

---

# Multiple Token Divergence: A Measure of In-Context Computation Density

---

**Vincent Herrmann**<sup>\*†</sup>

The Swiss AI Lab IDSIA/USI/SUPSI  
Lugano, Switzerland

**Eric Alcaide**<sup>\*</sup>

The Swiss AI Lab IDSIA/USI/SUPSI  
Lugano, Switzerland

**Jürgen Schmidhuber**

KAUST, The Swiss AI Lab IDSIA/USI/SUPSI  
Thuwal, Saudi-Arabia / Lugano, Switzerland

## Abstract

1 Measuring the in-context computational effort of language models is a key chal-  
2 lenge, as standard metrics like next-token loss fail to capture the complexity of the  
3 underlying reasoning. Prior work based on latent state compression is promising  
4 but can be invasive and unstable to train. In this paper, we propose Multiple Token  
5 Divergence (MTD), a simple and direct measure of computational effort that quan-  
6 tifies the KL divergence between the full model’s output distribution and that of a  
7 shallow, auxiliary prediction head. An MTD module can easily be inserted into a  
8 language model. Alternatively, pre-trained multiple token prediction heads that are  
9 included in some state-of-the-art models can be used directly, requiring no further  
10 training. We empirically show that MTD is more effective than prior methods  
11 at distinguishing complex tasks from simple ones. On mathematical reasoning  
12 benchmarks, we find that MTD correlates positively with problem difficulty—in  
13 direct contrast to next-token loss—and that lower MTD is associated with more  
14 accurate self-generated reasoning. MTD provides a practical, lightweight tool for  
15 analyzing and understanding the computational dynamics of language models.

## 16 1 Introduction

17 To solve unfamiliar and challenging reasoning problems, language models must perform sophisticated  
18 in-context computation. But how can we tell whether a model is making full use of its computational  
19 capacity at any given moment? It is well-established that the next-token prediction loss offers little  
20 insight [Schmidhuber, 1991a,b], as any particular reduction in loss can, in principle, be arbitrarily  
21 difficult to achieve [Bennett, 1988].

22 A more promising approach is to quantify meaningful computation by measuring the entropy, or  
23 incompressibility, of a model’s latent representations [Skean et al., 2025, Herrmann et al., 2025].  
24 This concept is rooted in the minimum description length (MDL) principle [Grünwald, 2007, Vitányi,  
25 2006, Elmoznino et al., 2024]: if the most compact description of a sequence’s structure is still  
26 complex, then predicting that sequence is computationally demanding. Based on this, the Prediction  
27 of Hidden States (PHi) loss was proposed as a measure of in-context computational complexity,  
28 quantifying the per-token information gain in a model’s latent space [Herrmann et al., 2025]. While  
29 promising, the PHi framework introduces significant practical challenges: it requires inserting a noisy

---

<sup>\*</sup>Equal Contribution

<sup>†</sup>Correspondence to [vincent.herrmann@idsia.ch](mailto:vincent.herrmann@idsia.ch)

information bottleneck that can degrade model performance, lead to unstable training, and is highly sensitive to its precise placement and the weighting of multiple loss terms.

In this work, we propose a simplified and more direct measure, Multiple Token Divergence (MTD), which quantifies information gain in the model’s output distribution. The core insight is simple: if a shallow computational shortcut (e.g., a single transformer block) can approximate the full model’s prediction, then the model is not performing particularly complex computation. If, however, there is a significant divergence between these two predictions, we can conclude that the model is leveraging its deeper computational capacity. MTD is straightforward to implement and can even be computed directly by leveraging the Multiple Token Prediction (MTP) modules that some modern pre-trained models already possess, requiring no additional fine-tuning.

In summary, our contributions are as follows:

- We introduce a conceptual distinction between information gain in a latent space (as in PHi) versus in the final output distribution.
- We propose giving auxiliary prediction modules access to the latest token embedding as a technique to isolate computationally non-trivial information gain.
- We present Multiple Token Divergence as a practical, effective, and less invasive measure of in-context computational effort.
- We empirically demonstrate that MTD outperforms PHi in distinguishing simple from complex tasks and investigate its properties in model-generated reasoning chains.

## 2 Background

This section outlines two related concepts for auxiliary prediction tasks in sequence models: predicting a model’s own future hidden states, as in the PHi framework, and predicting future output tokens, as in Multiple Token Prediction (MTP).

### 2.1 Prediction of Hidden States (PHi)

The Prediction of Hidden states (PHi) method, introduced by Herrmann et al. [2025], creates an information bottleneck within a sequence model to measure the complexity of its in-context computation. A PHi layer is inserted between a model’s “bottom” and “top” layers.

The model consists of the following modules: The **Bottom Layers** ( $B_\beta$ ) are the initial Transformer blocks that process the input sequence embeddings. The **PHi Layer** contains three key components: (1) An **encoder** ( $q_\psi$ ) that, at time step  $t$ , maps the hidden state  $g_t$  from the bottom layers to a posterior distribution over a latent variable  $z_t$ . This distribution is typically a diagonal Gaussian. (2) A **decoder** ( $a_\xi$ ) that reconstructs the hidden state, creating  $g'_t$  from a sample of the latent variable  $z_t$ . (3) An **autoregressive prior** that predicts the distribution of the current latent  $z_t$  using only the history of previous latents,  $z_{<t}$ . This can be implemented with a single Transformer block ( $M_\mu$ ) and two additional linear transforms,  $b_\kappa$  and  $p_\chi$ , which maps the output of the transformer to a prior distribution. The **Top Layers** ( $T_\tau$ ) are the remaining Transformer blocks that process the sequence of reconstructed hidden states  $g'$ . Finally, we have the standard token **Embedding** ( $E_\epsilon$ ) and the **Output Layer** ( $O_\omega$ ).

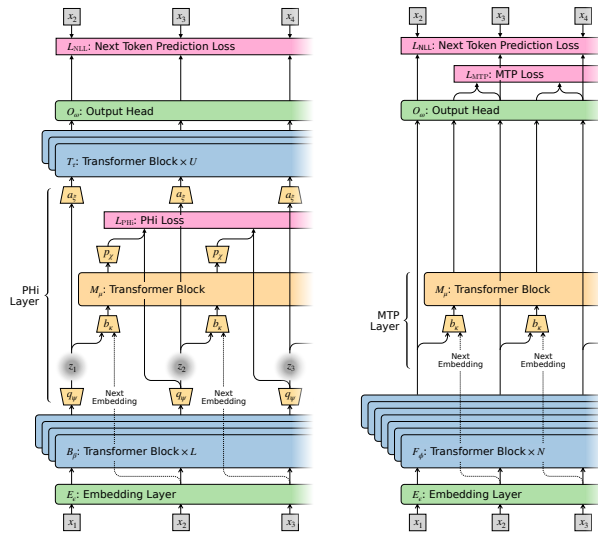


Figure 1: Comparison between the architecture of a PHi model (left) and of a MTP model (right).

81 Here and in the remainder of the paper, Greek subscript letter indicate learnable neural network  
82 parameters.

83 The forward pass processing tokens  $x_1, x_2, \dots$  is described by these equations:

$$\begin{aligned}
e_t &= E_\epsilon(x_t) && \text{Token embedding} \\
g_t &= B_\beta(e_1, \dots, e_t) && \text{Output from bottom layers} \\
z_t &\sim q_\psi(\cdot | g_t) && \text{Sample latent from the posterior} \\
g'_t &= a_\xi(z_t) && \text{Reconstruct the hidden state} \\
h_t &= T_\tau(g'_1, \dots, g'_t) && \text{Output from top layers} \\
\pi(\cdot | x_{<t}) &= O_\omega(h_{t-1}) && \text{Predict next token distribution} \\
L_{\text{NLL}}(t) &= -\log \pi(x_t | x_{<t}) && \text{Next-token prediction loss}
\end{aligned}$$

The Phi loss ( $L_{\text{Phi}}$ ) is the KL divergence between the posterior  $q_\psi$  (which has access to the current input  $x_t$  via  $g_t$ ) and the prior  $p_\chi$  (which only has access to past latents  $z_{<t}$ ). We assume an initial latent  $z_0$  is given.

$$c_t = b_\kappa(z_{t-1}) \quad \text{Linear projection of the last latent} \quad (1)$$

$$d_t = M_\mu(c_1, \dots, c_t) \quad \text{Output from Phi Transformer block}$$

$$\begin{aligned}
L_{\text{Phi}}(t) &= D_{\text{KL}}\left(q_\psi(\cdot | g_t) \parallel p_\chi(\cdot | d_t)\right) && \text{Phi Loss} \\
&= D_{\text{KL}}\left(q_\psi(\cdot | x_1, \dots, x_t) \parallel p_\chi(\cdot | z_0, z_1, \dots, z_{t-1})\right) && (2)
\end{aligned}$$

85 The model is trained to jointly minimize both  $L_{\text{NLL}}$  and  $L_{\text{Phi}}$ . The Phi loss quantifies the information  
86 gain at each timestep—the amount of new useful information present in the current input token  $x_t$   
87 that was not predictable from the history of latent states. It has been shown [Herrmann et al., 2025]  
88 that this value correlates well with the complexity and “interestingness” of tasks.

## 89 2.2 Multiple Token Prediction (MTP)

90 Multiple Token Prediction (MTP) is a technique used to improve model performance and enable  
91 faster inference via methods like speculative decoding [Cai et al., 2024, Gloeckle et al., 2024, Liu  
92 et al., 2024, Xiaomi et al., 2025]. In this setup, a computationally cheap auxiliary module is trained  
93 to directly predict the main model’s future output distribution.

94 First, consider a standard autoregressive model’s forward pass:

$$\begin{aligned}
e_t &= E_\epsilon(x_t) && \text{Token embedding} \\
h_t &= F_\phi(e_1, \dots, e_t) && \text{Output of all main Transformer blocks} \\
\pi(\cdot | x_{<t}) &= O_\omega(h_{t-1}) && \text{Predict next token distribution} \\
L_{\text{NLL}}(t) &= -\log \pi(x_t | x_{<t}) && \text{Next token prediction loss}
\end{aligned}$$

The goal of MTP is to approximate the main model’s prediction for the *next* token,  $x_{t+1}$ . Note that  
often in MTP, there are additional modules that approximate predictions for tokens even further in  
95 advance, i.e.,  $x_{t+n}$  for  $n > 1$ . We will not use them in this work.

A separate, smaller MTP module (e.g., a single Transformer block  $M_\mu$ ) generates its own prediction  
without access to the full model’s current hidden state  $h_t$  and is usually trained with a negative  
log-likelihood loss, which we call  $L_{\text{MTP}}$ . Optionally, the MTP module can be given access to the  
96 current token embedding  $e_t$ , as indicated by the square brackets.

$$\begin{aligned}
c_t &= b_\kappa(h_{t-1}, [e_t]) && \text{Input to MTP module (projection of } h_{t-1} \text{ and possibly } e_t) \\
d_t &= M_\mu(c_1, \dots, c_t) && \text{Output from MTP Transformer block} \\
\pi_{\text{MTP}}(\cdot | x_{<t}) &= O_\omega(d_{t-1}) && \text{MTP’s prediction for token } x_{t+1} \\
L_{\text{MTP}}(t) &= -\log \pi_{\text{MTP}}(x_t | x_{<t}) && \text{MTP Loss}
\end{aligned} \quad (3)$$

97 In order to predict the next token  $x + 1$ , the MTP module has access to the model’s history via  $h_{t-1}$   
98 (and optionally the current embedding  $e_t$ ), but it crucially lacks the result of the main model’s full  
99 computation at step  $t$  (i.e.,  $h_t$ ).

### 3 Multiple Token Divergence

Observe the similarity between the PHi framework’s autoregressive prior  $p_\chi$  and posterior  $q_\psi$  on the one hand, and the MTP prediction  $\pi_{\text{MTP}}$  and the full model prediction  $\pi$  on the other (Figure 1): In both cases, a computationally and informationally constrained module approximates the prediction from a full model. The key difference is that for PHi, this approximation occurs in a continuous latent space, whereas MTP operates directly on the discrete token distribution.

Based on this analogy, we propose the *Multiple Token Divergence* (MTD) as an alternative to the PHi loss. It is defined as the KL divergence between the full model’s next-token prediction  $\pi$  and the MTP module’s prediction  $\pi_{\text{MTP}}$ :

$$\begin{aligned} L_{\text{MTD}}(t) &= D_{\text{KL}}\left(\pi(\cdot | x_{\leq t}) \parallel \pi_{\text{MTP}}(\cdot | x_{< t})\right) \\ &= D_{\text{KL}}\left(\pi(\cdot | F_\phi(e_1, \dots, e_t)) \parallel \pi_{\text{MTP}}(\cdot | F_\phi(e_1, \dots, e_{t-1}), e_t)\right). \end{aligned} \quad (4)$$

The MTD loss,  $L_{\text{MTD}}$ , can either be optimized directly in conjunction with the standard next-token loss  $L_{\text{NLL}}$ , or it can be calculated post-hoc using an MTP module trained with  $L_{\text{MTP}}$  loss (see Section 4.2). For now, we ignore the optional access of the MTP module to the latest embedding  $e_t$ ; this will be addressed in Section 3.2.

#### 3.1 On the Difference between PHi and MTD

While PHi introduced a powerful conceptual framework for analyzing a model’s internal processing, its implementation can be complex. The stochasticity introduced by the variational information bottleneck can also interfere with the main sequence prediction task. In contrast, MTD is significantly simpler to implement as it functions as a non-invasive auxiliary task; providing a more direct and less disruptive method to obtain similar insights into the model’s per-token computational effort.

One interpretation of the PHi loss is that it measures changes in a particular point of the ‘latent program’ that the model synthesizes in-context to perform next-token prediction. The MTD loss, however, measures changes directly at the level of the output predictions. This distinction can lead to significant differences.

For instance, a small change in the latent program could result in a large shift in the output predictions. As an illustrative example, consider a model trained on two distinct types of sequences: one type consists of uniformly random tokens, while the other is a specific single, repeated token. To distinguish between these two cases, the latent program only needs to gain one bit of information. The PHi loss would therefore be low. However, the resulting change in the output distribution is large—shifting from a uniform distribution to a one-hot distribution. In this scenario, the MTD can be as high as  $D_{\text{KL}}(\text{one-hot} \parallel \text{uniform}) = \log_2(\text{vocabulary size})$  bits. In such cases, we expect  $L_{\text{MTD}}$  to be significantly higher than  $L_{\text{PHi}}$ , an effect we observe in our experiments in Section 4.1.1.

Conversely, one imagine cases where the latent program changes significantly while the output predictions remain stable. However, the PHi training objective penalizes encoding such changes, as they do not sufficiently improve downstream predictions and thus represent an inefficient use of the information bottleneck.

#### 3.2 Access to the Latest Token Embedding

A interesting nuance for both PHi loss and MTD is that the measured information gain at step  $t$  can originate from two distinct sources: (1) novel information contained within the current token  $x_t$  itself, and (2) complex computation performed by the model’s main layers ( $B_\beta$  for PHi,  $F_\phi$  for MTP), which cannot be easily approximated by the simpler prior or MTP module ( $M_\mu$ ).

To disentangle these two sources, we can provide the prior/MTP module with direct access to the latest token embedding,  $e_t$ . This is a common practice in MTP models (see Equation 3) and effectively isolates the second source of information gain. For the PHi framework, this modification involves updating Equation 1 to concatenate the previous latent state with the current embedding:

$$c_t = b_\kappa(z_{t-1}, e_t).$$

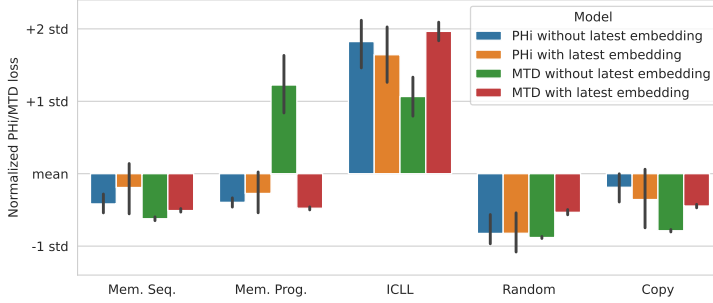


Figure 2: Normalized PHI or MTD loss of the four different model types on each of the five tasks. Only in-context language learning (ICLL) requires sophisticated in-context computation. This is reflected by the scores, with the exception of the MTD model without access to the latest embedding, which assigns high MTD also to the memorized programs task (see the discussion in Sections 3.1 and 4.1.1). Bootstrapped mean with 95% confidence intervals across 8 runs.

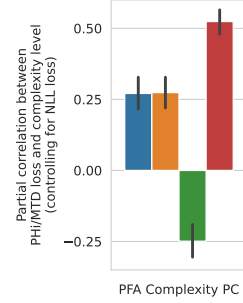


Figure 3: Partial correlation of PHI or MTD loss with the complexity of the modelled PFA, controlling for NLL. Also here, MTD without latest embedding access is the outlier.

With this change, the PHI prior has access to the same input token as the bottom layers,  $B_\beta$ . Consequently, the modified PHI loss, which we denote as  $\hat{L}_{\text{PHI}}$ , isolates the information gain attributable solely to the computation performed by  $B_\beta$ :

$$\hat{L}_{\text{PHI}}(t) = D_{\text{KL}}\left(q_\psi(\cdot | x_1, \dots, x_t) \parallel p_\chi(\cdot | z_0, \dots, z_{t-1}, e_t)\right).$$

A PHI layer modified in this way acts as an information bottleneck that specifically measures computational effort. Information that the prior can easily extract from the input embedding  $e_t$  is allowed to pass freely, while information that is computationally non-trivial for  $B_\beta$  to extract is quantified by  $\hat{L}_{\text{PHI}}$ . The same logic applies to the MTD module when it is given access to the latest embedding.

Arguably, PHI and MTD loss with access to the latest embedding provide a better measure of dense in-context computation. This access allows the prior/MTD module to account for trivial shifts in the predictions—like the one described in Section 3.1—thereby reducing the effective difference between the two metrics. In essence, providing access to the latest embedding allows us to quantify the information gain per step that is due to significant computational effort, whether measured in the latent space (PHI) or in the output distribution (MTD).

## 4 Experiments

We compare the different versions of PHI and MTD loss on transformer models trained from scratch on various tasks, before turning to pre-trained language models to further investigate the properties of MTD. For details on the exact training setup, please see Appendix A.

### 4.1 Training Sequence Models from Scratch

The considerations from Section 3 leave us with four different model configurations to compare: PHI and MTD models, each with and without access to the latest token embedding. The PHI model without this access corresponds to the original method proposed in prior work [Herrmann et al., 2025]. To compare these different setups, we train transformer models from scratch on several tasks and evaluate them in settings similar to those in [Herrmann et al., 2025].

#### 4.1.1 Evaluation on Different Tasks

The four model types are trained on five different tasks: (1) reciting memorized sequences, (2) modeling sequences from a small set of known formal languages (memorized programs), (3) in-context language learning (ICLL), where the formal language is unknown [Akyürek et al., 2024], (4)

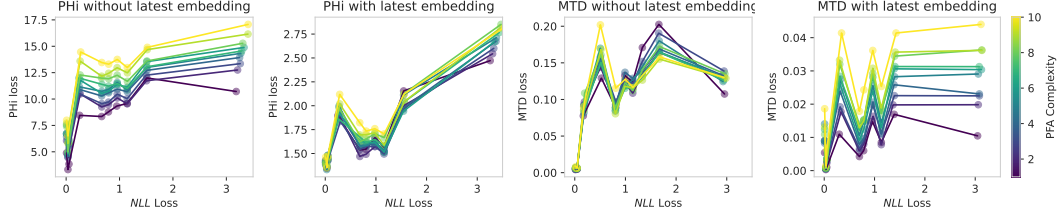


Figure 4: Token-wise PHi and MTD loss against binned NLL loss, for the different modeled PFA complexities. PHi loss without access to the latest embedding, and MTD loss with access to the latest embedding both show a clear correlation with complexity level, across NLL bins.

165 modeling random token sequences, and (5) a copying task that involves modeling random tokens  
 166 where subsequences appear twice. Of these, only ICLL—which requires inferring the structure of an  
 167 unknown probabilistic finite automaton (PFA) in-context—involves meaningful computation, in the  
 168 sense that a non-trivial latent program must be synthesized by the model.

169 Figure 2 shows a comparison of the normalized PHi and MTD losses for each task. The MTD with  
 170 latest embedding access shows the clearest distinction between the one complex task and the four  
 171 “boring” ones; note the high value for ICLL and the consistently low values for all other tasks. For  
 172 the MTD model *without* latest embedding access, we see the effect alluded to in Section 3.1: the  
 173 loss is high for both ICLL *and* the memorized programs. For the memorized programs task, the  
 174 actual in-context program required is minimal (only  $\sim \log_2(10)$  bits to identify any one out of the  
 175 ten memorized automata). However, the lack of information from the latest token causes a significant  
 176 shift in the model’s output distribution, resulting in a high MTD. Finally, giving the PHi layer access  
 177 to the latest embedding does not appear to improve its ability to distinguish boring from interesting  
 178 tasks.

#### 179 4.1.2 Task Complexity

180 Focusing on the ICLL task, we investigate the relationship between the models’ PHi or MTD  
 181 losses and the complexity of the underlying language, as measured by the description length of  
 182 the PFA. Figure 3 displays the partial correlation between the mean PHi or MTD loss across a  
 183 sequence and the language’s complexity. We control for the mean next-token prediction loss, as  
 184 it is positively correlated with language complexity ( $r = 0.367$ , 95% CI [0.315, 0.424]). Here  
 185 again, we find that MTD with latest embedding access shows the strongest positive correlation  
 186 ( $r = 0.524$ , [0.480, 0.565]). In contrast, MTD without access to the latest embedding is negatively  
 187 correlated with language complexity when controlling for NLL, confirming that it is not a reliable  
 188 measure for this purpose.

189 Figure 4 shows the token-wise PHi or MTD loss against binned NLL loss, broken down by PFA  
 190 complexity (from 1, simple, to 10, complex). This analysis reveals a more nuanced picture: only the  
 191 original PHi loss (without embedding access) and the MTD loss with embedding access show a clear,  
 192 positive token-wise relationship with language complexity after controlling for NLL.

#### 193 4.2 Pre-Trained Language Models

194 To validate our hypotheses on existing large-scale models, we leverage the pre-trained, open-source  
 195 MiMo-7B model [Xiaomi et al., 2025]. We chose this model for two reasons. First, as a modern,  
 196 high-quality 7B parameter model, its base pre-training incorporates an MTP objective, providing the  
 197 built-in auxiliary prediction heads necessary for calculating MTD without any post-hoc modification.  
 198 Second, its compact size, comparable to Llama 3 8B [Dubey et al., 2024a], allows for the efficient,  
 199 large-scale experimentation required to statistically validate our hypotheses across diverse tasks.

200 We note that while MiMo-7B was trained with an MTP objective from the outset, a similar setup  
 201 could be achieved for other models by keeping the base model frozen and training an MTP head  
 202 using standard teacher-student distillation [Schmidhuber, 1992, Hinton et al., 2015] with a fraction of  
 203 the original data and compute.

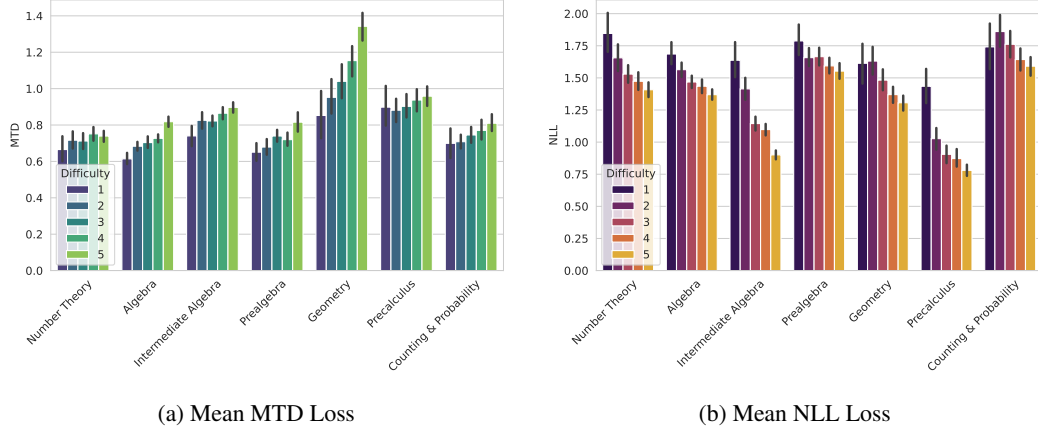


Figure 5: Mean losses of the MiMo model across the provided step-by-step solutions to the problems of the MATH test set, grouped by category and difficulty level. MTD clearly grows with difficulty, suggesting that the model is making more use of its computational capacity when processing more challenging problems. NLL loss, on the other hand, goes down with increasing complexity. Figure 7 shows similar results for self-generated chains of thought.

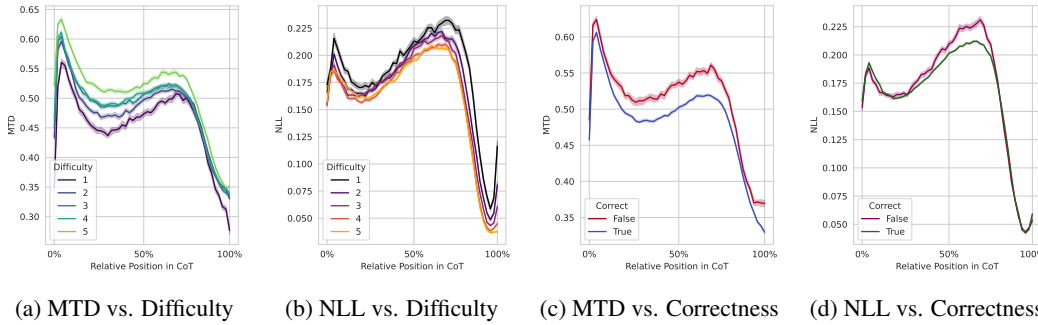


Figure 6: Token-wise losses against relative positions in self-generated CoT for the MATH test dataset. MTD shows a clear correlation with difficulty across the full CoT (a), the relationship between NLL and difficulty is less clear (b). Similarly, correct CoTs show higher MTD over all relative positions (c), which is not the case for NLL (d). Figure 10 (Appendix) shows similar results for the GSM-8k dataset.

#### 204 4.2.1 Reasoning Difficulty

205 We employ the MATH dataset [Hendrycks et al., 2021], which provides mathematics problems labeled  
 206 from Level 1 (easy) to Level 5 (hard), along with detailed reasoning solutions. We first compute the  
 207 mean MTD for the provided step-by-step solution for each problem in the dataset and find that it  
 208 clearly correlates with the difficulty level ( $r = 0.179$ , 95% CI [0.152, 0.203]). This correlation is  
 209 consistent across most problem categories, as shown in Figure 5a.

210 Interestingly, the NLL loss *negatively* correlates with problem difficulty ( $r =$   
 211  $-0.249$ ,  $[-0.274, -0.224]$ , Figure 5b). This suggests that from the model’s perspective,  
 212 reasoning chains for difficult problems are no less plausible or predictable. However, the higher  
 213 MTD indicates that the model must exert more computational effort to process and generate them.

214 We have the model generate ten different chains-of-thought (CoTs) for each problem and repeat  
 215 the analysis on these self-generated solutions. Also there, we observe very similar results (see  
 216 Figure 7 in the Appendix). The partial correlation between MTD and difficulty level, controlling  
 217 for NLL, is  $r = 0.199$ , [0.189, 0.208], while the correlation between NLL and difficulty is  $r =$   
 218  $-0.158$ ,  $[-0.168, -0.149]$ .

**Token-wise Development Across the Response** We also track the token-wise values for MTD and NLL across each generated CoT. As seen in Figure 6a, the positive correlation between MTD and problem difficulty holds consistently from the first tokens of the response to the last. Likewise, the negative correlation for NLL persists throughout the generation, even though not as pronounced (Figure 6b). The difference between MTD and NLL in their correlation with the problem difficulty is notable because, at a global level, MTD and NLL are positively correlated with each other ( $r = 0.255$ ,  $[0.246, 0.265]$ ). This highlights that MTD captures a distinct signal related to computational effort that is not present in the standard NLL loss.

#### 4.2.2 Reasoning Accuracy

For the self-generated CoTs, we also investigate the relationship between MTD values and the correctness of the final answer. Figure 6c plots the token-wise MTD, stratified by whether the rationale was correct or incorrect. We observe that correct responses are consistently associated with lower MTD. The relationship between NLL and correctness is less consistent (see Figure 6d).

Following the methodology from prior work [Herrmann et al., 2025], we randomly assemble pairs of one correct and one incorrect CoT for each math problem. The probability of choosing the correct CoT when picking the one with the lower mean MTD is 67.1% (95% CI:  $[65.4\%, 68.7\%]$ ). When selecting the one with the lower NLL, the probability is 73.3%  $[71.9\%, 74.8\%]$ .

We repeat these experiments on the GSM-8k dataset [Cobbe et al., 2021], where we find that selecting CoTs with lower MTD yields 66.0%  $[62.9\%, 69.2\%]$  correct answers, while lower NLL yields 72.2%  $[69.1\%, 75.0\%]$ . For token-wise curves, please see Appendix C.

These findings stand in contrast to the results for PHi loss, where, for a Llama 3B model, correct answers are associated with a *high* PHi loss [Herrmann et al., 2025]. While further investigation is needed, we hypothesize that different models may have different tendencies to either overly simplify or overly complicate their reasoning process. This tendency could determine whether computationally intensive answers—as opposed to more straightforward ones—are more or less likely to be correct for a given model architecture or training regime.

## 5 Discussion & Future Work

The experimental results establish MTD as an effective, granular measure of in-context computational effort—quantifying when a model is using its full processing capacity. In our experiments, it outperforms PHi loss in differentiating “boring” from “interesting” tasks and simple from complex ones. It successfully isolates the per-token information gain attributable to non-trivial, or “irreducible” [Wolfram, 2002], computation by the model. It should be noted that the MTD values depends on both the absolute and the relative capacities of the full model, and of the MTP module (especially  $M_\mu$ ). If the MTP module becomes similarly powerful to the full model, the MTD goes to zero. On the other hand, if the MTP module has too little capacity, MTD gives no information beyond the standard NLL loss.

While the empirical results are encouraging, several open questions remain, especially for models generating their own chains-of-thought. The fact that a CoT’s mean MTD correlates positively with problem complexity, while its mean NLL correlates negatively, is an interesting finding that warrants further investigation. Is this a unique property of the model we tested, or a general pattern across architectures and scales? Similarly, our finding that lower MTD is associated with more correct reasoning on two datasets is intriguing, especially as it contrasts with prior findings for PHi loss [Herrmann et al., 2025]. This relationship between computational effort and correctness is likely complex and model-dependent.

The MTD signal has the potential for many applications in training and inference. Examples could be **Dynamic Compute Allocation**: MTD could be monitored in real-time during generation. A prolonged period of low MTD might trigger early stopping for a simple task, while a sudden spike in MTD could activate more powerful components (e.g., additional Mixture-of-Experts layers) for a difficult step. **Solution Convergence**: The transition from a high-MTD processing phase to a low-MTD conclusion could act as a signal that the model has “settled” on a solution, potentially allowing for more efficient decoding. **Intrinsic Motivation**: In agent-based settings, MTD could serve as an intrinsic reward. This would encourage an agent to pursue policies that lead to computationally interesting states (high information gain), fostering the development of more sophisticated behaviors.

271 **Creative Decoding:** It could be explored whether decoding methods biased towards high-MTD  
 272 tokens might improve a model’s creativity and problem-solving abilities by encouraging it to venture  
 273 into less predictable computational paths. Conversely, could biasing towards easily predictable  
 274 low-MTD tokens lead to more robust reasoning?

## 275 6 Conclusion

276 In this work, we introduce Multiple Token Divergence, a practical and direct measure for quantifying  
 277 the computational effort of language models. We demonstrate that MTD, by measuring information  
 278 gain directly in the output distribution, serves as a more robust and stable metric than prior methods  
 279 that rely on latent state compression. Our findings show that MTD successfully distinguishes complex  
 280 in-context reasoning from simpler tasks and reveals a nuanced relationship between computational  
 281 effort and predictive loss. As a non-invasive and easily implemented metric, MTD provides a valuable  
 282 new tool for analyzing, evaluating, and ultimately steering the sophisticated reasoning processes of  
 283 modern sequence models.

## 284 References

- 285 Ekin Akyürek, Bailin Wang, Yoon Kim, and Jacob Andreas. In-context language learning: Architectures and  
 286 algorithms. In *Proc. Int. Conf. on Machine Learning (ICML)*, Vienna, Austria, July 2024.
- 287 C. H. Bennett. Logical depth and physical complexity. In *The Universal Turing Machine: A Half Century*  
 288 *Survey*, pages 227–258. Oxford University Press, 1988.
- 289 Tianle Cai, Yuhong Li, Zhengyang Geng, Hongwu Peng, Jason D. Lee, Deming Chen, and Tri Dao. Medusa:  
 290 Simple llm inference acceleration framework with multiple decoding heads. *arXiv preprint arXiv: 2401.10774*,  
 291 2024.
- 292 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert,  
 293 Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to  
 294 solve math word problems. *ArXiv*, abs/2110.14168, 2021.
- 295 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil  
 296 Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra,  
 297 Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien Rodriguez, Austen  
 298 Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux,  
 299 Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang  
 300 Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle  
 301 Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,  
 302 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip  
 303 Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Grégoire  
 304 Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan  
 305 Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan  
 306 Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny  
 307 Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton,  
 308 Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala,  
 309 Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and et al. The llama 3 herd of  
 310 models. *CoRR*, abs/2407.21783, 2024a. URL <https://doi.org/10.48550/arXiv.2407.21783>.
- 311 Abhimanyu Dubey et al. The Llama 3 herd of models. *ArXiv*, abs/2407.21783, 2024b.
- 312 Eric Elmoznino, Tom Marty, Tejas Kasetty, Léo Gagnon, Sarthak Mittal, Mahan Fathi, Dhanya Sridhar, and  
 313 Guillaume Lajoie. In-context learning and occam’s razor. *ArXiv*, abs/2410.14086, 2024.
- 314 Fabian Gloeckle, Badr Youbi Idrissi, Baptiste Roziere, David Lopez-Paz, and Gabriel Synnaeve. Better &  
 315 faster large language models via multi-token prediction. In Ruslan Salakhutdinov, Zico Kolter, Katherine  
 316 Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp, editors, *Proceedings of the 41st*  
 317 *International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*,  
 318 pages 15706–15734. PMLR, 21–27 Jul 2024.
- 319 Peter D. Grünwald. *The Minimum Description Length Principle*. Springer, 2007.
- 320 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob  
 321 Steinhardt. Measuring mathematical problem solving with the math dataset. *NeurIPS*, 2021.

322 Vincent Herrmann, Róbert Csordás, and Jürgen Schmidhuber. Measuring in-context computation complexity via  
323 hidden state prediction. In *Forty-second International Conference on Machine Learning*, 2025.

324 Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *stat*, 1050:9, 2015.

325 Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng,  
326 Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.

327 Sachit Menon, David Blei, and Carl Vondrick. Forget-me-not! Contrastive critics for mitigating posterior  
328 collapse. In *Uncertainty in Artificial Intelligence*, pages 1360–1370. PMLR, 2022.

329 Jürgen Schmidhuber. Curious model-building control systems. In *Proc. Int. Joint Conf. on Neural Networks*,  
330 pages 1458–1463, 1991a.

331 Jürgen Schmidhuber. A possibility for implementing curiosity and boredom in model-building neural controllers.  
332 In *Proc. Int. Conf. on Simulation of Adaptive Behavior: From Animals to Animats*, pages 222–227, 1991b.

333 Jürgen Schmidhuber. Learning complex, extended sequences using the principle of history compression. *Neural  
334 Computation*, 4(2):234–242, 1992. doi: 10.1162/neco.1992.4.2.234.

335 Oscar Skean, Md Rifat Arefin, Dan Zhao, Niket Nikul Patel, Jalal Naghiyev, Yann LeCun, and Ravid Shwartz-  
336 Ziv. Layer by layer: Uncovering hidden representations in language models. In *Forty-second International  
337 Conference on Machine Learning*, 2025.

338 Paul M Vitányi. Meaningful information. *IEEE Transactions on Information Theory*, 52(10), 2006.

339 Stephen Wolfram. A new kind of science. 2002.

340 LLM-Core Xiaomi, :, Bingquan Xia, Bowen Shen, Cici, Dawei Zhu, Di Zhang, Gang Wang, Hailin Zhang,  
341 Huaqiu Liu, Jiebao Xiao, Jinhao Dong, Liang Zhao, Peidian Li, Peng Wang, Shihua Yu, Shimao Chen, Weikun  
342 Wang, Wenhan Ma, Xiangwei Deng, Yi Huang, Yifan Song, Zihan Jiang, Bowen Ye, Can Cai, Chenhong He,  
343 Dong Zhang, Duo Zhang, Guoan Wang, Hao Tian, Haochen Zhao, Heng Qu, Hongshen Xu, Jun Shi, Kainan  
344 Bao, Kai Fang, Kang Zhou, Kangyang Zhou, Lei Li, Menghang Zhu, Nuo Chen, Qiantong Wang, Shaohui Liu,  
345 Shicheng Li, Shuhao Gu, Shuhuai Ren, Shuo Liu, Sirui Deng, Weiji Zhuang, Weiwei Lv, Wenyu Yang, Xin  
346 Zhang, Xing Yong, Xing Zhang, Xingchen Song, Xinzhe Xu, Xu Wang, Yihan Yan, Yu Tu, Yuanyuan Tian,  
347 Yudong Wang, Yue Yu, Zhenru Lin, Zhichao Song, and Zihao Yue. Mimo: Unlocking the reasoning potential  
348 of language model – from pretraining to posttraining, 2025. URL <https://arxiv.org/abs/2505.07608>.

## A Details on Sequence Models Trained from Scratch

We train all models using the Adam optimizer, a batch size of 16 and gradient norm clipping of 1.0. The learning rate is 0.0003, with a 500 step linear warm-up from zero and no decay. All losses are weighted equally, for the PHi loss we take the mean of the element-wise KL-Divergence for  $z$ , not the sum. Every model variation is trained 8 times with different random seeds for the initial weights and the procedurally generated data (which results in different memorized sequences and programs). The training of a model can be done on a single consumer-grade GPU (e.g., NVIDIA RTX 4090).

The base model is based on the Llama 3.2 architecture [Dubey et al., 2024b].

- Number of layers: 12
- Model dimensionality: 768
- Number of attention heads: 6
- MLP intermediate size: 2048
- Position of PHi layer: After layer 6

### PHi models:

To prevent posterior collapse, we employ an additional contrastive self-critic loss [Menon et al., 2022].

- Training steps: 30,000
- Placement of the PHi Layer: After the 10th layer
- $z$  dimensionality: 768
- $q_\psi$ : Linear transform
- $a_\xi$ : Linear transform
- $b_\kappa$ : Linear transform
- $M_\mu$ : One transformer block like the ones in the rest of the model
- $p_\psi$ : Linear transform

### MTD models:

- Training steps: 10,000
- $b_\kappa$ : Linear transform
- $M_\mu$ : One transformer block like the ones in the rest of the model

For generation of training and testing data, we follow Herrmann et al. [2025]. The only difference is that we do not perturb any tokens during training, and that we use the same models for the experiments in Sections 4.1.1 and 4.1.2.

## B Details on Pre-Trained Language Models

For our experiments, we use the SFT version of the MiMo-7B model [Xiaomi et al., 2025]. To calculate the MTD, we use the included MTP head that predicts one token in advance.

All experimental results include bootstrapped 95% confidence intervals.

## 384 C Additional Experimental Results

385 Figure 7 shows MTD and NLL for self-generated CoTs, broken down by category and difficulty  
 386 level. The results are comparable to the ones from the provided step-by-step solution, although the  
 387 differences between categories are less pronounced.

388 Figures 8 and 9 use the cumulative instead of the mean losses. Due to the fact that the provided  
 389 solutions as well as the generated ones grow in length as the problems become more difficult,  
 390 cumulative NLL also correlates positively with difficulty level.

391 Figure 10 shows the development of MTD and NLL across self-generated CoTs for the problems  
 392 of the GSM-8k test dataset (analogous to Figures 6c and 6d for MATH). Correct CoTs clearly have  
 393 lower MTD, and lower NLL. Interestingly, for the GSM-8k dataset, the shapes of the NLL curves  
 394 differ significantly from the shapes of the MTD, missing the prominent initial bump. Currently, we  
 395 have no explanation for this.

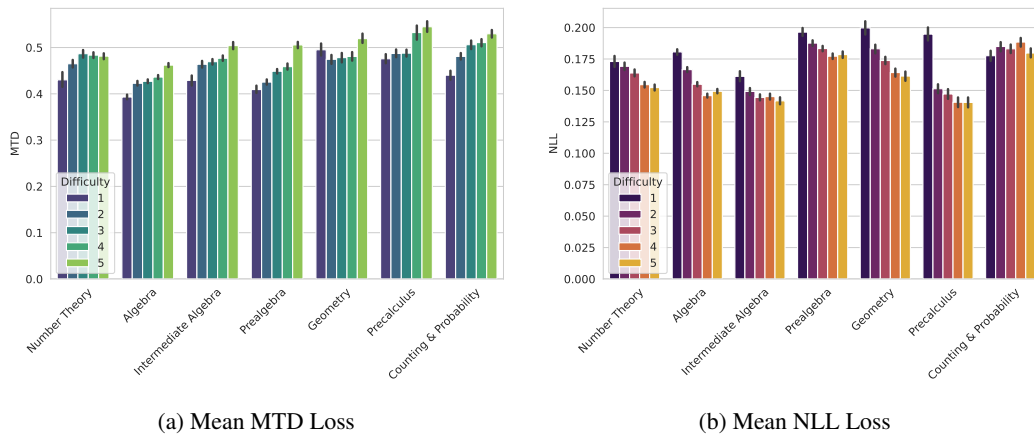
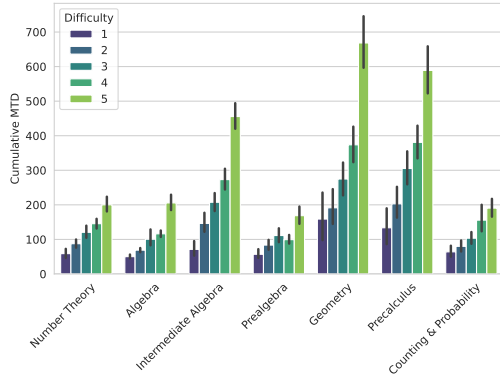
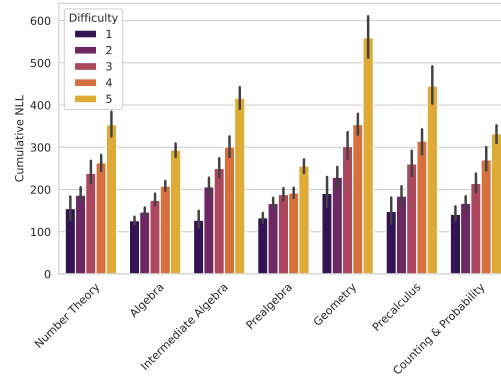


Figure 7: Mean losses of the MiMo model across self-generated CoTs for the problems of the MATH test set, grouped by category and difficulty level. Similarly as in Figure 5, we observe that MTD clearly grows with difficulty, as the model is making more use of its computational capacity when generating the solutions to more challenging problems. Also here, the mean NLL goes down with problem difficulty.

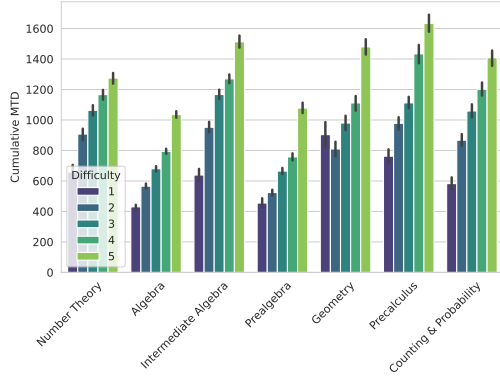


(a) Cumulative MTD Loss

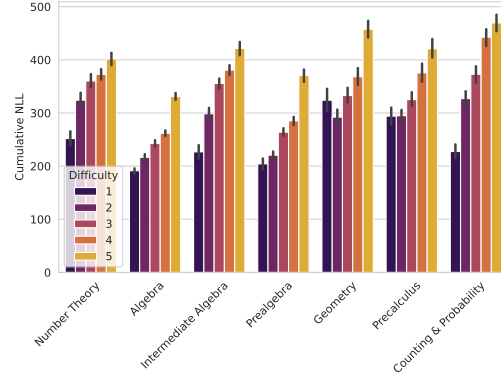


(b) Cumulative NLL Loss

Figure 8: Cumulative losses of the MiMo model across provided solutions from the MATH test set. Since more difficult problems have longer solutions, both cumulative MTD and cumulative NLL correlate with problem difficulty.

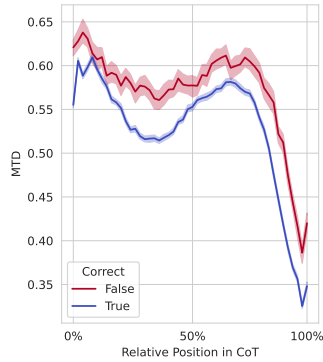


(a) Cumulative MTD Loss

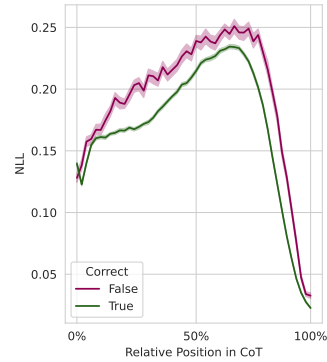


(b) Cumulative NLL Loss

Figure 9: Cumulative losses of the MiMo model across self-generated CoTs for the problems of the MATH test set. We observe a similar effect as in Figure 8.



(a) MTD vs. Correctness



(b) NLL vs. Correctness

Figure 10: Token-wise losses against relative positions in self-generated CoTs for the GSM-8k test dataset. Lower MTD and lower NLL are both associated with more correct reasoning.

## NeurIPS Paper Checklist

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: All concrete claims are supported by empirical results.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: See Section 5

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

### 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[NA\]](#)

Justification: No theoretical results in this paper.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: The information in the experiments section and the appendix is sufficient to reproduce all results.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

#### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: See above, code will be made public upon acceptance.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: See Appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: All experimental results include 95% confidence intervals

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.

- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: The experiments are relatively small scale and can be run on a single GPU.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification:

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA] .

Justification: We expect no direct societal impact from this work.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to

generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.

- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA] .

Justification:

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The model creates are given proper credit.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

## 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA] .

Justification: No new assets in this paper.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

#### 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification:

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

#### 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification:

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

#### 16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigor, or originality of the research, declaration is not required.

Answer: [NA] .

709

Justification:

710

Guidelines:

711

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.

712

713

- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.

714