
Complementarity-Driven Distillation from Multiple Foundation Models for DNA Classification

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 DNA sequence modeling has advanced with specialized foundation models such
2 as HyenaDNA, yet these models capture only partial genomic cues. In this work,
3 we investigate whether large language models (LLMs)—both subword-tokenized
4 (LLaMA) and byte-level (EvaByte)—provide complementary perspectives when
5 applied to DNA classification. Through experiments on the Human Enhancer
6 Cohn benchmark, we find that DNA-pretrained models and LLMs succeed on
7 largely disjoint subsets of data, revealing genuine cross-family complementarity.
8 Building on this insight, we propose a confidence-guided distillation framework
9 that aggregates supervision only from correct and confident teachers, producing soft
10 labels that safely transfer diverse knowledge. Our method consistently improves
11 both compact DNA-specific models and large byte-level LLMs, achieving gains of
12 up to +2.34 accuracy points while remaining robust against overfitting even under
13 near-perfect training accuracy. These findings highlight that DNA and language
14 models encode orthogonal yet synergistic representations, and that principled
15 distillation can unify them into a single model for robust genomic prediction.

16 1 Introduction

17 DNA sequence analysis lies at the heart of modern genomics. With the rapid advances in artificial
18 intelligence, deep learning has emerged as a powerful paradigm for tackling DNA-related tasks
19 such as enhancer detection, and regulatory element classification [Zhou and Troyanskaya, 2015,
20 Kelley et al., 2016]. Motivated by the success of large-scale language modeling [Devlin et al., 2019,
21 Touvron et al., 2023], a growing body of research has introduced DNA-specific foundation models
22 such as DNABERT [Ji et al., 2021], DNABERT2 [Zhou et al., 2023], and HyenaDNA [Nguyen et al.,
23 2023], which adapt Transformer to genomic sequences [Dalla-Torre et al., 2023, Avsec et al., 2021].
24 More recently, this line of work has been extended with models like Caduceus, which incorporates
25 bi-directional equivariant long-range modeling [Cao et al., 2024], and DNABERT-S, which introduces
26 species-aware embeddings for improved cross-species generalization [Zhang et al., 2024].

27 Despite these advances, current approaches remain limited. In particular, our understanding of
28 what features these models capture from DNA sequences—and how these features contribute to
29 classification performance—remains incomplete. More critically, DNA-pretrained models, though
30 specialized, may fail to capture alternative structural or semantic cues that could be extracted if DNA
31 sequences are treated as symbolic strings, analogous to natural language [Malusare et al., 2023].

32 In this work, we explore this perspective by hypothesizing that DNA sequences, viewed as ordered
33 character strings, can be effectively modeled not only by DNA-specific pretraining but also by
34 general-purpose large language models (LLMs). We investigate two complementary directions:
35 (i) a subword-tokenized LLM such as LLaMA [Touvron et al., 2023], which leverages Byte-Pair
36 Encoding (BPE), and (ii) a byte-level pre-trained model such as EvaByte [Zheng et al., 2025], which

37 processes inputs at the character level. The latter is particularly appealing, as it naturally aligns with
38 the character-based nature of DNA sequences, enabling finer-grained representations without reliance
39 on arbitrary subword segmentation [Malusare et al., 2023, Xue et al., 2022, Tay et al., 2022].

40 Through experiments, we demonstrate that DNA-pretrained models, standard LLMs, and byte-level
41 LLMs capture different aspects of DNA sequences. Importantly, these models succeed and fail on
42 complementary subsets of samples, both in test and training data. Building on this observation, we
43 propose a simple yet safe knowledge distillation framework [Hinton et al., 2015, Furlanello et al.,
44 2018, Zhang et al., 2019] that leverages model confidence to construct a soft-labeled distillation
45 dataset. By averaging the confidences of models that correctly classify a sample, we generate training
46 signals that preserve complementary knowledge. Distilling across DNA-specific and language-based
47 models consistently improves performance, and notably, the resulting models remain robust even
48 when trained to near-perfect accuracy on the distillation set.

49 Our contributions are threefold:

- 50 • We reveal that LLMs, including byte-level models, extract complementary representations
51 from DNA sequences compared to DNA-pretrained models.
- 52 • We introduce a confidence-guided distillation dataset construction method that is both simple
53 and resistant to overfitting.
- 54 • We empirically validate that this approach enhances DNA classification performance by
55 safely integrating complementary knowledge across model families.

56 **2 Method**

57 **2.1 Preliminary**

58 We consider three types of models: a DNA-pretrained model (HyenaDNA), a subword-based large
59 language model (LLaMA), and a byte-level language model (EvaByte). For LLaMA and EvaByte,
60 we follow a generative prompt-based setup: each DNA sequence is provided in the form of Use
61 your background knowledge about DNA enhancers. Classify the following DNA
62 sequence as an enhancer (0) or not (X):{DNA sequence}, Answer:., and the model is
63 trained to generate a binary answer (0 for positive and X for negative). In contrast, HyenaDNA is
64 trained following the original design in Nguyen et al. [2023], where a classifier head is attached to
65 the backbone and optimized with supervised cross-entropy loss.

66 **2.2 Motivation**

67 During evaluation, we observe that the three models (HyenaDNA, LLaMA, and EvaByte) often
68 succeed on different subsets of test samples, as shown in Figure 1. Although there exists a substantial
69 overlap among correctly predicted samples, the disagreement across models is considerably larger
70 than the discrepancy observed when comparing HyenaDNA models trained with different random
71 seeds, as shown in Appendix B. This indicates that LLaMA and EvaByte capture aspects of DNA
72 sequences that HyenaDNA fails to recognize, and vice versa.

73 Interestingly, such complementary behavior is not restricted to test data. Even among training samples,
74 the models show differences in which instances they predict correctly or incorrectly. This finding
75 suggests that each model focuses on distinct features of the input sequences during training, leading
76 to complementary strengths. Building on this observation, we propose to leverage training samples
77 themselves for distillation, enabling integration of the diverse knowledge captured by different
78 models.

79 **2.3 Distillation Framework**

80 Based on the above observation, we construct a confidence-guided distillation dataset. Specifically,
81 from the training set, we extract all samples that are correctly predicted by at least one of the three
82 models. For each such sample, we record the confidence scores of the models that made correct
83 predictions and compute the average of these scores to form a soft label. This process yields a
84 soft-labeled dataset that captures complementary signals across models.

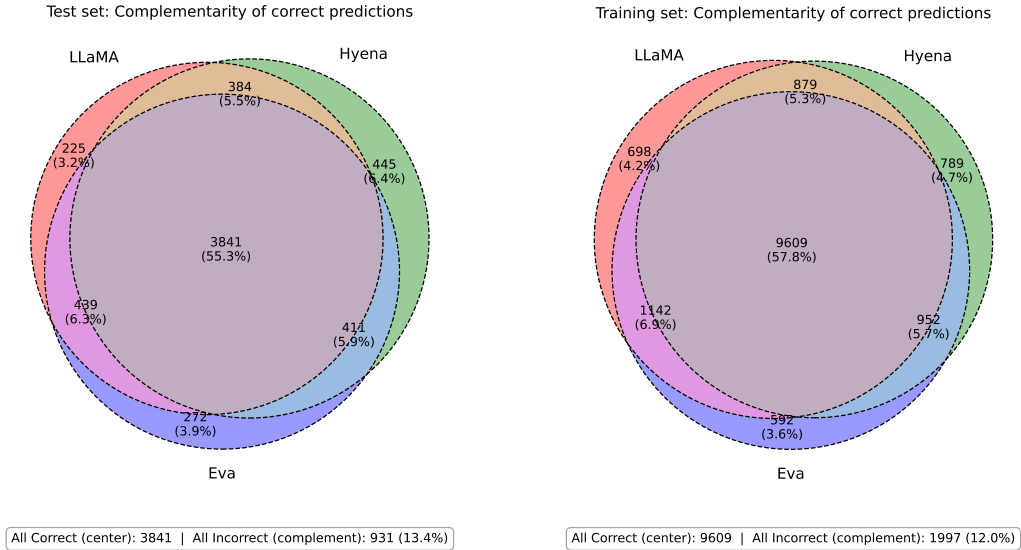


Figure 1: Complementarity of correct predictions among LLaMA, Hyena, and Eva on the test (left) and training (right) sets. Numbers indicate counts with percentages relative to the total; the complement region (all incorrect) is shown below each plot. Note that pairwise disagreements across model families substantially exceed the seed variance of a single family, evidencing genuinely complementary competence.

85 Finally, we use this distillation dataset to further train the models (HyenaDNA, LLaMA, and EvaByte).
 86 This simple yet effective framework allows each model to benefit from the unique perspectives of the
 87 others. The overall procedure is summarized in Algorithm 1.

Algorithm 1 Confidence-Guided Distillation

- 1: **Input:** Training set $\mathcal{D}_{\text{train}}$, teacher models \mathcal{M}
 - 2: Initialize distillation set $\mathcal{D}_{\text{distill}} \leftarrow \emptyset$
 - 3: **for** each sample $(x, y) \in \mathcal{D}_{\text{train}}$ **do**
 - 4: Collect predictions and confidences from all teachers
 - 5: Keep only teachers that predict y correctly
 - 6: **if** at least one correct teacher exists **then**
 - 7: Soft label $\tilde{y}(x) \leftarrow$ average confidence of correct teachers
 - 8: Add $(x, \tilde{y}(x))$ to $\mathcal{D}_{\text{distill}}$
 - 9: **end if**
 - 10: **end for**
 - 11: **Return** $\mathcal{D}_{\text{distill}}$
-

88 **3 Experiments**

89 **3.1 Experimental Setup**

90 We conduct all experiments on the Human Enhancer Cohn dataset from GenomicBenchmarks
 91 [Gresova et al., 2022], a widely used benchmark for enhancer identification. The task is formulated
 92 as a binary classification problem, distinguishing enhancer from non-enhancer sequences, with an
 93 equal number of samples per class to ensure balance. The dataset contains 10,421 positive and 10,422
 94 negative samples in the training set, and 3,474 samples per class in the test set.

95 We evaluate three model families: HyenaDNA (1.6M parameters), a DNA-pretrained model with a
 96 classifier head; EvaByte (6.5B parameters), a byte-level large language model; and LLaMA3-8B (8B
 97 parameters), a subword-tokenized large language model. All models are first trained and evaluated
 98 on the Human Enhancer Cohn dataset using the same train/test split.

99 After this supervised training stage, we perform confidence-guided distillation. Specifically, distilla-
 100 tion is applied to HyenaDNA, EvaByte, and LLaMA separately, starting from their baseline models
 101 trained on the Human Enhancer Cohn dataset, while keeping all other settings unchanged.

102 3.2 Baselines

103 As baselines, we use the three model families introduced above—HyenaDNA, EvaByte, and LLaMA3-
 104 8B—each trained individually on the dataset. In addition, we consider simple ensemble heuristics: (i)
 105 majority voting across the three models and (ii) a confidence-based strategy that selects the prediction
 106 from the most confident model. Finally, we compare these baselines against our confidence-guided
 107 distillation applied to each model family individually, as summarized in Table 1.

108 3.3 Main Results

109 Table 1 summarizes the effect of distillation compared to ensemble heuristics. Both HyenaDNA and EvaByte
 110 improve after distillation; LLaMA3-8B shows the largest gains. Majority voting and confidence-based en-
 111 sembling achieve slightly higher accuracy in aggregate, but they require multiple models at inference time. By
 112 contrast, our confidence-guided distillation compresses complementary knowledge into a single student, offer-
 113 ing voting and confidence-based ensembling achieve slightly higher accuracy in aggregate, but they require
 114 multiple models at inference time. By contrast, our confidence-guided distillation compresses complementary
 115 knowledge into a single student, offering comparable accuracy with much greater efficiency and deployability.
 116
 117
 118
 119
 120
 121
 122

Table 1: Effect of confidence-guided distillation on the Human Enhancer Cohn dataset. Distilled rows indicate models fine-tuned from the corresponding baseline. Δ denotes the accuracy gain over baseline.

Model	Variant	Pretraining	Acc. (%)	Δ
HyenaDNA	Baseline	DNA	73.06	–
	Distilled	DNA	73.33	+0.27
EvaByte	Baseline	Natural Language	71.60	–
	Distilled	Natural Language	73.50	+1.90
LLaMA3-8B	Baseline	Natural Language	70.06	–
	Distilled	Natural Language	72.40	+2.34
Majority Voting Ensemble			73.07	–
Confidence-based Ensemble			73.70	–

123 3.4 No Collapse under High Memorization

124 A key finding is that even as both EvaByte and LLaMA achieve very high training accuracy (98.36% and 95.6%, respectively), their test accuracy does not collapse and stays stable around 72–74%. Even though the models memorize the training data almost perfectly, their test accuracy remains stable without degradation. This indicates that our distillation targets do not induce overfitting collapse, but instead provide robust generalization under high memorization. Figure 2 illustrates the training trajectory.
 125
 126
 127
 128
 129
 130
 131
 132
 133
 134

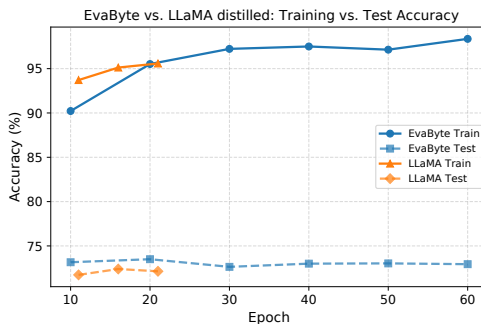


Figure 2: EvaByte and LLaMA under confidence-guided distillation.

135 4 Conclusion and Discussion

136 We investigated the complementarity between DNA-pretrained models and language-pretrained
 137 LLMs for DNA classification, showing that they capture orthogonal genomic signals. Building
 138 on this, we proposed a confidence-guided distillation strategy that integrates only reliable teacher
 139 signals, yielding consistent improvements—modest for HyenaDNA and substantial for EvaByte and
 140 LLaMA—while remaining robust against overfitting.

141 A key limitation of our approach is that distillation requires an additional training phase, which adds
 142 computational overhead and may constrain scalability. Looking forward, it would be valuable to move
 143 beyond post-hoc distillation and design unified architectures that can directly fuse heterogeneous
 144 sources of complementary knowledge. More broadly, our findings suggest that language-based
 145 LLMs can provide non-redundant genomic insights, pointing toward new opportunities for leveraging
 146 foundation models to advance trustworthy biomedical AI.

147 **References**

- 148 Žiga Avsec, Vikram Agarwal, Daniel Visentin, Joseph R Ledsam, Agnieszka Grabska-Barwińska,
149 Kristin R Taylor, Yannis Assael, John Jumper, Pushmeet Kohli, and David R Kelley. Effective
150 gene expression prediction from sequence by integrating long-range interactions. *Nature Methods*,
151 18(10):1196–1203, 2021.
- 152 Shuxiao Cao, Zhiwei Hong, Stephen A. Baccus, Christopher Ré, et al. Caduceus: Bi-directional
153 equivariant long-range dna sequence modeling. *arXiv preprint arXiv:2405.07990*, 2024.
- 154 Hugo Dalla-Torre, Tristan Bepler, and et al. Nucleotide transformer: building and evaluating robust
155 foundation models for human genomics. *bioRxiv*, 2023.
- 156 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
157 bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference
158 of the North American Chapter of the Association for Computational Linguistics (NAACL-HLT)*,
159 pages 4171–4186, 2019.
- 160 Tommaso Furlanello, Zachary C Lipton, Michael Tschannen, Laurent Itti, and Anima Anandkumar.
161 Born-again neural networks. In *Proceedings of the 35th International Conference on Machine
162 Learning (ICML)*, 2018.
- 163 Katarina Gresova, Vlastimil Martinek, David Cechak, Petr Simecek, Panagiotis Alexiou, et al.
164 Genomic benchmarks: A collection of datasets for genomic sequence classification. *bioRxiv*, 2022.
165 URL <https://www.biorxiv.org/content/10.1101/2022.06.08.495248>.
- 166 Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv
167 preprint arXiv:1503.02531*, 2015.
- 168 Yanrong Ji, Zhihan Zhou, Han Liu, and Ramana V Davuluri. Dnabert: pre-trained bidirectional
169 encoder representations from transformers model for dna-language in genome. *Bioinformatics*, 37
170 (15):2112–2120, 2021.
- 171 David R Kelley, Jasper Snoek, and John L Rinn. Basset: learning the regulatory code of the accessible
172 genome with deep convolutional neural networks. *Genome Research*, 26(7):990–999, 2016.
- 173 Aditya Malusare, Harish Kothandaraman, et al. Understanding the natural language of dna using
174 encoder-decoder foundation models with byte-level precision. *arXiv preprint arXiv:2311.02333*,
175 2023.
- 176 Eric Nguyen, Michael Poli, Marjan Faizi, Armin Thomas, Callum Birch-Sykes, Michael Wornow,
177 Aman Patel, Clayton Rabideau, Stefano Massaroli, Yoshua Bengio, Stefano Ermon, Stephen A
178 Baccus, and Christopher Ré. Hyenadna: Long-range genomic sequence modeling at single
179 nucleotide resolution. *arXiv preprint arXiv:2306.15794*, 2023.
- 180 Yi Tay, Mostafa Dehghani, Vinh Q Tran, Xavier Garcia, Dara Bahri, Zhen Xiao, Naveen Arivazha-
181 gan, and Donald Metzler. Charformer: Fast character transformers via gradient-based subword
182 tokenization. In *International Conference on Learning Representations (ICLR)*, 2022.
- 183 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
184 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand
185 Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language
186 models. *arXiv preprint arXiv:2302.13971*, 2023.
- 187 Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya
188 Barua, and Colin Raffel. Byt5: Towards a token-free future with pre-trained byte-to-byte models.
189 In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (ACL)*,
190 2022.
- 191 Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Your classifier is
192 secretly an energy based model and you should treat it like one. In *Advances in Neural Information
193 Processing Systems (NeurIPS)*, 2019.

- 194 Hao Zhang, Zhihan Zhou, Yanrong Ji, and Han Liu. Dnabert-s: Pioneering species differentiation
195 with species-aware dna embeddings. *arXiv preprint arXiv:2406.01914*, 2024.
- 196 Lin Zheng, Xueliang Zhao, Guangtao Wang, Chen Wu, David Dong, Angela Wang, Mingran Wang,
197 Yun Du, Haige Bo, Amol Sharma, Bo Li, Kejie Zhang, Changran Hu, Urmish Thakker, and
198 Lingpeng Kong. Evabyte: Efficient byte-level language models at scale. [https://hkunlp.
199 github.io/blog/2025/evabyte](https://hkunlp.github.io/blog/2025/evabyte), 2025.
- 200 Jian Zhou and Olga G Troyanskaya. Predicting effects of noncoding variants with deep learning-based
201 sequence model. *Nature Methods*, 12(10):931–934, 2015.
- 202 Zhihan Zhou, Yanrong Ji, Weijian Li, Pratik Dutta, Ramana Davuluri, and Han Liu. Dnabert-2: Effi-
203 cient foundation model and benchmark for multi-species genome. *arXiv preprint arXiv:2306.15006*,
204 2023.

205 A Experimental Details

206 We follow the standard GenomicBenchmarks split: training and test sets are provided, and we further
207 reserve 20% of the training set for validation. All models are trained with AdamW optimizer, a linear
208 warmup of 100 steps, and early stopping with patience of 3 based on validation performance. Final
209 results are reported on the held-out test split using the checkpoint with the best validation score.

210 A.1 HyenaDNA: Training and Evaluation Details

211 For HyenaDNA, DNA sequences are tokenized at the character level ($\{A,C,G,T,N\}$) with maximum
212 length 500. The model is initialized from a pretrained HyenaDNA checkpoint. Training uses AdamW
213 (learning rate 6×10^{-4} , weight decay 0.1, batch size 256) and cross-entropy loss. Unlike the original
214 HyenaDNA setup that trains for 100 epochs, we apply early stopping on validation loss and select the
215 best checkpoint.

216 A.2 LLaMA3-8B: Training and Evaluation Details

217 For LLaMA3-8B, we adopt a generative prompt-based setup, where each DNA sequence
218 is presented as: Use your background knowledge about DNA enhancers. Classify
219 the following DNA sequence as an enhancer (O) or not (X): {DNA sequence},
220 Answer: and the model is trained to autoregressively predict the label (O/X). We apply LoRA
221 (rank $r = 8$, $\alpha = 16$, dropout 0.05) to the attention projection modules (q_proj , k_proj , v_proj ,
222 up_proj , $down_proj$). Optimization uses learning rate 1×10^{-5} and batch size 16. We train for up
223 to 100 epochs with early stopping on validation accuracy.

224 A.3 EvaByte: Training and Evaluation Details

225 EvaByte is trained with the same prompt-based causal language modeling setup as LLaMA3-8B,
226 where each DNA sequence is followed by a binary label (O/X) to be predicted autoregressively. The
227 main differences are in the hyperparameters: we use batch size 4 (instead of 16), and learning rate
228 1×10^{-5} without weight decay. LoRA (rank $r = 8$, $\alpha = 16$, dropout 0.05) is applied to the attention
229 projection modules as in the LLaMA3-8B experiments.

230 B Complementarity Beyond Seed Variance

231 To verify that the observed complementarity is not simply due to randomness in training or seed
232 variance, we compared overlaps between HyenaDNA models trained with different seeds and between
233 HyenaDNA and other model families. Table 2 shows that the disagreement between HyenaDNA runs
234 (seed0 vs. seed42) is relatively small, while the disagreement between HyenaDNA and EvaByte is
235 substantially larger. This confirms that the complementarity arises from model family differences
236 rather than stochastic variation.

237 Compared to the $\sim 10\%$ disagreement across HyenaDNA seeds, cross-family comparisons with
238 EvaByte exhibit nearly double the disagreement rate (around 20%). This demonstrates that large

Table 2: Comparison of prediction overlaps on the Human Enhancer Cohn dataset. “Both Correct” denotes samples where both models predict correctly, “Disagree” denotes samples where one model is correct and the other is not, and “Both Wrong” denotes samples where both fail. Percentages are relative to the full test set.

Pair of Models	Both Correct	Disagree	Both Wrong
Hyena (seed0, 72.25%) vs. Hyena (seed42, 73.13%)	67.7%	9.9%	22.3%
Hyena (seed42, 73.13%) vs. EvaByte (71.43%)	61.2%	22.2%	16.6%
Hyena (seed0, 72.25%) vs. EvaByte (71.43%)	62.4%	19.9%	18.7%

239 language models such as EvaByte capture distinct genomic cues that DNA-pretrained models like
 240 HyenaDNA do not, indicating genuine complementarity rather than random variation. Therefore, our
 241 distillation framework leverages orthogonal knowledge across model families, beyond what can be
 242 obtained by simply retraining the same architecture with different seeds.

243 **NeurIPS Paper Checklist**

244 **1. Claims**

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

247 Answer: [Yes]

248 Justification: Abstract and introduction includes paper's contributions.

249 Guidelines:

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

259 **2. Limitations**

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

261 Answer: [Yes]

262 Justification: Section 4 Conclusion and Discussion includes limitations of the work.

263 Guidelines:

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

290 **3. Theory assumptions and proofs**

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

293 Answer:[NA]

294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347

Justification: This paper does not include theorem.

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: Section 3 Experiments include those information.

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?

348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399

Answer: [No]

Justification: All experiments and method are explained in detail in the manuscript, so that they can be reproduced easily. Furthermore, the code will be available in camera-ready version.

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: Section 3 Experiments include those information.

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: [No]

Justification: The experiments are conducted on a single seed but the results are justified with same results across various model architecture.

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.)

- 400 • The assumptions made should be given (e.g., Normally distributed errors).
- 401 • It should be clear whether the error bar is the standard deviation or the standard error
- 402 of the mean.
- 403 • It is OK to report 1-sigma error bars, but one should state it. The authors should
- 404 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
- 405 of Normality of errors is not verified.
- 406 • For asymmetric distributions, the authors should be careful not to show in tables or
- 407 figures symmetric error bars that would yield results that are out of range (e.g. negative
- 408 error rates).
- 409 • If error bars are reported in tables or plots, The authors should explain in the text how
- 410 they were calculated and reference the corresponding figures or tables in the text.

411 8. Experiments compute resources

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

415 Answer: [Yes]

416 Justification: Section 3 Experiments include those information.

417 Guidelines:

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

426 9. Code of ethics

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

429 Answer: [Yes]

430 Justification: This research follows the guidelines.

431 Guidelines:

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

437 10. Broader impacts

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

440 Answer: [No]

441 Justification: There is no societal impact of the work performed.

442 Guidelines:

- 443 • The answer NA means that there is no societal impact of the work performed.
- 444 • If the authors answer NA or No, they should explain why their work has no societal
 445 impact or why the paper does not address societal impact.
- 446 • Examples of negative societal impacts include potential malicious or unintended uses
 447 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
 448 (e.g., deployment of technologies that could make decisions that unfairly impact specific
 449 groups), privacy considerations, and security considerations.

- 450 • The conference expects that many papers will be foundational research and not tied
451 to particular applications, let alone deployments. However, if there is a direct path to
452 any negative applications, the authors should point it out. For example, it is legitimate
453 to point out that an improvement in the quality of generative models could be used to
454 generate deepfakes for disinformation. On the other hand, it is not needed to point out
455 that a generic algorithm for optimizing neural networks could enable people to train
456 models that generate Deepfakes faster.
- 457 • The authors should consider possible harms that could arise when the technology is
458 being used as intended and functioning correctly, harms that could arise when the
459 technology is being used as intended but gives incorrect results, and harms following
460 from (intentional or unintentional) misuse of the technology.
- 461 • If there are negative societal impacts, the authors could also discuss possible mitigation
462 strategies (e.g., gated release of models, providing defenses in addition to attacks,
463 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
464 feedback over time, improving the efficiency and accessibility of ML).

465 11. Safeguards

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

469 Answer: [NA]

470 Justification: This research does not include risk for misuse.

471 Guidelines:

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

482 12. Licenses for existing assets

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

486 Answer: [Yes]

487 Justification: They are properly cited in references.

488 Guidelines:

- 489 • The answer NA means that the paper does not use existing assets.
- 490 • The authors should cite the original paper that produced the code package or dataset.
- 491 • The authors should state which version of the asset is used and, if possible, include a
492 URL.
- 493 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 494 • For scraped data from a particular source (e.g., website), the copyright and terms of
495 service of that source should be provided.
- 496 • If assets are released, the license, copyright information, and terms of use in the
497 package should be provided. For popular datasets, paperswithcode.com/datasets
498 has curated licenses for some datasets. Their licensing guide can help determine the
499 license of a dataset.
- 500 • For existing datasets that are re-packaged, both the original license and the license of
501 the derived asset (if it has changed) should be provided.

502 • If this information is not available online, the authors are encouraged to reach out to
503 the asset’s creators.

504 **13. New assets**

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

507 Answer: [NA]

508 Justification: This paper does not release new assets.

509 Guidelines:

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

518 **14. Crowdsourcing and research with human subjects**

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

522 Answer: [NA]

523 Justification: This paper does not involve crowdsourcing nor research with human subjects.

524 Guidelines:

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

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

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

539 Answer: [NA]

540 Justification: This paper does not involve crowdsourcing nor research with human subjects.

541 Guidelines:

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

552 **16. Declaration of LLM usage**

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

557 Answer: [NA]

558 Justification: The core method development in this research does not involve LLMs as any
559 important, original, or non-standard components.

560 Guidelines:

- 561 • The answer NA means that the core method development in this research does not
562 involve LLMs as any important, original, or non-standard components.
- 563 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)
564 for what should or should not be described.