

# A KL Lens on Quantization: Fast, Forward-Only Sensitivity for Mixed-Precision SSM–Transformer Models

Anonymous CVPR submission

Paper ID 21

## Abstract

001 *Deploying Large Language Models (LLMs) on edge devices*  
002 *faces severe computational and memory constraints, limit-*  
003 *ing real-time processing and on-device intelligence. Hy-*  
004 *brid architectures combining Structured State Space Mod-*  
005 *els (SSMs) with transformer-based LLMs offer a balance*  
006 *of efficiency and performance. Aggressive quantization*  
007 *can drastically cut model size and speed up inference, but*  
008 *its uneven effects on different components require care-*  
009 *ful management. In this work, we propose a lightweight,*  
010 *backpropagation-free, surrogate-based sensitivity analysis*  
011 *framework to identify hybrid SSM–Transformer compo-*  
012 *nents most susceptible to quantization-induced degrada-*  
013 *tion. Relying solely on forward-pass metrics, our method*  
014 *avoids expensive gradient computations and retraining,*  
015 *making it suitable for situations where access to in-domain*  
016 *data is limited due to proprietary restrictions or privacy*  
017 *constraints. We also provide a formal analysis showing that*  
018 *the Kullback–Leibler (KL) divergence metric better cap-*  
019 *tures quantization sensitivity for Language modeling tasks*  
020 *than widely adopted alternatives such as mean squared er-*  
021 *ror (MSE) and signal-to-quantization-noise ratio (SQNR).*  
022 *Through extensive experiments on SSM and hybrid archi-*  
023 *tectures, our ablation studies confirm that KL-based rank-*  
024 *ings align with observed performance drops and outper-*  
025 *form alternative metrics. This framework enables the prac-*  
026 *tical deployment of advanced hybrid models on resource-*  
027 *constrained edge devices with minimal accuracy loss. We*  
028 *further validate our approach with real-world on-device*  
029 *profiling on Intel Lunar Lake hardware, demonstrating that*  
030 *KL-guided mixed-precision achieves near-FP16 perplexity*  
031 *with model sizes and throughput competitive with Uniform*  
032 *INT4 on both CPU and GPU execution modes.*

## 033 1. Introduction

034 Large Language Models (LLMs) based on the Transformer  
035 architecture [17] have achieved remarkable success across

NLP tasks, but their deployment on edge devices is con- 036  
strained by high memory and computation demands. Re- 037  
cently, State Space Models (SSMs) [7] have emerged as 038  
efficient alternatives for long sequence modeling, offering 039  
 $\mathcal{O}(n)$  sequence processing and favorable memory scaling. 040  
SSM-based architectures can maintain strong accuracy on 041  
sequence tasks while being more hardware-friendly, which 042  
makes them attractive for resource-constrained platforms. 043  
In parallel, *hybrid models* [4, 6] combining SSM and Trans- 044  
former components have been proposed to capture the best 045  
of both worlds. These hybrids achieve competitive perfor- 046  
mance with improved efficiency, exemplified by models like 047  
Zamba [6] and Hymba [4] that interleave SSM and attention 048  
modules. 049

Despite their efficiency, even SSM-based hybrids at scale 050  
can contain hundreds of millions to billions of parameters, 051  
remaining impractical for direct deployment on mobile and 052  
IoT devices. Quantization [10, 13] is a proven approach 053  
to compress models by reducing the numerical precision 054  
of weights and activations. 4-bit uniform quantization of 055  
Transformers has become standard in industry for inference 056  
speed-ups, and recent studies demonstrate that even lower- 057  
bit weight quantization has become feasible for LLMs given 058  
proper handling of outlier features [2, 20, 22]. However, 059  
applying uniform quantization to all parts of a hybrid SSM 060  
model can lead to disproportionate accuracy degradation in 061  
certain layers. Prior work [15] notes that SSMs exhibit 062  
different quantization behavior than Transformers: for in- 063  
stance, the state transition in selective SSMs has highly sen- 064  
sitive internal activations, and SSM outputs contain large 065  
magnitude outliers not seen in Transformer outputs. These 066  
differences call for a more *nuanced*, layer-wise approach to 067  
quantization. 068

In this paper, we demonstrate that quantizing SSM– 069  
Transformer models for edge deployment is both necessary 070  
and tractable. While prior mixed-precision techniques [5, 071  
14] have shown success in CNNs and Transformers by us- 072  
ing gradient-based methods or signal-to-quantization-noise 073  
ratio (SQNR) for sensitivity analysis, these approaches fall 074  
short in the context of language modeling. In particu- 075

lar, we find that SQNR—though effective for convolutional architectures—fails to capture the true sensitivity of components in language models. Our key contribution is a novel, gradient-free sensitivity analysis method tailored for SSM–Transformer architectures. It operates entirely via forward-pass signals and reveals which layers truly require higher precision. This insight enables a post-training quantization pipeline that achieves substantial compression with minimal accuracy degradation, addressing a gap where existing mixed-precision strategies break down.

The remainder of the paper is organized as follows. Section 2 reviews prior work on quantization for LLMs, SSMs, and hybrid architectures. Section 3 presents our forward-pass sensitivity analysis, including per-layer metrics, mixed-precision assignment, and Kendall’s  $\tau$  ranking. Section 4.1 formally shows that KL divergence is a tighter proxy for perplexity than SQNR. Section 5.1 reports ablations on the HYMBA hybrid model at the layer and parameter level. Section 4.2 validates layer-level correlations across Mamba and hybrid architectures, confirming that  $KL_{\text{student} \rightarrow \text{teacher}}$  outperforms SQNR and other metrics. Finally, Section 7 demonstrates end-to-end deployment on Intel Lunar Lake hardware, achieving near-FP16 perplexity with latency and throughput competitive with Uniform INT4 on CPU and GPU.

## 2. Related Work

Prior work on efficient sequence modeling can be viewed along two largely orthogonal axes: **model architecture** and **quantization strategy**. Along the first axis, recent systems include Transformer-based large language models (LLMs), state-space models (SSMs), and hybrids that combine the two. Along the second axis, quantization ranges from **homogeneous** schemes that use a single bit-width throughout the model to **mixed-precision** schemes that assign different precisions to different components based on sensitivity. Together, these two axes provide a useful way to understand how prior work trades off efficiency, accuracy, and hardware friendliness.

Early quantization work focused mainly on *homogeneous quantization*, especially for Transformer models. Uniform low-bit quantization is useful because it maps well to integer hardware and simplifies deployment. Early work in vision showed that carefully calibrated 8-bit quantization can preserve accuracy [10], and later work adapted similar ideas to LLMs using outlier-aware calibration and improved rounding methods [13]. However, large Transformers show strong layerwise variation in sensitivity, with some activations and weights exhibiting strong outlier behavior. As a result, forcing all layers to use the same precision can lead to noticeable degradation unless extra techniques such as per-layer clipping are used [21]. This motivates going beyond one-size-fits-all precision assignment.

Mixed-precision quantization addresses this issue by allocating bit-widths non-uniformly across the network. The main idea is to keep sensitive parts at higher precision while quantizing more robust parts more aggressively. HAWQ is an early example, using Hessian-based sensitivity estimates to assign layerwise precisions and showing strong compression–accuracy trade-offs in CNNs [5]. Similar ideas were later extended to LLMs. For example, Pandey et al. [14] propose a post-training method that uses a calibration set to allocate 4/8-bit weights across LLM layers, outperforming uniform int8 baselines on GLUE [18] and SQuAD [16]. Related work also explores low-bit fine-tuning [3], adaptive post-training rounding [12], and hardware-aware bit allocation [19]. Overall, these works show that mixed precision often gives a better accuracy–efficiency trade-off than homogeneous quantization for LLMs.

Quantization research has also started to extend beyond Transformers to *state-space models*. SSMs offer a different efficiency profile, with linear-time sequence processing and constant-memory recurrent state updates, making them attractive for long-context and edge settings. Mamba shows that modern SSMs can achieve competitive sequence modeling performance at lower compute cost than Transformers on several long-range tasks [7]. Quantization studies for these models are still limited, but existing results already suggest different sensitivity patterns from Transformers. Pierro and Abreu [15] provide one of the first systematic studies of post-training quantization for Mamba-style recurrent LLMs, showing that uniform int8 quantization can retain strong performance when sensitive state-transition components are treated carefully, in some cases with quantization-aware training. At the same time, pushing all components to int4 leads to much larger degradation, suggesting that homogeneous quantization is mainly viable for SSMs at moderate precision.

These observations naturally motivate *hybrid architectures* that combine SSM and Transformer components, where both architectural heterogeneity and precision heterogeneity become important. Recent models such as Hymba and Zamba interleave recurrent state-space computation with attention and feed-forward blocks to capture the efficiency advantages of SSMs while retaining Transformer flexibility [4, 6]. Such models are especially well suited to mixed-precision deployment because their sub-modules play different computational roles. In practice, recurrent state-transition matrices may need higher precision to preserve long-range dynamics, while feed-forward and attention-related components can often be quantized more aggressively. This makes hybrid models a natural setting for sensitivity-guided precision allocation.

Overall, prior work points to two consistent themes. First, quantization sensitivity is highly non-uniform, mak-

181 ing homogeneous precision increasingly suboptimal as  
182 models become larger and more diverse. Second, this  
183 non-uniformity becomes even stronger in hybrid SSM–  
184 Transformer architectures, where recurrent and attention-  
185 based components have different numerical requirements.  
186 Our work builds on this intersection and studies how  
187 sensitivity-guided precision assignment can better exploit  
188 the structure of hybrid sequence models.

### 189 3. Methodology

190 We propose a lightweight sensitivity analysis method that  
191 evaluates each layer or module of a model for quantization  
192 robustness *without requiring backpropagation*. In contrast,  
193 our method uses forward-pass computations only, making  
194 it scalable to very large models and avoiding any need for  
195 retraining or fine-tuning during the analysis phase.

196 Our core idea is to measure how much quantization of a  
197 single layer affects the model’s output or performance met-  
198 ric. Given a trained model, we perform a series of modified  
199 inference passes over a representative dataset:

---

#### Algorithm 1 KL-PPL CORRELATION VIA KENDALL $\tau$

---

**Require:** Model  $\mathcal{M}$  with  $L$  layers; quantization function  
QUANTIZE( $\cdot$ ); evaluation dataset  $\mathcal{D}$

**Ensure:** Kendall  $\tau$  and  $p$ -value between KL and PPL rank-  
ings

```

1: Initialize lists: KL  $\leftarrow []$ , PPL  $\leftarrow []$ 
2: for  $\ell = 1$  to  $L$  do
3:   Teacher  $\leftarrow \mathcal{M}$ 
4:   Student  $\leftarrow$  QUANTIZE( $\mathcal{M}, \ell$ )
5:   logits $_T$ , PPL $_T \leftarrow$  EVALUATE(Teacher,  $\mathcal{D}$ )
6:   logits $_S$ , PPL $_S \leftarrow$  EVALUATE(Student,  $\mathcal{D}$ )
7:   KL $_{\ell} \leftarrow$  COMPUTEKL(logits $_T$ , logits $_S$ )
8:   Append KL $_{\ell}$  to KL; append PPL $_S$  to PPL
9: end for
10:  $\pi_{\text{KL}} \leftarrow$  ARGSORT(KL,  $\downarrow$ )
11:  $\pi_{\text{PPL}} \leftarrow$  ARGSORT(PPL,  $\downarrow$ )
12:  $(\tau, p) \leftarrow$  KENDALLTAU( $\pi_{\text{KL}}, \pi_{\text{PPL}}$ )
13: return  $(\tau, p)$ 

```

---

Table 1. Sort direction for each metric ( $\uparrow$ : ascending,  $\downarrow$ : descend-  
ing).

Metric	Order
Perplexity	$\downarrow$
SQNR (dB)	$\uparrow$
KL $_{\text{teacher} \rightarrow \text{student}}$	$\downarrow$
KL $_{\text{student} \rightarrow \text{teacher}}$	$\downarrow$
$\Delta$ Cross Entropy	$\downarrow$

200 This requires no gradients and is easily parallelizable.

### 3.1. Efficient Mixed-Precision Assignment 201

202 Given sensitivity scores, we assign a low bit-width to most  
203 layers and a higher bit-width to the top- $k$  sensitive ones.  
204 This yields significant compression with minimal accuracy  
205 loss. Given model output and input logits, we perform a  
206 comprehensive analysis to discover a correlation between  
207 perplexity and metrics discussed in Sec. 4.1.

### 3.2. Overview of Kendall’s $\tau$ Ranking 208

209 Kendall’s  $\tau$  [9] is a non-parametric statistical measure that  
210 quantifies the ordinal correlation between two rankings.  
211 Specifically, given two rankings  $R_1$  and  $R_2$  of  $n$  elements,  
212 Kendall’s  $\tau$  is defined as:

$$\tau(R_1, R_2) = \frac{|\mathcal{C}| - |\mathcal{D}|}{\frac{1}{2}n(n-1)}, \quad 213$$

214 where  $\mathcal{C}$  and  $\mathcal{D}$  represent the sets of concordantly and dis-  
215 concordantly ranked pairs, respectively. The coefficient ranges  
216 from  $-1$  (complete disagreement) to  $+1$  (complete agree-  
217 ment), with 0 indicating no correlation. Unlike Pearson’s  
218 correlation, Kendall’s  $\tau$  relies solely on rank order, making  
219 it particularly suitable for assessing ranking-based sensitiv-  
220 ity metrics.

## 4. Sensitivity Scoring in SSMs 221

### 4.1. Analytical Comparison of $D_{\text{KL}}$ and SQNR as Sensitivity Metrics for Mixed-Precision Lan- guage Models 222

223 Recent works [14, 21] establish the *signal-to-quantization-*  
224 *noise ratio* (SQNR) as a reliable sensitivity measure for  
225 CNNs and transformer-based models. SQNR quantifies the  
226 distortion introduced by quantizing model logits compared to  
227 the original (teacher) logits, expressed in decibels (dB):  
228

$$\text{SQNR} = 10 \log_{10} \frac{\mathbb{E}[\ell_{\text{orig}}^2]}{\mathbb{E}[(\ell_{\text{orig}} - \ell_{\text{quant}})^2]}. \quad 229$$

230 Higher SQNR values signify that quantization retains the  
231 original activations more accurately.  
232

233 However, for autoregressive language models, we  
234 demonstrate that the *Kullback–Leibler divergence* [11]  
235 ( $D_{\text{KL}}$ ) generalizes better than SQNR. KL divergence mea-  
236 sures distributional shifts due to quantization in both direc-  
237 tions to account for potential asymmetries:

$$\text{KL}(P||Q) = \mathbb{E}_P[\log P - \log Q], \quad 238$$

239 where  $P = \text{softmax}(\ell_{\text{orig}})$  and  $Q = \text{softmax}(\ell_{\text{quant}})$ .  
240 Higher KL values indicate greater divergence between the  
241 original and quantized distributions.

242 While SQNR captures signal fidelity at the logit level,  
243 it does not reliably correspond to downstream performance  
244 in language tasks. In contrast, KL divergence provides a  
245 tighter link to the change in perplexity ( $\downarrow$ ), making it more  
246 suitable for sensitivity scoring in autoregressive settings.

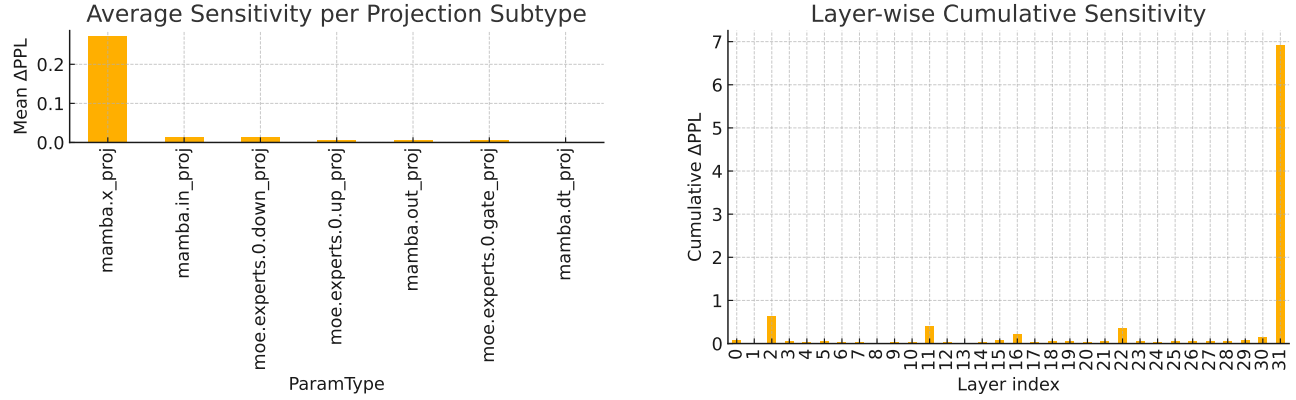


Figure 1. Quantization-sensitivity characteristics of Hymba. (a) Average ablation sensitivity  $\bar{\Delta}\text{PPL}$  for each projection subtype, with `mamba.x_proj` dominating the profile. (b) Layer-wise cumulative sensitivity  $\Sigma\Delta\text{PPL}$ , showing a sharp peak at block 31.

247 **Notation.** Let  $x$  be a context and  $y \in V$  a token. Through-  
248 out,

$$249 \quad \begin{aligned} p(y | x) &: \text{full-precision } \textit{teacher} \text{ model,} \\ q_\theta(y | x) &: \text{quantized model.} \end{aligned}$$

250 Expectations  $\mathbb{E}_p[\cdot]$  are taken over  $(x, y) \sim p$ . As all expecta-  
251 tions  $\mathbb{E}_p[\cdot]$  are taken with respect to the teacher, the results  
252 are *relative* to that uncompressed model.<sup>1</sup>

253 [Cross-entropy split]

$$254 \quad (q_\theta, p) = -\mathbb{E}_p[\log q_\theta(y | x)] \quad (1)$$

$$255 \quad = H(p) + D_{\text{KL}}(p \| q_\theta). \quad (2)$$

256 [PPL factorization]

$$257 \quad (q_\theta) = \exp((q_\theta, p)) \quad (3)$$

$$258 \quad = (p) \exp(D_{\text{KL}}(p \| q_\theta)), \quad (4)$$

259 where  $(p) = \exp(H(p))$ .

260 [Teacher-relative PPL bound] If  $D_{\text{KL}}(p \| q_\theta) \leq \varepsilon$  then

$$261 \quad (q_\theta) \leq (p) e^\varepsilon.$$

262 [SQNR Is Not Monotonic in ] There exist logit pairs  
263  $(z^{(1)}, \hat{z}^{(1)})$  and  $(z^{(2)}, \hat{z}^{(2)})$  such that

$$264 \quad \text{SQNR}(z^{(1)}, \hat{z}^{(1)}) < \text{SQNR}(z^{(2)}, \hat{z}^{(2)}), \quad (5)$$

$$265 \quad (1) < (2). \quad (6)$$

266 *Constant-shift example:* Let  $\hat{z} = z + c\mathbf{1}$  with any  $c \neq 0$ .

<sup>1</sup>If we use the test-set distribution for  $p$ , the formulas stay the same. The only change is that  $H(p)$  now depends on that specific test data.

Since  $(z) = (\hat{z})$ , the models share the same , but

$$267 \quad \text{SQNR}(z, \hat{z}) = 10 \log_{10} \frac{\|z\|_2^2}{\|z - \hat{z}\|_2^2} \quad (7) \quad 268$$

$$269 \quad = 10 \log_{10} \frac{\|z\|_2^2}{c^2 \|\mathbf{1}\|_2^2} \quad (8)$$

$$270 \quad \xrightarrow{|c| \rightarrow \infty} -\infty. \quad (9)$$

271 So SQNR can be made arbitrarily small while is un-  
272 changed.

## 4.2. Metric Asymmetry

273 We observe a strong asymmetry between the two KL  
274 directions used to compare student and teacher distri-  
275 butions. While  $KL_{\text{student} \rightarrow \text{teacher}}$  aligns well with the  
276 perplexity-based sensitivity ranking, the reverse direction  
277  $KL_{\text{teacher} \rightarrow \text{student}}$  shows a weak negative correlation ( $\tau \approx$   
278  $-0.14$ ).

279 This difference arises because the two directions pe-  
280 nalize different types of errors.  $KL_{\text{student} \rightarrow \text{teacher}}$  strongly  
281 penalizes probability mass that the student assigns to to-  
282 kens the teacher considers unlikely—precisely the type  
283 of distortion introduced by quantization. In contrast,  
284  $KL_{\text{teacher} \rightarrow \text{student}}$  focuses on whether the student underesti-  
285 mates the teacher’s high-probability tokens, which tends to  
286 remain relatively stable under quantization. As a result, the  
287 reverse KL is less sensitive to the perturbations that drive  
288 perplexity degradation. 289

290 **Practical Implication.** As shown above, the *KL diver-*  
291 *gence gap*  $\Delta D_{\text{KL}}(p \| q_\theta)$  provides an upper bound proxy  
292 for perplexity, making it the preferred layer-sensitivity  
293 metric for language modelling tasks. Hence, in mixed-  
294 precision quantization of language models, precision allo-  
295 cation should be driven by  $\Delta D_{\text{KL}}$ .

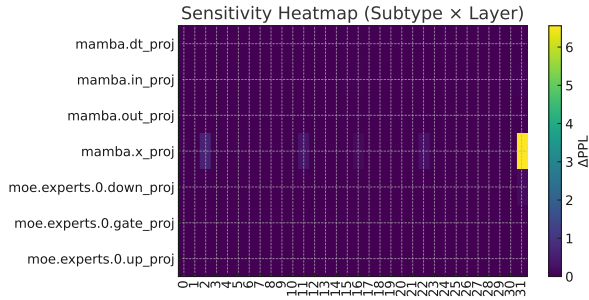


Figure 2. Sensitivity heat-map ( $\Delta$ Perplexity; rows: projection subtype, columns: layer). Late-stage `x_proj` peaks coincide with moderate MoE up/down sensitivities in the same block, underscoring their coupled importance for accurate mixed-precision deployment.

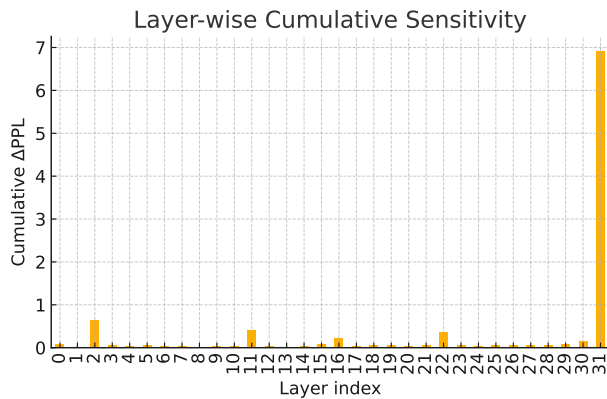


Figure 3. Layer-wise cumulative sensitivity  $\Sigma\Delta$ PPL, showing a sharp peak at block 31.

## 296 5. Ablations Studies on Hymba Hybrid Model

### 297 5.1. Layer- and Parameter-Level Quantization Sensitivity in HYMBA

299 We present a systematic, tensor-level ablation study of the  
300 32-layer HYMBA model [4]. Although similar analyses  
301 are common for CNNs and pure-Transformer models, we  
302 are not aware of a prior, publicly documented examination  
303 of quantization sensitivity in hybrid SSM-Transformer archi-  
304 tectures; our findings therefore offer a reference point  
305 for subsequent research on this emerging model class. All  
306 experiments are conducted under a **uniform 4-bit (INT4)**  
307 **quantization scheme**, a setting that is now widely sup-  
308 ported on commercial mobile and server-class accelerators.

309 **Subtype imbalance.** Among the seven projection sub-  
310 modules, the cross-projection `mamba.x_proj` exhibits the  
311 highest average sensitivity ( $\bar{\Delta} = 0.27$ ), exceeding all other  
312 subtypes by more than an order of magnitude ( $\bar{\Delta} \leq 0.014$ ).  
313 Conversely, `mamba.dt_proj` is nearly insensitive ( $\bar{\Delta} =$

$3.6 \times 10^{-4}$ ), suggesting that extremely low-precision rep-  
314 resentations are feasible (Figure 2). 315

**Localized layer hotspots.** Aggregating  $\Delta$ Perplexity  
316 across sub-modules reveals that block 31 alone accounts for  
317 over 70% of the aggregate sensitivity budget ( $\Sigma\Delta = 6.9$ ),  
318 with smaller peaks in blocks 2, 11, 16, and 22 (Figure 3).  
319 Most remaining blocks tolerate aggressive quantization. 320

**Coupled `x_proj`-MoE behaviour.** The heat-map in Fig-  
321 ure 2 shows that late-stage `x_proj` sensitivities co-occur  
322 with moderate MoE up/down sensitivities in the same  
323 block, indicating that quantization noise in mixing path-  
324 ways can be amplified by expert routing. 325

## 5.2. Ranking Correlation Protocol 326

We evaluate how effectively each sensitivity metric predicts  
327 the quantization impact at a per-layer granularity, defined  
328 by the increase in perplexity ( $\uparrow$ ) when that layer alone is  
329 quantized to `int8`. Following prior practice [14], we es-  
330 tablish the ground-truth ranking based on  $\Delta$  and measure  
331 the alignment with proxy metrics—namely, SQNR, KL di-  
332 vergences (teacher-to-student and student-to-teacher), and  
333  $\Delta$ CCE—using Kendall’s  $\tau$  coefficient (see Sec. 3.2). Met-  
334 rics with higher  $\tau$  values better reflect the true sensitivity of  
335 layers. 336

## 6. Experiments & Results 337

### 6.1. Layer-Level Correlation Analysis 338

Tab. 2 summarizes the Kendall’s  $\tau$  results across several hy-  
339 brid SSM models of varying scales. Notably, the student-  
340 to-teacher KL divergence consistently achieves the highest  
341 correlation, averaging a Kendall’s  $\tau$  of **0.79**, surpassing the  
342 widely used SQNR metric (average  $\tau = 0.76$ ). Statistical  
343 significance tests (paired one-sided t-tests) further vali-  
344 date that KL divergence significantly outperforms SQNR  
345 ( $p < 10^{-6}$ ). This improved correlation is explained by KL  
346 divergence operating directly in probability space, aligning  
347 closely with the downstream metric of interest (perplexity),  
348 whereas SQNR measures only logit-level fidelity and can  
349 overlook distributional shifts (Sec. 4.1). 350

### 6.2. Ranking-Correlation Algorithm 351

We detail the Kendall’s  $\tau$  ranking-correlation algorithm  
352 used to measure the alignment between our sensitivity met-  
353 rics and the observed perplexity ( $\uparrow$ ) change. 354

**Ranking Protocol.** We first sort layers by the observed  $\Delta$   
355 as the ground truth ranking. We then compare this ranking  
356 against rankings obtained from each sensitivity metric. 357

Table 2. Kendall’s  $\tau$  values for each metric across models.

Metric	Mamba-130M	Mamba-380M	Mamba-1.4B	Mamba2-130M	Hymba	Zamba	Avg.
SQNR (dB)	0.6967	0.7314	0.7941	0.4898	0.7124	0.8419	0.7111
$KL_{teacher \rightarrow student}$	-0.1980	-0.2316	-0.2607	-0.0272	-0.1617	0.1142	-0.1275
$KL_{student \rightarrow teacher}$	<b>0.8419</b>	<b>0.7936</b>	<b>0.8327</b>	<b>0.8078</b>	<b>0.6646</b>	<b>0.8060</b>	<b>0.7911</b>
$\Delta CE$	-0.1400	-0.1662	-0.0679	0.2364	-0.1260	-0.1234	-0.0645

Table 3. p-values for each metric across models.

Metric	Mamba-130M	Mamba-380M	Mamba-1.4B	Mamba2-130M	Hymba	Zamba
SQNR (dB)	$4.97 \times 10^{-24}$	$1.5 \times 10^{-51}$	$1.91 \times 10^{-60}$	$6.869 \times 10^{-7}$	$5.87 \times 10^{-57}$	$4.57 \times 10^{-52}$
$KL_{teacher \rightarrow student}$	$4.06 \times 10^{-3}$	$1.729 \times 10^{-6}$	$7.30 \times 10^{-8}$	$7.827 \times 10^{-1}$	$3.07 \times 10^{-4}$	$3.95 \times 10^{-2}$
$KL_{student \rightarrow teacher}$	$2.51 \times 10^{-34}$	$2.29 \times 10^{-60}$	$2.84 \times 10^{-66}$	$2.637 \times 10^{-16}$	$8.31 \times 10^{-50}$	$7.1 \times 10^{-48}$
$\Delta CE$	$4.21 \times 10^{-2}$	$5.97 \times 10^{-4}$	$1.61 \times 10^{-1}$	$1.656 \times 10^{-2}$	$4.92 \times 10^{-3}$	$2.61 \times 10^{-2}$

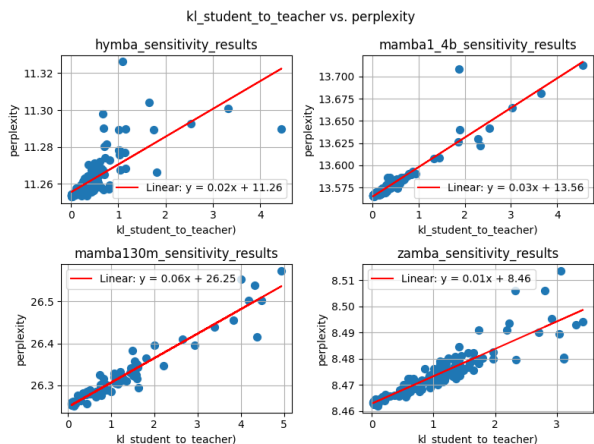


Figure 4. KL-Divergence vs. Perplexity correlation.

In this paper, we demonstrate that quantizing SSM-Transformer models for edge deployment is both necessary and tractable. While prior mixed-precision techniques [5, 14] have shown success in CNNs and Transformers using gradient-based sensitivity analysis or signal-to-quantization-noise ratio (SQNR), these approaches fall short in the context of language modelling. In particular, SQNR—though effective for convolutional architectures—fails to capture the true sensitivity of components in autoregressive LLMs, as we formally demonstrate in Sec. 4.1. Our key contribution is a novel, gradient-free sensitivity analysis method tailored for SSM-Transformer architectures that operates entirely via forward-pass signals. We validate this approach end-to-end with real-world profiling on **Intel Lunar Lake** hardware [8], covering CPU and GPU execution across two Mamba model scales, and show that KL-guided mixed-precision achieves near-FP16/INT8 perplexity while matching or exceeding INT4 throughput—a combination

that homogeneous quantization cannot provide.

In this paper, we demonstrate that quantizing SSM-Transformer models for edge deployment is both necessary and tractable. While prior mixed-precision techniques [5, 14] have shown success in CNNs and Transformers using gradient-based sensitivity analysis or signal-to-quantization-noise ratio (SQNR), these approaches are generally less aligned with the linear projection layers within the SSM block than KL divergence.

In particular, SQNR—though effective for convolutional architectures—often lags behind KL-divergence in capturing the true sensitivity of components in autoregressive LLMs, as we formally demonstrate in Sec. 4.1. Our key contribution is a novel, gradient-free sensitivity analysis method tailored for SSM-Transformer architectures that operates entirely via forward-pass signals. We validate this approach end-to-end with real-world profiling on **Intel Lunar Lake** hardware [8], covering CPU and GPU execution across two Mamba model scales, and show that KL-guided mixed-precision achieves near-FP16/INT8 perplexity while matching or exceeding INT4 throughput, a combination that homogeneous quantization cannot provide.

## 7. On-Device Profiling of Mixed-Precision Quantization

### 7.1. Experimental Setup

We profile models end-to-end on an **Intel Lunar Lake** platform [8], which integrates a CPU, dedicated GPU, and Neural Processing Unit (NPU) on the same die, making it representative of next-generation AI-capable edge clients. Models are converted to OpenVINO IR via the XAMBA framework [1] and benchmarked with the OpenVINO `benchmark_app` (latency mode, 100 iterations). We evaluate **Mamba-130M** and **Mamba-1.4B** [7] on CPU, and **Mamba-2 130M** [7] on GPU, comparing against **FP16**,

410 **Uniform INT8**, and **Uniform INT4** baselines. We report  
411 model size (MB), latency (ms), throughput (FPS), and per-  
412 perplexity (WikiText-103).

413 **Mixed-precision configurations.** Each configuration  $p_k$   
414 ( $k = 1, \dots, 10$ ) corresponds to a distinct KL threshold  $\kappa_k$ ,  
415 above which layers are retained at FP16 and below which  
416 they are compressed to INT4 (CPU) or INT8 (GPU). For-  
417 mally, given per-layer KL sensitivity scores  $\{s_\ell\}_{\ell=1}^L$ , con-  
418 figuration  $p^k$  assigns:

$$419 \quad b_\ell = \begin{cases} \text{FP16} & \text{if } s_\ell \geq \kappa_k, \\ \text{INT4/INT8} & \text{otherwise,} \end{cases} \quad (10)$$

420 where  $\kappa_1 > \kappa_2 > \dots > \kappa_{10}$ . **p01** applies the most con-  
421 servative threshold, quantizing only the least sensitive lay-  
422 ers; **p10** applies the most aggressive threshold, approaching  
423 near-uniform quantization. The configurations thus trace a  
424 smooth compression–accuracy curve controlled entirely by  
425 the KL sensitivity scores.

## 426 7.2. CPU Profiling: Mamba-130M and Mamba- 427 1.4B

428 Figures 5 to 7 show perplexity, latency, and throughput plot-  
429 ted against model size for both CPU models. KL-guided  
430 quantization progressively compresses Mamba-130M from  
431 493 MB to 84 MB (5.9 $\times$ ) and Mamba-1.4B from 5.2 GB to  
432 723 MB (7.2 $\times$ ), matching Uniform INT4 at the most ag-  
433 gressive setting. Perplexity remains near the FP16 baseline  
434 through **p06** for Mamba-130M and **p09** for Mamba-1.4B;  
435 beyond these thresholds, a sharp sensitivity cliff emerges—  
436 consistent with the ablation finding that a single high-  
437 sensitivity layer can dominate the entire quantization error  
438 budget (Sec. 5.1). **p05** represents the best operating point  
439 for both models, achieving INT8-level or better throughput  
440 (35.9 vs. 35.5 FPS; 5.6 vs. 5.1 FPS) with near-FP16 per-  
441 vexity, and delivering the lowest observed latency for Mamba-  
442 1.4B (178 ms)—below both Uniform INT8 (196 ms) and  
443 INT4 (190 ms).

## 444 7.3. GPU Profiling: Mamba-2 130M

445 Figure 8 reports GPU results as two complementary views:  
446 perplexity vs. throughput (left) and perplexity vs. size.

447 **Why GPU memory savings are negligible.** Un-  
448 like the CPU models, KL-MP configurations cluster  
449 around  $\approx 45$  MB regardless of  $k$ , versus 81 MB for  
450 the FP16 baseline. This is a consequence of two  
451 pipeline constraints: (1) the XAMBA exporter applies  
452 `compress_to_fp16=True` by default, so the FP16  
453 baseline is already half the FP32 size before any further  
454 quantization; and (2) the GPU pipeline applies INT8 rather  
455 than INT4, yielding only a 2 $\times$  reduction over FP16 on

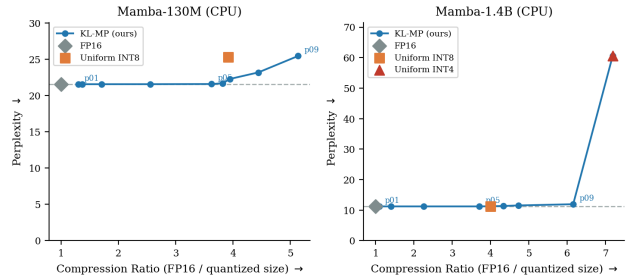


Figure 5. Perplexity vs. model size (CPU, Intel Lunar Lake). The x-axis is inverted so more compressed models appear to the right. KL-MP maintains near-FP16 accuracy across a wide compression range. The sharp cliff at **p10** for Mamba-130M (diverged) and **p10** for Mamba-1.4B (PPL = 60.6) reflects the quantization of a single high-sensitivity layer identified in our ablation (Sec. 5.1); Uniform INT4 crosses this boundary by construction.

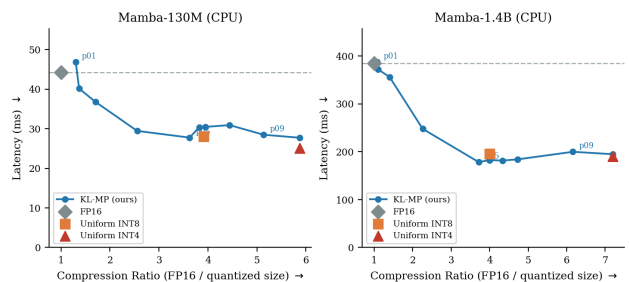


Figure 6. Inference latency vs. model size (CPU, Intel Lunar Lake). KL-MP **p05** achieves the lowest latency for Mamba-1.4B (178 ms), undercutting both Uniform INT8 and INT4, demonstrating that sensitivity-guided precision allocation reduces not only memory but also compute overhead.

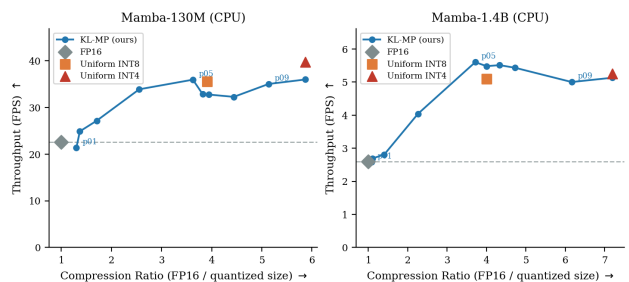


Figure 7. Throughput vs. model size (CPU, Intel Lunar Lake). KL-MP **p05** matches or exceeds INT4 throughput for both models while preserving near-FP16 perplexity, occupying a region of the accuracy–efficiency space inaccessible to homogeneous quantization.

the quantized projection layers, while the non-quantizable  
embedding table dominates the file size throughput. Con-  
sequently, efficiency gains on GPU manifest as **la-  
tency and throughput improvements** from faster INT8  
arithmetic—not memory reduction.

456  
457  
458  
459  
460

Table 4. On-device results: best KL-MP point vs. baselines on Intel Lunar Lake.  $\uparrow$  higher is better;  $\downarrow$  lower is better. “—” \* The 130M model exhibited substantially higher perplexity under uniform INT4 quantization.  $\dagger$ FPS and latency are not reported, as INT8 execution is not supported for the certain relevant operations on Intel Lunar Lake

Model	Config	PPL $\downarrow$	Size	FPS $\uparrow$	Latency $\downarrow$
Mamba-130M (CPU)	FP16	21.59	493 MB	22.6	44 ms
	INT8	25.32	126 MB	35.5	28 ms
	INT4*	3e35	84 MB	39.7	25 ms
	<b>KL-MP p05 (ours)</b>	<b>21.61</b>	<b>136 MB</b>	<b>35.9</b>	<b>28 ms</b>
Mamba-1.4B (CPU)	FP16	11.22	5.2 GB	2.6	384 ms
	INT8	11.25	1.3 GB	5.1	196 ms
	INT4	60.55	723 MB	5.3	190 ms
	<b>KL-MP p05</b>	<b>11.22</b>	<b>1.4 GB</b>	<b>5.6</b>	<b>178 ms</b>
Mamba-2 130M (GPU)	FP16	46.45	81 MB	0.02	60006 ms
	INT8 $\dagger$	53.03	125 MB	—	—
	<b>KL-MP p02</b>	<b>46.46</b>	<b>45 MB</b>	<b>0.29</b>	<b>3417 ms</b>

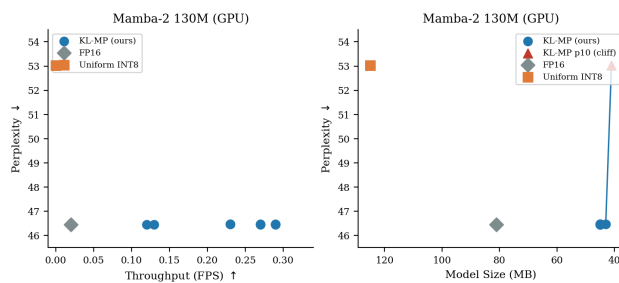


Figure 8. GPU profiling of Mamba-2 130M on Intel Lunar Lake. **(Left)** Perplexity vs. throughput: KL-MP configurations **p01–p09** achieve up to  $14.5\times$  throughput improvement over FP16 with negligible perplexity cost. **(Right)** Perplexity vs. model size: file sizes cluster near 45 MB due to FP16 export compression and INT8-only quantization; only **p10** (red triangle) crosses the precision cliff, collapsing to Uniform INT8 perplexity.

**Latency and accuracy.** The FP16 baseline incurs 60,006 ms due to iGPU memory-transfer overhead, while KL-MP reduces this by up to  $17.6\times$  (**p02**: 3,417 ms, 0.29 FPS). Perplexity holds at 46.45–46.48 for **p01–p09**, matching the FP16 baseline; only **p10** crosses the precision cliff to 53.03—identical to Uniform INT8—confirming that the KL top- $k$  selection correctly identifies the critical layer boundary.

## 8. Conclusion

In this work, we introduced a practical and scalable sensitivity analysis framework tailored specifically for hybrid SSM–Transformer models, enabling efficient mixed-precision quantization without requiring backpropagation

or retraining. By systematically evaluating quantization sensitivity through forward-pass metrics, we identified critical layers where quantization-induced errors significantly impact model performance. Our results on language modeling tasks demonstrate that selectively maintaining higher precision in these sensitive layers effectively mitigates accuracy degradation, achieving up to  $4\times$  reduction in model size with minimal perplexity loss. Future directions include extending our analysis framework to broader hybrid architectures and exploring dynamic quantization strategies that adapt precision allocation during inference, further enhancing model efficiency on resource-constrained edge devices. Real-world profiling on Intel Lunar Lake confirms these findings end-to-end: KL-guided mixed-precision reduces Mamba-1.4B from 5.2 GB to 1.4 GB with no measurable perplexity loss, and cuts Mamba-2 130M GPU latency by up to  $17.6\times$  over the FP16 baseline.

## References

- 491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546
- [1] XAMBA: Enabling efficient state space models on resource-constrained neural processing units. Software repository, 2024. 6
- [2] Tim Dettmers, Mike Lewis, Sam Shleifer, and Luke Zettlemoyer. LLM.int8(): 8-bit matrix multiplication for transformers at scale. *arXiv preprint arXiv:2208.07339*, 2022. 1
- [3] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. QLoRA: Efficient finetuning of quantized LLMs. *arXiv preprint arXiv:2305.14314*, 2023. 2
- [4] Xin Dong, Yonggan Fu, Shizhe Diao, Wonmin Byeon, Zijia Chen, Ameya Sunil Mahabaleshwar, Shih-Yang Liu, Matthijs Van Keirsbilck, Min-Hung Chen, Yoshi Suhara, et al. Hymba: A hybrid-head architecture for small language models. *arXiv preprint arXiv:2411.13676*, 2024. 1, 2, 5
- [5] Zhen Dong, Zhewei Yao, Amir Gholami, Michael W. Mahoney, and Kurt Keutzer. Hawq: Hessian aware quantization of neural networks with mixed-precision. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 293–302, 2019. 1, 2, 6
- [6] Paolo Glorioso, Quentin Anthony, Yury Tokpanov, James Whittington, Jonathan Pilault, Adam Ibrahim, and Beren Millidge. Zamba: A compact 7b ssm hybrid model. *arXiv preprint arXiv:2405.16712*, 2024. 1, 2
- [7] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023. 1, 2, 6
- [8] Intel Corporation. Intel core ultra 200v series processors (lunar lake) product brief. <https://www.intel.com/content/www/us/en/products/docs/processors/core-ultra/core-ultra-200v-product-brief.html>, 2024. Accessed: 2025. 6
- [9] Maurice G. Kendall. A new measure of rank correlation. *Biometrika*, 30(1–2):81–93, 1938. 3
- [10] Raghuraman Krishnamoorthi. Quantizing deep convolutional networks for efficient inference: A whitepaper. *arXiv preprint arXiv:1806.08342*, 2018. 1, 2
- [11] Solomon Kullback and Richard A. Leibler. On information and sufficiency. *The Annals of Mathematical Statistics*, 22(1):79–86, 1951. 3
- [12] Markus Nagel, Rana Ali Amjad, Mart Van Baalen, Christos Louizos, and Tijmen Blankevoort. Up or down? adaptive rounding for post-training quantization. In *International Conference on Machine Learning (ICML)*, 2020. 2
- [13] Markus Nagel, Marios Fournarakis, Rana Ali Amjad, Yelysei Bondarenko, Mart Van Baalen, and Tijmen Blankevoort. A white paper on neural network quantization. *arXiv preprint arXiv:2106.08295*, 2021. 1, 2
- [14] Nilesh Prasad Pandey, Markus Nagel, Mart van Baalen, Yin Huang, Chirag Patel, and Tijmen Blankevoort. A practical mixed precision algorithm for post-training quantization. *arXiv preprint arXiv:2302.05397*, 2023. 1, 2, 3, 5, 6
- [15] Alessandro Pierro and Steven Abreu. Mamba-ptq: Outlier channels in recurrent large language models. *arXiv preprint arXiv:2407.12397*, 2024. 1, 2
- [16] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of EMNLP*, 2016. 2
- [17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. 1
- [18] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. GLUE: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*, 2018. 2
- [19] Kuan Wang, Zhijian Liu, Yujun Lin, Ji Lin, and Song Han. HAQ: Hardware-aware automated quantization. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [20] Guangxuan Xiao, Ji Lin, Mathieu Seznec, Julien Demouth, and Song Han. SmoothQuant: Accurate and efficient post-training quantization for large language models. *arXiv preprint arXiv:2211.10438*, 2023. 1
- [21] Yuewei Yang, Xiaoliang Dai, Jialiang Wang, Peizhao Zhang, and Hongbo Zhang. Efficient quantization strategies for latent diffusion models. *arXiv preprint arXiv:2312.05431*, 2023. 2, 3
- [22] Zhewei Yao, Yuxiong He, et al. ZeroQuant: Efficient and affordable post-training quantization for large-scale transformers. *arXiv preprint arXiv:2206.01861*, 2022. 1
- 547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572

**573 A. On-Device Profiling: Reproducibility Guide**

574 All profiling scripts target Intel Lunar Lake via OpenVINO.  
575 The pipeline runs in three stages.

**576 Step 1 | Convert and Quantize**

```
577 python convert.py # base models
578 python quantize_uniform.py # INT8 / INT4
579 python quantize_mixed.py # KL M-P (CPU)
580 python quantize_mixed_gpu.py # KL M-P (GPU)
581
```

583 The sensitivity metric is controlled by a constant inside  
584 each script. All output filenames automatically inherit the  
585 selected tag (e.g., mamba-130m-hf\_kl\_point05.  
586 xml):

```
587 SENSITIVITY_METRIC = "kl_student_to_teacher"
588 METRIC_TAG = "kl"
589
```

591 The ten configurations **p01–p10** are produced by split-  
592 ting the sensitivity-ranked layer list into equal segments.  
593 Configuration **p01** quantizes only the least sensitive layers,  
594 while **p10** approaches near-uniform quantization.

**595 Step 2 | Benchmark**

```
596 python benchmark.py # CPU latency +
597 throughput
598 python benchmark_gpu.py # GPU latency +
599 throughput
600
```

602 Each model is evaluated using the OpenVINO  
603 benchmark\_app tool:

```
604 benchmark_app -m <model>.xml \
605 -d CPU -hint latency \
606 -t 60 -niter 50 \
607 --inference_only TRUE
608
```

610 The flags `-t 60` and `-niter 50` ensure that each run  
611 executes for at least 60 seconds and at least 50 iterations,  
612 whichever takes longer.

613 Results are written to `log/benchmark_log/` as per-  
614 model text logs and a summary CSV file: `latency_  
615 throughput_{device}_{tag}_report.csv`.

```
616 python eval_perplexity.py # CPU, WikiText-2
617 PPL
618 python eval_perplexity_gpu.py # GPU, WikiText-2
619 PPL
620
```

622 Perplexity is computed on the WikiText-2 test set. Fake  
623 quantization is applied to replicate each configuration's  
624 weight precision during evaluation.

```
625 ds = load_dataset(
626     "wikitext", "wikitext-2-raw-v1",
627     split="test")
628
629 text = "\n\n".join(
630     t for t in ds["text"] if t.strip())
631
```

Results are stored in `perplexity_results_  
633 {tag}.json` with the structure  
634

```
635 {model: {point: perplexity}}
636
```

This format allows direct comparison across sensitivity  
638 metrics and model scales.  
639