ GEM \\ TS^2(\alpha=0.5) & \\ \hline Method & CE+WD \\ NEFT \\ CE \\ GEM \\ Method & CE+WD \\ \hline Method & CE+WD \\ \hline \end{tabular}$	ARC@2 84.61 84.94 84.61 84.94 83.27 Wino@2 70.79 71.19 70.63 73.63 71.11 Truth@2 45.65	ARC@4 87.62 87.29 87.95 88.96 89.63 Wino@4 V 77.34 77.26 77.26 81.76 84.37 Truth@4 49.69	ARC@8  88.96 88.62 89.62 92.30 92.97  Vino@8 W 81.84 82.00 82.16 87.05 91.31  Truth@8  52.75	(b)  ARC@16  89.96  89.96  90.30  92.60  93.97  ino@16 Wi 83.34 83.58 83.66 89.50  95.26  Truth@16  55.07	Qwen2-71 ARC@32 89.96 89.96 90.30 93.64 94.31  no@32 MM 83.89 83.97 84.13 90.37 95.97  Truth@32 55.93	B Hella@2 80.20 80.20 80.45 80.41 75.39  LU@2 M 69.42 69.28 69.37 69.91 69.01  GSM@2	Hella@4 85.50 85.39 85.60 <b>87.54</b> 87.36 MLU@4 N 77.86 80.36 <b>82.11</b> GSM@4	Hella@8 88.76 88.67 91.94 94.18 MMLU@8 84.53 79.83 83.89 88.86 90.67 GSM@8	Hella@16 90.67 90.59 90.62 94.52 96.97  MMLU@16 90.37 85.97 89.46 93.86 94.84  GSM@16	MML
$\begin{aligned} & \text{Method} \\ & \text{CE+WD} \\ & \text{NEFT} \\ & \text{CE} \\ & \text{GEM} \\ & \text{TS}^2(\alpha=0.5) \\ & \text{Method} \\ & \text{CE+WD} \\ & \text{NEFT} \\ & \text{CE} \\ & \text{GEM} \\ & \text{TS}^2(\alpha=0.5) \\ & \text{Method} \\ & Meth$	ARC@2  84.61  84.94  84.61  84.94  83.27  Wino@2  70.79  71.19  70.63  73.63  71.11  Truth@2	ARC@4  87.62  87.99  87.95  88.96  89.63  Wino@4 V  77.34  77.26  81.76  84.37	71.85  ARC@8  88.96 88.62 89.62 92.30 92.97  Vino@8 W 81.84 82.00 82.16 87.05 91.31  Truth@8  52.75 52.50	(b)  ARC@16  89.96 89.96 90.30 92.60 93.97  7ino@16 Wi 83.34 83.58 83.66 89.50 95.26	Qwen2-71 ARC@32 89.96 89.96 90.30 93.64 94.31  no@32 MM 83.89 83.97 84.13 90.37 95.97	B Hella@2 80.20 80.45 80.01 75.39  LU@2 M 69.42 69.28 69.37 69.91 69.01	Hella@4 85.50 85.39 85.60 87.54 87.36 MLU@4 M 77.86 71.49 76.65 80.36 82.11 GSM@4 77.86 71.49	Hella@8 88.76 88.67 88.79 91.94 94.18 4MLU@8 84.53 79.83 83.89 88.86 90.67	Hella@16 90.67 90.59 90.62 94.52 96.97  MMLU@16 90.37 85.97 89.46 93.86 94.84  GSM@16 90.37 85.97	Hell:
$\begin{split} & \text{Method} \\ & \text{CE+WD} \\ & \text{NEFT} \\ & \text{CE} \\ & \text{GEM} \\ & \text{TS}^2(\alpha=0.5) \end{split}$ $& \text{Method} \\ & \text{CE+WD} \\ & \text{NEFT} \\ & \text{CE} \\ & \text{GEM} \\ & \text{TS}^2(\alpha=0.5) \end{split}$ $& \text{Method} \\ & \text{CE+WD} \\ & \text{NEFT} \\ & \text{CE} \\ & \text{GEM} \\ & \text{TS}^2(\alpha=0.5) \end{split}$ $& \text{Method} \\ & \text{CE+WD} \\ & \text{NEFT} \\ & \text{Method} \\ & \text{CE+WD} \\ & \text{NEFT} \\ & \text{Method} \\ & \text{CE-WD} \\ & \text{CE-WD}$	ARC@2  84.61  84.94  84.61  84.94  83.27  Wino@2  70.79  71.19  70.63  73.63  71.11  Truth@2  45.65  45.41	ARC@4  87.62 87.29 87.95 88.96 89.63  Wino@4 V 77.34 77.26 81.76 84.37  Truth@4  49.69 49.32	ARC@8  88.96 88.62 89.62 92.30 92.97  Vino@8 W 81.84 82.00 82.16 87.05 91.31  Truth@8  52.75	(b)  ARC@16  89.96 89.96 90.30 92.60 93.97  Fino@16 Wi 83.34 83.58 83.66 89.50 95.26  Truth@16  55.07 54.46	Qwen2-71 ARC@32 89.96 89.96 90.30 93.64 94.31  no@32 MM 83.89 83.97 84.13 90.37 95.97  Truth@32 55.93 55.69	B Hella@2 80.20 80.20 80.45 80.61 75.39  LU@2 M 69.42 69.28 69.37 69.91 69.01  GSM@2 66.72 67.32	Hella@4 85.50 85.39 85.60 87.54 87.36 MLU@4 N 77.86 71.49 76.65 80.36 82.11 GSM@4 77.86 71.49 76.65 80.36 82.11	Hella@8 88.76 88.67 91.94 94.18 MMLU@8 84.53 79.83 83.89 88.86 90.67 GSM@8 84.53 79.83	Hella@16 90.67 90.59 90.62 94.52 96.97  MMLU@16 90.37 85.97 89.46 93.86 94.84	Hella MML

Table 5: Pass Rate (%) of Different Methods on 6 OpenLLM leaderboard tasks under Various Sampling Budgets.

0.8), while variance increases, revealing a more balanced allocation of probability mass to plausible alternatives.

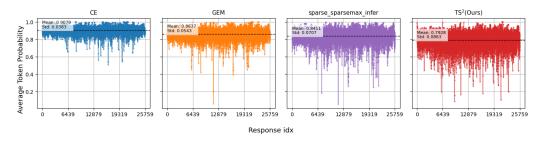


Figure 5: Macro-level analysis: average selected token probability distribution on AlpacaEval.

**Micro-level probing.** To complement the macro-level view, we design a controlled probing task to test whether models can distribute probability mass across relevant candidates. We prompt the model with the few-shot instruction to generate a single-digit number. Each model is queried 100 times. Whenever a digit is generated, we record the probability distribution of the top-300 tokens. Finally, we compute the average probability of each token across the 100 trials, resulting in a fine-grained view of how probability mass is allocated.

```
You're an AI assistant, I will give you an example of following question.

Example:
User: Give me a word of fruit.

Assistant: Apple.

Now you follow the format of the example,

Give me a single-digit number,

Answer:
```

The results, shown in Figure 1, reveal stark differences. **CE** collapses to a one-hot distribution: the chosen digit monopolizes probability, while the tail is filled with irrelevant tokens. **GEM** retains a few candidate digits but remains nearly one-hot, yielding limited diversity. **Sparsemax** (**Sparsemax-infer**) distributes mass across more digits, but still assigns non-negligible probability to spurious tokens. In contrast,  $TS^2$  combines sparsemax, which preserves probability on relevant digits, with the suppressing loss, which eliminates unrelated characters. This synergy results in distributions that are both diverse and accurate.

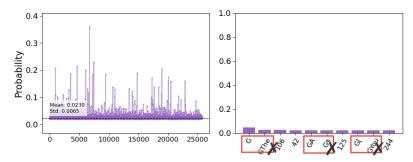


Figure 6: Macro- and Micro-level probling: Sparsemax Training and Softmax Inference

As a special case, we also examine the strategy of **sparsemax training with softmax inference**(shown in Figure 6). In the *macro-level probing*, this setting produces a distribution that is close to uniform, suggesting that the model does not exhibit clear preferences over candidate tokens. In the *micro-level probing task*, we observe that although some valid numerical answers (such as "42" or "125") appear, a large number of irrelevant tokens also receive comparable probability mass. As a result, the model's outputs become difficult to interpret, and its effective generation ability is diminished. This illustrates why conventional diversity metrics may report artificially high scores in this case: while probability is spread across many tokens, much of it corresponds to spurious rather than meaningful outputs.

#### C.5 HARD THRESHOLD

CE	0.25-top3	0.5-top3	1.0-top3	0.5-top5	0.5-top10	1.0-top10	$0.25$ -TS $^{2}$	$0.5$ -TS $^2$
50.00	50.16	52.58	48.80	49.86	47.58	46.61	53.73	51.03

Table 6: Results on Llama-3.2-1B with different  $\alpha$  and sparsification strategies. All tokens except target are thrown out the support set.

We also evaluate the Llama-3.2-1B under different values of  $\alpha$  and supersession strategies. Here, "top-k" means only the largest k logits are preserved in the defined support set, while all other tokens (except the target) are thrown out the support set, and "TS²" denotes our proposed two-stage suppression method. We take the vanilla cross-entropy (CE) training as the baseline, which yields a score of 50.

From the results, we observe two main trends: (i) Increasing  $\alpha$  from 0.25 to 1.0 generally decreases performance, indicating that larger  $\alpha$  values reduce the model's output diversity. (ii) Within the top-k setting, smaller k (e.g., top-3 vs. top-10) leads to higher diversity and better scores, while larger k values dilute the distribution and hurt performance. Overall, both reducing  $\alpha$  and carefully selecting smaller k encourage the model to maintain useful diversity, while our  $\mathrm{TS}^2$  method further boosts results beyond simple top-k truncation.

#### D DETAILED PROOFS

Corollary 1 If Definition 1 holds and  $\varepsilon_{\text{tail}} < \varepsilon_{\text{head}}$ , then  $\max_{j \notin \mathcal{S}} p_j \le \varepsilon_{\text{tail}} < \varepsilon_{\text{head}} \le \min_{i \in \mathcal{S}} p_i$ , so each plausible sample has strictly higher probability than any tail sample.

*Proof.* From tail suppression,  $\sum_{j \notin \mathcal{S}} p_j \leq \varepsilon_{\text{tail}}$ , hence  $\max_{j \notin \mathcal{S}} p_j \leq \varepsilon_{\text{tail}}$ . From head preservation,  $\min_{i \in \mathcal{S}} p_i \geq \varepsilon_{\text{head}}$ . Combine with the condition  $\varepsilon_{\text{tail}} < \varepsilon_{\text{head}}$ , we complete the proof.

**Corollary** 2 If all probability mass collapses onto the ground-truth token, i.e.,  $p_y=1$  and  $p_{y'}=0 \ \forall y'\neq y$ , then p fails to qualify the TSPD  $(m(\geq 2), \varepsilon_{\rm head}, \varepsilon_{\rm tail})$ .

*Proof.* For  $m \geq 2$ ,  $S = \text{Top}_m(\mathbf{p})$  contains y and some  $y' \neq c$  with  $p_{y'} = 0$ , violating  $\min_{j \in S} p_j \geq \varepsilon_{\text{head}} > 0$ .

**Lemma** 3 [Gradients vanish outside the sparsemax support] Let  $p = \operatorname{sparsemax}(z)$  and  $S^{\operatorname{sp}}(z)$  be its support. Consider the sparsemax loss  $\mathcal{L}_{\operatorname{sp}}(z,y)$  with target y. If  $y \in S^{\operatorname{sp}}(z)$ , then  $\forall i \notin S^{\operatorname{sp}}(z)$ ,  $\frac{\partial \mathcal{L}_{\operatorname{sp}}(z,y)}{\partial z_i} = 0$ .

*Proof.* The gradient satisfies  $\nabla_{\mathbf{z}} \mathcal{L}_{sp}(\mathbf{z}, y) = \mathbf{p} - \mathbf{e}_y$ . For  $i \notin S^{\mathsf{sp}}(\mathbf{z})$  we have  $p_i = 0$ , and under the assumption  $y \in S^{\mathsf{sp}}(\mathbf{z})$  we have  $i \neq y$ , hence  $\partial \mathcal{L}_{sp}(\mathbf{z}, y) / \partial z_i = 0$ .

**Theorem** 4 [Sparsemax expands pairwise gaps faster than softmax] Let  $z \in \mathbb{R}^K$ ,  $p^{sf} = \operatorname{softmax}(z)$ , and  $p^{sp} = \operatorname{sparsemax}(z)$ . For any indices  $i \neq j$ , let  $u := z_i - z_j$  and we have

$$\begin{split} \frac{\partial}{\partial u} (p_i^{\rm sp} - p_j^{\rm sp}) &= 1 \quad \forall \, i, j \in \mathcal{S}^{\rm sp} \\ \frac{\partial}{\partial u} \left( p_i^{\rm sf} - p_j^{\rm sf} \right) &< 1 \end{split} \qquad \qquad \text{softmax} \end{split}$$

*Proof.* Inside the sparsemax support, we have  $p_j^{\sf sp} = z_j - \tau(\boldsymbol{z})$  and  $p_i^{\sf sp} - p_j^{\sf sp} = (z_i - \tau(\boldsymbol{z})) - (z_j - \tau(\boldsymbol{z})) = z_i - z_j$ , thus  $\frac{\partial}{\partial u}(p_i^{\sf sp} - p_j^{\sf sp}) = 1$ . For softmax, using the Jacobian  $\nabla \boldsymbol{p}^{\sf sf} = \operatorname{diag}(\boldsymbol{p}^{\sf sf}) - \boldsymbol{p}^{\sf sf}(\boldsymbol{p}^{\sf sf})^{\top}$  and differentiating only in the direction  $z_i \uparrow$ ,  $z_j \downarrow$  (other logits fixed) yields  $\frac{\partial}{\partial u}(p_i^{\sf sf} - p_j^{\sf sf}) = p_i^{\sf sf} + p_j^{\sf sf} - (p_i^{\sf sf} - p_j^{\sf sf})^2$ , which is strictly < 1 for finite  $\boldsymbol{z}$ .

**Corollary** 5 [Softmax remains TSPD-valid when sparsemax is one-hot] Let  $z \in \mathbb{R}^K$  with  $y = \arg\max_j z_j$ , and  $\delta_j := z_y - z_j$ . Assume sparsemax is one-hot at y, i.e.,  $\delta_{\min} := \min_{j \neq y} \delta_j \geq \gamma > 0$  (e.g.,  $\gamma = 1$ ), and the top-m head is bounded:  $\delta_{(k)} := z_c - z_{(k)} \leq B \ \forall k = 2, \ldots, m$ . Set

 $A_m = m + (K - m)e^{-\gamma}$ . Then for  $p^{sf} = \operatorname{softmax}(z)$  we have

$$p_y^{\mathrm{sf}} \geq \frac{1}{A_m}, \quad p_{(k)}^{\mathrm{sf}} \geq \frac{e^{-B}}{A_m} \; (\forall k = 2, \dots, m), \quad \sum_{k > m} p_{(k)}^{\mathrm{sf}} \leq \frac{(K - m)e^{-\gamma}}{A_m}.$$

Consequently,  $p^{sf}$  satisfies TSPD of order m with any thresholds  $0 < \varepsilon_{\text{head}} \le \frac{e^{-B}}{A_m}$ ,  $\frac{(K-m)e^{-\gamma}}{A_m} \le \varepsilon_{\text{tail}} \le 1 - m \, \varepsilon_{\text{head}}$ .

*Proof.* For any j,

$$p_j^{\rm sf} \ = \ \frac{e^{z_j}}{\sum_k e^{z_k}} \ = \ \frac{e^{-(z_y-z_j)}}{1+\sum_{k\neq y} e^{-(z_y-z_k)}} \ = \ \frac{e^{-\delta_j}}{\Omega}, \quad \text{where } \Omega := 1+\sum_{k\neq y} e^{-\delta_k}.$$

Then,  $\forall 2 \leq k \leq m$ , we have  $e^{-\delta_{(k)}} \in [e^{-B},1]$  according to the head bound  $\delta_{(k)} \leq B$ ;  $\forall k > m$ , we have  $e^{-\delta_{(k)}} \leq e^{-\gamma}$  according to the sparsemax one-hot margin  $\delta_{(k)} \geq \gamma$ .

To lower-bound  $p_i^{sf}$ , we upper-bound  $\mathcal{C}$  by taking the largest possible contributions in each group:

$$C = 1 + \sum_{k=2}^{m} e^{-\delta_{(k)}} + \sum_{k>m} e^{-\delta_{(k)}} \le 1 + (m-1) \cdot 1 + (K-m) e^{-\gamma} = A_m.$$

Therefore, we have

$$p_y^{\mathsf{sf}} = \frac{1}{\mathcal{C}} \geq \frac{1}{A_m}, \quad p_{(k)}^{\mathsf{sf}} = \frac{e^{-\delta_{(k)}}}{\mathcal{C}} \geq \frac{e^{-B}}{A_m} \quad (k = 2, \dots, m).$$

For k>m,  $\delta_{(k)}\geq \gamma$  gives  $\sum_{k>m}p_{(k)}^{\rm sf}\leq \frac{(K-m)e^{-\gamma}}{A_m}$ . We complete the proof.

# E ADDITIONAL TECHNICAL ANALYSIS

#### E.1 A NEW TRAINING PARADIGM

Training with CE loss leads to distribution collapse: under gradient descent, the predictive distribution p converges to the target y. This causes over-confident and degenerate predictions at inference.

To address this issue, we discuss a new paradigm consisting of three steps:

1. Inflation during training. Given  $p^{sf} = \operatorname{softmax}(z)$ , we define an inflated distribution

$$\tilde{p}_i = \frac{f(p_i^{\mathsf{sf}})}{\sum_{j=1}^K f(p_j^{\mathsf{sf}})}, \ i = 1, 2, \dots, K,$$

where  $f:[0,1]\to\mathbb{R}_+$  is strictly increasing and satisfies a ratio amplification property:

$$p_i^{\mathsf{sf}} > p_j^{\mathsf{sf}} \implies \frac{f(p_i)}{f(p_j)} > \frac{p_i^{\mathsf{sf}}}{p_j^{\mathsf{sf}}}.$$

- 2. Loss applied on the inflated distribution. We train by minimizing a tailored loss  $\ell(\tilde{p}, y)$ . Ratio-amplifying inflation accelerates the collapse of  $\tilde{p}$  to one-hot.
- 3. **Softmax inference.** At test time, predictions are made with the original *p*, which remains smooth and calibrated.

This paradigm improves optimization dynamics while preserving smooth probabilistic predictions.

**Theorem 6** (Invertible  $\phi$ -mappings training prevents collapse at inference). Given the predictive distribution p, let  $\ell$  be a strictly proper loss and  $f:[0,1]\to\mathbb{R}_+$  be strictly increasing and invertible. Define the inflated distribution

$$\tilde{\boldsymbol{p}} = \Phi(\boldsymbol{p}), \quad \Phi(p)_i = \frac{f(p_i)}{\sum_{j=1}^K f(p_j)}$$

where i = 1, 2, ..., K.

- 1. (Training) Under gradient descent,  $\tilde{p}_y = 1$  when this loss converges to 0, i.e., the inflated distribution collapses to the one-hot label.
- 2. (Inference) Define the recovered distribution

$$p_i^* = \frac{f^{-1}(\tilde{p}_i)}{\sum_{j=1}^K f^{-1}(\tilde{p}_j)}, i = 1, 2, \dots, K.$$

Then,  $p^*$  remains strictly inside the simplex, that is

$$\sum_{j=1}^{K} p_{j}^{*} = 1, \text{ and } 0 < p_{j}^{*} < 1 \ \forall j.$$

In particular,  $p^*$  never collapses to a one-hot vector.

*Proof.* (1) For any strictly proper loss  $\ell$ , the stationary condition

$$\nabla_{\boldsymbol{p}}\ell(\boldsymbol{p},y) = 0 \iff p_y = 1$$

- $\Longrightarrow$  the predictive distribution  $\boldsymbol{p}$  converges to the one-hot label. Since  $\Phi$  is a bijection onto the simplex (as f is strictly increasing), minimizing  $\ell(\Phi(\boldsymbol{p}),y)$  w.r.t.  $\boldsymbol{p}$  is equivalent to minimizing  $\ell(\tilde{\boldsymbol{p}},y)$  w.r.t.  $\tilde{\boldsymbol{p}}$ . Thus, under gradient descent in the inflated space, we obtain  $\tilde{p}_y=1$ .
- (2) For inference, we recover  $p^*$  from  $\tilde{p}$  via

$$p_i^* = \frac{f^{-1}(\tilde{p}_i)}{\sum_{j=1}^K f^{-1}(\tilde{p}_j)}, \quad i = 1, 2, \dots, K.$$

- Since  $\tilde{p} \in \Delta^{K-1}$ , we have  $0 \leq \tilde{p}_i \leq 1$  and  $\sum_i \tilde{p}_i = 1$ . Because  $f^{-1}$  is strictly increasing and continuous, we have  $f^{-1}(\tilde{p}_i)0 \ \forall i$ . Hence  $p_i^* \geq 0 \ \forall i$ , and normalization ensures  $\sum_i p_i^* = 1$ .
- To show non-collapse, suppose by contradiction that  $p_y^* = 1$ . Then  $p_j^* = 0 \ \forall j \neq y$ . But this would require  $f^{-1}(\tilde{p}_j) = 0 \ \forall j \neq y$ , i.e.  $\tilde{p}_j = f(0)$ . Since  $\tilde{p}_j > 0$  (strictly inside the simplex), this is impossible. Thus  $p^*$  cannot be a one-hot vector.
- Therefore,  $p^*$  remains a smooth distribution in the simplex, preventing distribution collapse at inference.
- **Limit analysis.** Suppose  $p_y = 1 \epsilon$  with  $\epsilon > 0$  distributed among other coordinates so that  $p_j > 0$  for some  $j \neq y$ . Then

$$\frac{\tilde{p}_y}{\tilde{p}_j} = \frac{f(1-\epsilon)}{f(\epsilon)}.$$

Since f is strictly increasing and satisfies ratio amplification, we have

$$\lim_{\epsilon \to 0^+} \frac{f(1-\epsilon)}{f(\epsilon)} = +\infty.$$

Therefore,

$$\lim_{\epsilon \to 0^+} \tilde{p}_y = 1, \quad \lim_{\epsilon \to 0^+} \tilde{p}_j = 0.$$

In contrast, for the original distribution p we only have

$$p_y = 1 - \epsilon < 1, \qquad p_j = \epsilon > 0.$$

- Thus, the inflated distribution  $\tilde{p}$  achieves the one-hot collapse strictly earlier, while the underlying p remains smooth with strictly positive mass on all coordinates.
- At inference time, we return to p by applying the original activation function (e.g., softmax). This ensures the predicted distribution is smoother and less degenerate than one-hot, even though the training dynamics in the inflated space enforced early collapse.

**Theorem 7** (Sparsemax as a piecewise ratio-amplifying  $\phi$ -mapping of softmax). Let  $z \in \mathbb{R}^K$  be a logit vector,  $p^{sf} = \operatorname{softmax}(z) \in \Delta^{K-1}$  with

$$p_i^{\rm sf} = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}, \quad i = 1, 2, \dots, K;$$

and  $p^{\sf sp} = \operatorname{sparsemax}(z) \in \Delta^{K-1}$  with

$$p_i^{\mathsf{sp}} = \max\{z_i - \tau(\boldsymbol{z}), 0\}, \quad \sum_i p_i^{\mathsf{sp}} = 1.$$

Then we define

$$p_i^{\text{sp}} = \Phi(p)_i = \frac{f(p_i)}{\sum_{j=1}^K f(p_j)},$$

where  $f:[0,1]\to\mathbb{R}_{>0}$  is the piecewise function  $f(x)=\max\{\log x-\theta,0\}, \forall \theta\in\mathbb{R}$ .

*Proof.* According to  $p_i^{sf} = \frac{e^{z_i}}{\sum_j^K e^{z_j}}$ , we have  $z_i = \log p_i^{sf} + C$  with  $C = \log \sum_j e^{z_j}$ . Substituting this into the definition of sparsemax,

$$p_i^{\mathsf{sp}} = \max\{\log p_i^{\mathsf{sf}} + C - \tau(\boldsymbol{z}), 0\}.$$

Letting  $\theta = \tau(z) - C$ , we obtain

$$p_i^{\mathsf{sp}} = \max\{\log p_i - \theta, 0\}.$$

Since  $\sum_i p_i^{\sf sp} = 1$ , normalizing yields

$$p_i^{\mathsf{sp}} = \frac{\max\{\log p_i - \theta, 0\}}{\sum_j \max\{\log p_j - \theta, 0\}}.$$

We now analyze the following two cases.

Case I (support set  $S^{\mathsf{sp}}(z) = \{i : \log p_i^{\mathsf{sf}} > \theta\}$ ). For  $i \in S$ ,  $f(p_i^{\mathsf{sf}}) = \log p_i^{\mathsf{sf}} - \theta > 0$ . On (0, 1],  $\log x$  is strictly increasing; subtracting  $\theta$  preserves this property. Therefore  $\forall i, j \in S$  with  $p_i^{\mathsf{sf}} > p_j^{\mathsf{sf}}$ , we obtain the ratio amplification property:

$$\frac{\Phi(\boldsymbol{p})_i}{\Phi(\boldsymbol{p})_j} = \frac{\log p_i^{\mathsf{sf}} - \theta}{\log p_i^{\mathsf{sf}} - \theta} > \frac{p_i^{\mathsf{sf}}}{p_i^{\mathsf{sf}}}.$$

Thus  $\Phi$  inflates the relative ratios within the support.

Case 2 (outside the support  $S^{sp}(z)$ ). For  $j \notin S^{sp}(z)$ , we have  $\log p_j^{sf} \leq \theta$  and hence  $f(p_j^{sf}) = 0$ . Therefore,

$$\Phi(\boldsymbol{p})_j = \frac{0}{\sum_{i \in S^{\mathrm{sp}}(\boldsymbol{z})} f(p_i^{\mathrm{sf}})} = 0.$$

By contrast,  $p_j^{\sf sf} > 0$  since  $p = \operatorname{softmax}(z)$  has full support. Thus  $\operatorname{sparsemax}(z)$  coincides with  $\Phi(p)$ , where  $\Phi$  is generated by the piecewise ratio-amplifying function f.

Overall,  $\operatorname{sparsemax}(z)$  is a piecewise ratio-amplifying inflation of  $\operatorname{softmax}(z)$ . Training on  $\Phi(p)$  drives the inflated distribution to collapse to one-hot on its support, while inference with the original softmax p preserves strictly positive mass on all coordinates. This prevents the predictive distribution from degenerating into an exact one-hot vector at inference.

Having established sparsemax as a concrete instance of ratio-amplifying inflation, it is natural to ask whether other mappings f might be equally effective, or perhaps even more suitable in specific contexts. To answer this, we next examine the general collapse condition in the binary case.

E.2 GENERAL COLLAPSE CONDITION IN THE BINARY CASE

Consider binary classification with p = (p, 1 - p) and label y = 1. The inflated distribution is

$$\tilde{p}_1 = \frac{f(p)}{f(p) + f(1-p)}, \quad \tilde{p}_2 = 1 - \tilde{p}_1.$$

Define the ratio

$$R(p) = \frac{f(1-p)}{f(p)}.$$

Then

$$\tilde{p}_1 = \frac{1}{1 + R(p)}.$$

For a precision parameter  $\epsilon > 0$ , we say collapse occurs if

$$\tilde{p}_1 \ge 1 - \epsilon \iff R(p) \le \frac{\epsilon}{1 - \epsilon}.$$

**1. Power inflation.** For  $f(x) = x^{\alpha}$ ,  $\alpha > 1$ ,

$$R(p) = \left(\frac{1-p}{p}\right)^{\alpha}.$$

Collapse condition:

$$p > \frac{1}{1 + \left(\frac{\epsilon}{1 - \epsilon}\right)^{1/\alpha}}.$$

**2. Exponential inflation.** For  $f(x) = e^{\gamma x}$ ,  $\gamma > 0$ ,

$$R(p) = e^{\gamma(1-2p)}.$$

Collapse condition:

$$p > \frac{1}{2} + \frac{1}{2\gamma} \log \frac{1 - \epsilon}{\epsilon}.$$

**3. Logarithmic inflation.** For  $f(x) = \log(x + \delta)$  with  $\delta > 0$ ,

$$R(p) = \frac{\log(1 - p + \delta)}{\log(p + \delta)}.$$

Collapse condition:

$$\frac{\log(1-p+\delta)}{\log(p+\delta)} < \frac{\epsilon}{1-\epsilon}.$$

#### E.3 GRADIENT DYNAMICS UNDER RATIO AMPLIFICATION

The ratio-amplifying property of  $\phi$ -mappings not only accelerates the collapse of  $\tilde{p}$ , but also reshapes the gradient dynamics during training. For a strictly proper loss  $\ell$ , the gradient w.r.t. logits z is assumed to be

$$\nabla_{\boldsymbol{z}}\ell(\boldsymbol{z};y) = \boldsymbol{p} - \boldsymbol{e}_{y}, \ \boldsymbol{p} = g(\boldsymbol{z}),$$

where  $g(\cdot)$  denotes a probability distribution obtained from the logits z and  $e_y$  is a one-hot vector with the y-th entry equals 1.

When training on the inflated distribution  $\tilde{p} = \Phi(p)$ , the chain rule gives

$$\nabla_{\boldsymbol{z}}\ell(\tilde{\boldsymbol{p}},y) = \frac{\partial \tilde{\boldsymbol{p}}}{\partial \boldsymbol{p}} \cdot (\tilde{\boldsymbol{p}} - y),$$

where  $\frac{\partial \tilde{p}}{\partial p}$  is the Jacobian of the inflation operator.

 **Effect of ratio amplification.** Suppose  $f:[0,1]\to\mathbb{R}_+$  is strictly increasing and ratio-amplifying, so that

$$\frac{\tilde{p}_y}{\tilde{p}_j} > \frac{p_y}{p_j}, \quad \forall j \neq y.$$

This guarantees that the *relative gap* between the correct and incorrect probabilities grows under  $\Phi$ . Hence, even if the exact magnitude of each gradient entry depends on the Jacobian structure, the ratio

$$\frac{\left|\nabla_{z_y}\right|}{\left|\nabla_{z_i}\right|}$$

is enlarged compared to the original probability space. In other words, the margin  $z_y-z_j$  receives stronger effective gradient pressure to grow. Intuitively, because  $\tilde{p}_y>p_y$  and  $\tilde{p}_j< p_j$  for  $j\neq y$ , the gradient signal on the correct logit  $z_y$  is reinforced, while the signals on the incorrect logits  $z_j$  are diminished. This rescaling accelerates the suppression of false classes and boosts the dominance of the true class. Although the absolute gradient values are determined by both  $\tilde{p}$  and the Jacobian  $\frac{\partial \tilde{p}}{\partial p}$ , the effective separation between correct and incorrect classes is consistently larger under ratio-amplifying mappings.

**Summary.** Any  $\phi$ -mapping with ratio amplification reshapes the optimization dynamics by preconditioning the gradient flow:

- The *relative strength* of gradients is tilted further in favor of the true class.
- Incorrect classes are suppressed earlier, as their probabilities are diminished more aggressively.

Consequently, the system reaches effective one-hot collapse earlier than when training directly on p. Crucially, since inference is carried out with the original distribution p, the final predictions remain smooth and non-degenerate, preserving diversity while benefiting from sharper supervision during training.