

Beyond fine-tuning: LoRA modules boost near-OOD detection and LLM security

Etienne Salimbeni^{1,2}, Francesco Craighero¹, Renata Khasanova², Milos Vasic², Pierre Vandergheynst¹

¹ EPFL, Lausanne, Switzerland, ² Oracle Labs, Zurich, Switzerland

{etienne.salimbeni, francesco.craighero, pierre.vandergheynst}@epfl.ch

{milos.vasic, renata.khasanova}@oracle.com

Abstract—Under resource constraints, LLMs are usually fine-tuned with additional knowledge using Parameter Efficient Fine-Tuning (PEFT), using Low-Rank Adaptation (LoRA) modules. In fact, LoRA injects a new set of small trainable matrices to adapt an LLM to a new task, while keeping the latter frozen. At deployment, LoRA weights are subsequently merged with the LLM weights to speed up inference. In this work, we show how to exploit the unmerged LoRA’s embedding to boost the performance of Out-Of-Distribution (OOD) detectors, especially in the more challenging near-OOD scenarios. Accordingly, we demonstrate how improving OOD detection also helps in characterizing wrong predictions in downstream tasks, a fundamental aspect to improve the reliability of LLMs. Moreover, we will present a use-case in which the sensitivity of LoRA modules and OOD detection are employed together to alert stakeholders about new model updates. This scenario is particularly important when LLMs are out-sourced. Indeed, test functions should be applied as soon as the model changes the version in order to adapt prompts in the downstream applications. In order to validate our method, we performed tests on Multiple Choice Question Answering datasets, by focusing on the medical domain as a fine-tuning task. Our results motivate the use of LoRA modules even after deployment, since they provide strong features for OOD detection for fine-tuning tasks and can be employed to improve the security of LLMs.

1. Introduction

Large Language Models (LLMs) are gaining popularity due to their general-purpose capabilities and are increasingly integrated into real-world applications, including medicine [31] and finance [20]. Their fast developing pace and ease of integration is alarming, since misconfiguration can be particularly damaging [4, 16, 34]. New government regulations are focusing on LLM-based applications. The *EU AI Act* [1] and the *White House Executive Order* on AI systems [3] are setting plans for their safe deployment, including the “robust monitoring of AI systems” [1]. Additionally, new *OWASP* [2] guidelines have been published, highlighting the security risks of integrating LLM into applications.

One major challenge of Machine Learning is protecting against unexpected behaviours of the model. Indeed, real-world applications might involve data that differs from the training one due to distributional shifts [36]. Coupled with random effects in the data, these shifts can make the model more uncertain about its predictions [13]. Consequently, detecting such Out-Of-Distribution (OOD) instances [36]

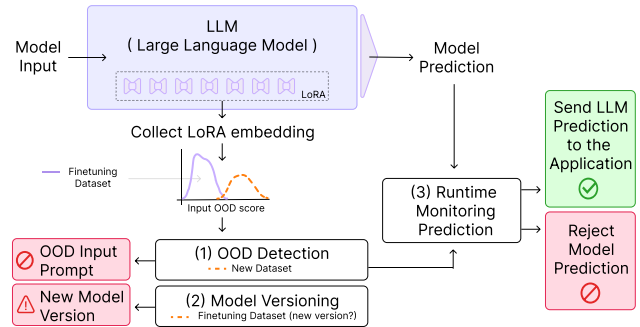


Figure 1: **Boosting LLM security with LoRA modules.** Given a fine-tuned LLM and the LoRA embeddings of the FT dataset, one can check: (1) if the LoRA embeddings of a new dataset are OOD, (2) if the model version has changed by detecting changes in LoRA embeddings, (3) if a prediction should be discarded due to an OOD input sample or low-confidence output.

is crucial to allow users to discard untrustworthy predictions. While OOD detection has been a fast-growing field, especially on classification tasks, these approaches have been also recently extended to LLMs and text generation [25]. In this paper, we will study OOD detection in the context of fine-tuned LLMs employing Low-Rank Adaptation (LoRA) modules [11].

Fine-tuning is a common practice for adapting a model to a specific domain. However, recent results raise new concerns on the reliability of fine-tuned LLMs. Indeed, fine-tuning can deteriorate previous safety alignments enforced during pre-training [23]. Moreover, it has been shown that fine-tuning worsens OOD robustness [7].

Low-Rank Adaptation (LoRA) modules [11] are commonly used to allow fine-tuning LLMs under resources constraints. Given a froze LLM, these small trainable modules are first injected for task adaptation and then merged with the model to reduce latency at inference time. Originally designed for fine-tuning purposes, LoRA modules are now being employed for greater control beyond their original purpose. Such applications include task arithmetic (add, combine or remove learned properties) [37], scaling the influence of the fine-tuned task at inference[27], and switching tasks using dynamic LoRA module routing [12, 28].

In the following, we show how *unmerged* LoRA modules can also be exploited to improve the security and reliability of LLMs. First, we show that LoRA embeddings are more sensitive to near-OOD samples, allowing simpler

OOD detectors such as the Mahalanobis Distance [18] to perform well in most scenarios. Second, we will present a novel use-case of OOD detection for model inspection. Model updates might in fact require version checking [9], to prevent major security flaws such as backdoor attacks [35] as well as simple misconfiguration in the LLM service supply chain [9]. With LoRA embeddings, one can easily detect model changes even under subtle updates. Last, we will test how LoRA embeddings improve runtime prediction monitoring, also known as selective prediction [8, 17, 33], when an LLM is employed for downstream tasks such as question answering. While OOD detection accounts for unintended inputs, a prediction might be uncertainty also due to random effects in the data. These two sources of uncertainty are usually referred to as epistemic, due to lack of knowledge, and aleatoric uncertainty, due to the stochastic nature of the data-generating process [13]. Similar to previous results on vision tasks [15], we will show how aggregating OOD detection and the entropy of the model confidence can improve the reliability of LLMs. This combined approach, accounting for the two sources of uncertainty, improves the detection of incorrect predictions compared to taking each individual metric alone. We will focus on decoder-only LLMs and medical Multiple Choice Question Answering (MCQA). However, it’s important to note that our methods are applicable to other large pre-trained models fine-tuned with LoRA and other generation tasks.

Contributions. In Section 3 we show that LoRA embedding boost the detection performance in near-OOD scenarios over the fine-tuning task. In Section 3 we present a novel use-case where OOD detection is employed to detect model updates. Last, in Section 3 we combine OOD detection and the output entropy to improve near-OOD runtime prediction monitoring in downstream tasks.

2. Methods

Datasets and Model. We integrated the previous results on OOD detection with RMD within the abstractive summarization and translation domains [25] by focusing on multiple question answering, which limits the number of generated token to 1. We selected three MCQA datasets and chose the medical domain as a fine-tuning task, by considering the MedMCQA [22] and the PubMedQA [14] datasets. Then, we employed the MMLU [10] multi-domain dataset to define both near- and far-OOD samples (refer to Section A.1 to get the subtasks assigned to each category). In contrast to [25], we use a decoder only language model: Llama2-7B [32]. Llama2-7B has a vocabulary size of 32 000, an embedding size of 4 096 and has 32 layers. We fine-tuned the model with LoRA [11] on the MedMCQA training split, using a batch size of 32, the Adam optimizer and a learning rate of $2e-4$. Moreover, we set LoRA to rank 16 and attached it to the query and value projections of each transformer layer. Concatenating all LoRA embeddings leads to a final embedding of size $32 \cdot 2 \cdot 16 = 2048$.

Embeddings. We compare two types of embeddings: last layer activations and LoRA embeddings. LoRA reparametrization of the i -th layer can be expressed as

$l^i(x) = W_0^i x + B^i A^i x$, where W_0^i is the pretrained frozen weights and $B^i A^i$ are two matrices of the LoRA module. Now, given an input of N tokens, we define the last layer activation embedding as $E_{LLA}(x) := \frac{1}{N} \prod_{i=1}^N l^i(x)$ and LoRA embeddings as $E_{LORA}(x) = \frac{1}{N} \prod_{i=1}^N \prod_{j=1}^L A^j x$. Where l^i is the final layer activation for token i , and $\prod_{j=1}^L A^j x$ denotes the concatenation of all LoRA modules intermediate activation $A^j x$ for the L layers (see Fig. A.1). Both embeddings are scaled through division by the maximum value. Importantly, we considered multiple layer embeddings only with LoRA, due to its reduced dimensionality compared to the full-rank layers. Both MD and RMD employed the embeddings on the fine-tuning dataset to compute train and train .

2.1. OOD Detection and Prediction Monitoring

In order to perform OOD detection, we selected three approaches with different requirements and ease of use. The Mahalanobis Distance (MD) [19] is a well-known approach for OOD detection that has the advantage of not requiring any hyperparameter tuning. The Relative Mahalanobis Distance (RMD) [24] is an improvement over MD that is expected to boost the performance in near-OOD scenarios, by correcting the MD with given a so-called background distribution. In our case, the embeddings of PubMedQA will be used as the background dataset. Last, KNN [30] is a recently proposed OOD detector, which requires the selection of the number of neighbours k (in our experiments, we used $k = 100$).

By detecting the embeddings that are far from the in-distribution ones, OOD detectors capture the epistemic uncertainty of the model. While the epistemic uncertainty is a result of lack of knowledge, the aleatoric uncertainty is related to the randomness in the data [13]. In order to estimate the aleatoric uncertainty, we simply compute the entropy of the token providing the answer to the question. Similar to [15], we will consider a combination of the two uncertainty in order to define a stronger approach to monitor the predictions of our model, also called selective prediction [8, 17, 33]. For this experiment, we will consider the MD approach, that thanks to the LoRA embeddings achieves good performance while having less requirements than RMD and KNN (see Table 1). While MD has no upper-bound, the entropy has range $[0; 1]$. Therefore, in order to combine them, it would be convenient to rescale the former. Since the squared MD follows a Chi-squared distribution with degrees of freedom equal to the number of dimensions [21], we can take the p -value of the Chi-squared distribution instead of the distance to obtain a normalized value. The final metric will just be the sum of the p -value and the entropy. More details about each method can be found in Section A.2.

3. Results

LoRA Modules Improve Near-OOD Detection. In Table 1 we compare the AUROC score for OOD detection of different embeddings (E_{LLA} , E_{LORA}) on both near- and far-OOD datasets, as defined in Section A.1, against the test dataset of MedMCQA (our in-distribution fine-tuning domain). In accordance with the results reported

Method	Near OOD				Far OOD
	clinical knowledge	anatomy	college biology	computer science	professional law
Perplexity	0.651	0.383	0.654	0.587	0.712
	Last Layer Activation (E_{LLA})				
KNN	0.387	0.296	0.786	0.997	0.999
MD	0.428	0.312	0.774	0.997	0.999
RMD* (baseline)	0.688	0.730	0.998	0.997	0.999
	LoRA (E_{LoRA})				
KNN	0.819	0.729	0.890	0.997	0.998
MD	0.814	0.733	0.890	0.996	0.994
RMD*	0.828	0.762	0.998	0.993	0.999

Table 1: OOD detection AUROCs. AUROCs distinguishing MMLU tasks from the MedMCQA dataset. * RMD requires a background dataset. (baseline) The approach of [25].

in [25], the perplexity proves to be a poor choice as an OOD score, as it struggles to distinguish even far-OOO datasets. When employing the last layer embeddings, all the methods perfectly discriminate far-OOO datasets. However, in near-OOO scenarios only RMD demonstrates positive performance, while KNN and MD fail completely. On the other hand, LoRA embeddings allow KNN and MD to perform on par with RMD on the near-OOO datasets, while keeping the same high performance on the far-OOO ones. As clearly emerges from Fig. A.2, LoRA embeddings boost the performance of the simpler MD approach, that neither requires hyperparameter tuning nor additional datasets like KNN and RMD, respectively. Indeed, RMD heavily depends on the goodness of the background dataset to perform well in the near-OOO dataset. Overall, LoRA's improvement over the last layer embedding is promising towards an improved and more efficient OOD detection in near-tuning tasks.

Detecting Model Updates. Given the good performance of the simple MD approach on LoRA embeddings, even in near-OOO scenarios, we investigate an interesting use-case to improve the security of near-tuned LLMs: detecting the degree of change of a model version update. This time, instead of checking if an external dataset is OOD, we aim to detect whether the embeddings of the in-distribution data have changed due to a (possibly unexpected) model update.

OpenAI's models endpoint degradation over time on some specific tasks [6] underlines the practical significance of this issue. Existing methods, such as verifying model weights hashes [9] or using zero-knowledge proofs [29], offer only a binary indication of model change. Given a dataset of interest and a model version, our approach is instead able to quantify model change. Such a scenario is relevant when a stakeholder out-sources LLM for a specific near-tuning task, where a model update might trigger a testing cascade on downstream tasks. Indeed, prompts may be invalidated on a different model version and malicious updates might inject backdoors in the model [35]. In Fig. 2, we present the MD AUROCs for discriminating between the embeddings of our LLM near-tuned for 500 steps on the MedMCQA training set (model version 0) and those obtained after near-tuning for the near-tuning task for security purposes. Note that our > 500 steps (next versions). Clearly, LoRA embeddings are much more sensitive to model updates than the last layer ones: while the latter has an AUROC: 8 1000

Figure 2: AUROCs distinguishing the embeddings at different near-tuning steps. AUROCs of the Mahalanobis Distance distinguishing MedMCQA embeddings (LoRA, last layer) after 500 near-tuning steps (model version 0) from the ones after > 500 steps (next model versions).

	MedMCQA	Near OOD	Far OOD
Entropy	0.554	0.541	0.547
MD Last Layer Activation	0.528	0.509	0.510
MD Last Layer Activation + Entropy	0.582	0.550	0.549
MD LORA	0.531	0.523	0.509
MD LORA + Entropy	0.589	0.576	0.543

Table 2: AUROCs scores when differentiating correct and incorrect predictions. We considered the MedMCQA validation dataset, and the near- and far-OOO datasets defined in Section A.1.

near-tuning steps after version 0 at 500 steps.

Runtime Monitoring Predictions. In Table 2 we report the AUROCs when detecting incorrect model predictions, i.e., wrong answer choices. We tested the output entropy, MD on the two types of embeddings and a combination of the two. The results show again how LoRA helps to improve MD in the near-OOO scenario, even if the setting is different than Section 3. Moreover, aggregating MD and entropy achieves the best performance, due to the different sources of uncertainty captured by the two metrics, i.e. epistemic and aleatoric [13].

4. Conclusion

In our experiments, we found compelling evidence supporting the hypothesis that LoRA embeddings possess stronger near-OOO properties compared to last layer activations and perplexity in near-tuning tasks, integrating our findings with previous research on OOD detection in LLMs [25]. This enables LLM-based applications to better monitor whether the model is being used for the intended task, to quantify the model version changes when the LLM is out-sourced, and to halt the model when the uncertainty about its predictions is too high. Importantly, LoRA modules allow us to employ simpler approaches for OOD detection, such as the Mahalanobis distance, that neither rely on additional data nor require hyperparameter tuning. Our findings suggest that LoRA weights should be kept near-tuned for 500 steps on the MedMCQA training set (model version 0) and those obtained after near-tuning for the near-tuning task for security purposes. Note that our > 500 steps (next versions). Clearly, LoRA embeddings are much more sensitive to model updates than the last layer ones: while the latter has an AUROC: 8 1000 practice in platforms like HuggingFace.

4.1. Funding

F. Craighero is funded by the Swiss National Science Foundation (SNSF) Sinergia grant CRSII5_205884.

References

- [1] European artificial intelligence act. <https://www.europarl.europa.eu/doceo/document/TA-9-2023-0236/EN.pdf>, 2023.
- [2] Owasp top 10 llm applications. <https://owasp.org/www-project-top-10-for-large-language-model-applications/>, 2023.
- [3] Us executive order on the safe, secure, and trustworthy development and use of artificial intelligence. <https://www.whitehouse.gov/briefing-room/presidential-actions/2023/10/30/executive-order-on-the-safe-secure-and-trustworthy-development-and-use-of-artificial-intelligence/>, 2023.
- [4] Timothy Bickmore, Ha Trinh, Steinn Ólafsson, Teresa O'Leary, Reza Asadi, Nathaniel Rickles, and Ricardo Cruz. Patient and consumer safety when using conversational assistants for medical information: Observational study (preprint). *Journal of Medical Internet Research*, 20, 07 2018.
- [5] Jiefeng Chen, Jinsung Yoon, Sayna Ebrahimi, Serkan Ö. Arik, Tomas Pfister, and Somesh Jha. Adaptation with self-evaluation to improve selective prediction in llms. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023*, Singapore, December 6-10, 2023, pages 5190–5213. Association for Computational Linguistics, 2023.
- [6] Lingjiao Chen, Matei Zaharia, and James Zou. How is chatgpt's behavior changing over time? *arXiv preprint arXiv:2307.09009*, 2023.
- [7] Sishuo Chen, Wenkai Yang, Xiaohan Bi, and Xu Sun. Fine-tuning deteriorates general textual out-of-distribution detection by distorting task-agnostic features. In Andreas Vlachos and Isabelle Augenstein, editors, *Findings of the Association for Computational Linguistics: EACL 2023*, Dubrovnik, Croatia, May 2-6, 2023, pages 552–567. Association for Computational Linguistics, 2023.
- [8] Yonatan Geifman and Ran El-Yaniv. Selective classification for deep neural networks. *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, page 4885–4894, Red Hook, NY, USA, 2017. Curran Associates Inc.
- [9] Wei Hao, Daniel Mendoza, Rafael da Silva, Deepak Narayanan, and Amar Phanishayee. Mgit: A model versioning and management system. 2023.
- [10] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *International Conference on Learning Representations*, 2021.
- [11] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. *International Conference on Learning Representations*, 2022.
- [12] Chengsong Huang, Qian Liu, Bill Yuchen Lin, Tianyu Pang, Chao Du, and Min Lin. LoraHub: Efficient cross-task generalization via dynamic lora composition. 2024.
- [13] Eyke Hüllermeier and Willem Waegeman. Aleatoric and epistemic uncertainty in machine learning: An introduction to concepts and methods. *Machine Learning* 110:457–506, 2021.
- [14] Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. PubMedQA: A dataset for biomedical research question answering. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2567–2577, Hong Kong, China, November 2019. Association for Computational Linguistics.
- [15] Ramneet Kaur, Susmit Jha, Anirban Roy, Sangdon Park, Oleg Sokolsky, and Insup Lee. Detecting odds as datapoints with high uncertainty. *CML 2021 Workshop on Uncertainty and Robustness in Deep Learning 2021*.
- [16] Leonie Koessler and Jonas Schuett. Risk assessment at agi companies: A review of popular risk assessment techniques from other safety-critical industries. 2023.
- [17] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017*, December 4-9, 2017, Long Beach, CA, USA, pages 6402–6413, 2017.
- [18] Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. *Advances in neural information processing systems*, 31, 2018.
- [19] Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicola Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018*, NeurIPS 2018, December 3-8, 2018, Montreal, Canada, pages 7167–7177, 2018.
- [20] Yinheng Li, Shaofei Wang, Han Ding, and Hang Chen. Large language models in finance: A survey. In *Proceedings of the Fourth ACM International Conference on AI in Finance*, pages 374–382, 2023.
- [21] Bryan F. J. Manly. *Multivariate Statistical Methods: A Primer*, Third Edition Chapman and Hall/CRC, New York, 3 edition, May 2014.
- [22] Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In Gerardo Flores, George H Chen, Tom Pollard, Joyce C Ho, and Tristan Naumann, editors, *Proceedings of the Conference on Health, In-*

- ference, and Learning volume 174 of Proceedings of Machine Learning Research pages 248–260. PMLR, 07–08 Apr 2022.
- [23] Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson. Fine-tuning aligned language models compromises safety, even when users do not intend to! 2023.
- [24] Jie Ren, Stanislav Fort, Jeremiah Liu, Abhijit Guha Roy, Shreyas Padhy, and Balaji Lakshminarayanan. A simple χ^2 to mahalanobis distance for improving near-ood detection arXiv preprint arXiv:2106.09022 2021.
- [25] Jie Ren, Jiaming Luo, Yao Zhao, Kundan Krishna, Mohammad Saleh, Balaji Lakshminarayanan, and Peter J. Liu. Out-of-distribution detection and selective generation for conditional language models. In The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023 OpenReview.net, 2023.
- [26] Jie Ren, Yao Zhao, Tu Vu, Peter J. Liu, and Balaji Lakshminarayanan. Self-evaluation improves selective generation in large language models arXiv, abs/2312.09300, 2023.
- [27] Viraj Shah, Nataniel Ruiz, Forrester Cole, Erika Lu, Svetlana Lazebnik, Yuanzhen Li, and Varun Jampani. Ziplora: Any subject in any style by effectively merging loras. 2023.
- [28] Ying Sheng, Shiyi Cao, Dacheng Li, Coleman Hooper, Nicholas Lee, Shuo Yang, Christopher Chou, Banghua Zhu, Lianmin Zheng, Kurt Keutzer, Joseph E. Gonzalez, and Ion Stoica. S-lora: Serving thousands of concurrent lora adapters. 2023.
- [29] Tobin South, Alexander Camuto, Shrey Jain, Shayla Nguyen, Robert Mahari, Christian Paquin, Jason Morton, and Alex 'Sandy' Pentland. Verifiable evaluations of machine learning models using zksnarks, 2024.
- [30] Yiyu Sun, Yifei Ming, Xiaojin Zhu, and Yixuan Li. Out-of-distribution detection with deep nearest neighbors. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato, editors, International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA volume 162 of Proceedings of Machine Learning Research pages 20827–20840. PMLR, 2022.
- [31] Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. Large language models in medicine. *Nature medicine* 29(8):1930–1940, 2023.
- [32] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and ne-tuned chat models. 2023.
- [33] Dustin Tran, Jeremiah Zhe Liu, Michael W Dusenberry, Du Phan, Mark Collier, Jie Ren, Kehang Han, Zi Wang, Zelda E Mariet, Huiyi Hu, Neil Band, Tim G. J. Rudner, Zachary Nado, Joost van Amersfoort, Andreas Kirsch, Rodolphe Jenatton, Nithum Thain, E. Kelly Buchanan, Kevin Patrick Murphy, D. Sculley, Yarin Gal, Zoubin Ghahramani, Jasper Snoek, and Balaji Lakshminarayanan. Plex: Towards reliability using pretrained large model extensions. *First Workshop on Pre-training: Perspectives, Pitfalls, and Paths Forward at ICML 2022* 2022.
- [34] Boxin Wang, Weixin Chen, Hengzhi Pei, Chulin Xie, Mintong Kang, Chenhui Zhang, Chejian Xu, Zidi Xiong, Ritik Dutta, Rylan Schaeffer, Sang T. Truong, Simran Arora, Mantas Mazeika, Dan Hendrycks, Zinan Lin, Yu Cheng, Sanmi Koyejo, Dawn Song, and Bo Li. Decodingtrust: A comprehensive assessment of trustworthiness in gpt models. 2024.
- [35] Haomiao Yang, Kunlan Xiang, Mengyu Ge, Hongwei Li, Rongxing Lu, and Shui Yu. A comprehensive overview of backdoor attacks in large language models within communication networks. 2023.
- [36] Jingkang Yang, Kaiyang Zhou, Yixuan Li, and Ziwei Liu. Generalized out-of-distribution detection: A survey. arXiv preprint arXiv:2110.11334 2021.
- [37] Jinghan Zhang, Shiqi Chen, Junteng Liu, and Junxian He. Composing parameter-efficient modules with arithmetic operations. *Advances in Neural Information Processing Systems* 2023.

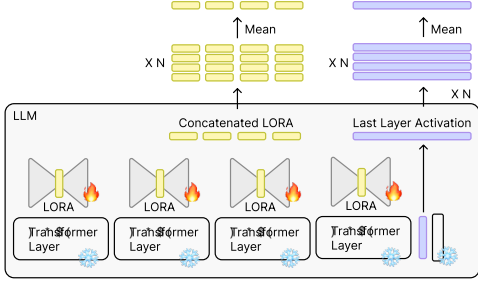


Figure A.1: **Embedding generation.** For each N input token, we collect a concatenation of the LoRA embeddings and the last layer activation. The final LoRA embeddings are an average of all the concatenations. For Last Layer Activation embeddings, the same averaging process is applied.

A. Appendix

A.1. Datasets (extended version)

We selected three Multiple Choice Question Answering (MCQA) datasets and chose the medical domain as a fine-tuning task, by considering the MedMCQA [22] and the PubMedQA [14] datasets. Then, we employed the MMLU [10] multi-domain dataset to define both near- and far-OD samples.

The MedMCQA dataset [22] contains around 194 000 multiple-choice questions, each with four options, derived from the Indian medical entrance exams (AIIMS and NEET). It includes 21 medical subjects and around 2 400 healthcare related topics.

The PubMedQA dataset [14] contains 1 000 expert-annotated and 211 300 artificially generated labelled Question Answering (QA) instances. The task involves generating a yes/no/maybe answers based on a context provided in the form of a PubMed abstract.

The MMLU dataset [10] includes questions from 57 different domains. As near-OD, we selected subtasks related to the medical domain such as “anatomy”, “clinical knowledge”, “college medicine”, “medical genetics”, “professional medicine”, and “college biology”. Conversely, as far-OD we picked: “professional law”, “international law”, “business ethics”, “computer security”, “college computer science”, “astronomy”, “abstract algebra” and “college chemistry”. These subtasks feature multiple-choice questions with four options and a known correct answer.

A.2. Related Work

OOD Detections Methods. Perplexity score assesses how effectively a language model predicts the next word in a sequence. Intuitively, a lower perplexity score suggests less uncertainty and better performance in prediction. Thus, it serves as an indicator of the input sequence’s proximity to the training dataset.

For Out-Of-Distribution (OOD) detection, if the embedding of a test input or output significantly deviates from the training data’s embedding distribution, it’s likely to be OOD. From the many existing OOD detection approaches [36], we selected three of them.

k -nearest neighbours (KNN) [30] computes the distance to the k -th nearest neighbours in the training set within a normalized embedding space, using this distance as an OOD indicator. Two key parameters, α and k , are involved. α determines the proportion of training data used for nearest neighbour calculations, and k specifies the particular nearest neighbour.

Mahalanobis distance (MD) [19] measures the OOD score by fitting a Gaussian, $\mathcal{N}(\mu; \Sigma)$; $\mu \in \mathbb{R}^d$; $\Sigma \in \mathbb{R}^{d \times d}$, to the training embeddings and using the Mahalanobis distance:

$$MD(x) := MD(x; \mu; \Sigma) := (x - \mu)^T \Sigma^{-1} (x - \mu) \quad (1)$$

Previous methods struggle to distinguish the fine-tuned task from other similar tasks. To address this, the Relative Mahalanobis Distance (RMD) [24] has been proposed to improve near-OD detection:

$$RMD_{\text{train}}(x) := MD_{\text{train}}(x) / MD_{bg}(x) \quad (2)$$

Where $MD_{bg}(x)$ is the Mahalanobis distance with regard to a background dataset, that in our scenario is the PubMedQA dataset.

OOD Detection in LLM. When it comes to LLMs, the RMD has been tested in [25] for conditional language models using the last layer activation of the encoder and decoder. They evaluated RMD with 2 experiments: distinguishing different version of newspaper summaries and types of translations as near-OD tasks. Note that this method involves using an additional background dataset bg , ideally encompassing the topics of near-OD tasks. The end goal is to further differentiate embeddings in the near-OD domain.

For an MCQA task, selective prediction have been tested by adding another answer “None Of the Above” [26], as well as training a classifier on top of the fine-tuned model [5].

A.3. Combining MD and Entropy

Given a question x , let $f_i(x)$ be the output confidence of an LLM for the i -th answer. Then, the Shannon entropy of the output is defined as:

$$H(x) = - \sum_i f_i(x) \log f_i(x) \quad (3)$$

Moreover, given a Mahalanobis distance MD for a data point x , the squared distance MD^2 follows a Chi-square (χ^2) distribution with d degrees of freedom [21], where d is the number of dimensions of the data point. The p -value associated with this Mahalanobis distance is calculated as follows:

$$p_{MD}(x) = 1 - \text{CDF}_{\chi^2, d}(MD^2(x)) \quad (4)$$

where $\text{CDF}_{\chi^2, d}$ represents the cumulative distribution function of the chi-square distribution with d degrees of freedom evaluated at the squared Mahalanobis distance MD^2 .

Last, given a question x with the associated LLM embeddings $E(x)$ (either E_{LLA} or E_{LORA}), we can compute the p -value p_{MD} for the Mahalanobis distance of the embeddings and the Shannon entropy $H(x)$ of the model prediction. The final combination is simply defined as: $H(x) + p_{MD}(E(x))$.

