# On the Effect of Uncertainty on Layer-wise Inference Dynamics

Sunwoo Kim [1]   Haneul Yoo [1]   Alice Oh [1]

## Abstract

Understanding how large language models (LLMs) internally represent and process their predictions is central to detecting uncertainty and preventing hallucinations. While several studies have shown that models encode uncertainty in their hidden states, it is underexplored how this affects the way they process such hidden states. In this work, we demonstrate that the dynamics of output token probabilities across layers for certain and uncertain outputs are largely aligned, revealing that uncertainty does not seem to affect inference dynamics. Specifically, we use the Tuned Lens, a variant of the Logit Lens, to analyze the layer-wise probability trajectories of final prediction tokens across 11 datasets and 5 models. Using incorrect predictions as those with higher epistemic uncertainty, our results show aligned trajectories for certain and uncertain predictions that both observe abrupt increases in confidence at similar layers. We balance this finding by showing evidence that more competent models may learn to process uncertainty differently. Our findings challenge the feasibility of leveraging simplistic methods for detecting uncertainty at inference. More broadly, our work demonstrates how interpretability methods may be used to investigate the way uncertainty affects inference.

## 1. Introduction

As capacities of LLMs grow, so does users' reliance on them. Hence, it is important to understand models' awareness of uncertainty and be able to detect it. To this end, many researchers have leveraged interpretability methods to study uncertainty calibration and quantification. Research using black-box methods have shown that models reflect uncertainty in their outputs (Kadavath et al., 2022; Wang et al.,

2024). Mechanistic interpretability (MI) methods, such as using sparse auto-encoders (Ferrando et al., 2024) and linear probes on embeddings of frozen, pretrained models (Ahdritz et al., 2024), have also found that models encode features for uncertainty in their internal representations. Although existing MI analyses provide ample evidence that *what* the model processes (*i.e.*, embeddings) is different for uncertain outputs, the way this affects *how* models process uncertain outputs differently is underexplored. Models could adapt their inference dynamics to uncertainty, perhaps processing outputs with greater uncertainty more, in terms of layers effectively used. If so, we may leverage these distinctly different inference dynamics for uncertainty detection. Additionally, it would seem natural that more competent models would exhibit more adaptive ways of processing uncertain predictions. Motivated by the thoughts above, we present two exploratory research questions:

- **RQ1:** How does uncertainty affect the inference dynamics of models? (*Figure 1, Figure 2*)

- **RQ2:** Does model competence affect the ability to adapt its inference dynamics to uncertainty? (*Figure 3*)

To answer the research questions, we analyze inference dynamics by using the levels of confidence the model has on possible prediction tokens throughout its layers. We extract such confidence levels from the hidden states passed in the residual stream between layers by employing the Tuned Lens (Belrose et al., 2023), a variant of the Logit Lens (nostalgebraist, 2020). We gather data across 11 datasets with 5 LLMs. We focus on *epistemic uncertainty* (Ahdritz et al., 2024; Hou et al., 2024), uncertainty due to lack of a model's knowledge and training, which is more easily observable compared to its counterpart, *aleatoric uncertainty*; the model exhibits epistemic uncertainty when it incorrectly answers a question from a set of answer choices that includes the correct answer.

Through our extensive experiments, we observe that probability trajectories are strikingly aligned and that models decide on their outputs at certain layers largely independent of uncertainty. This inference pattern means more complex and nuanced interpretability approaches may be required to extract model uncertainty at inference. While Jiang et al. (2024) examined inference dynamics in the con-

text of hallucination of known facts, our research extends this understanding across diverse tasks and models for epistemic uncertainty. Further, we add nuance to our findings by showing evidence that adaptive inference dynamics may arise with greater model competence.

## 2. Using Tuned Lens to Extract Inference Dynamics

We use the Tuned Lens (Belrose et al., 2023) to convert the hidden states in the residual stream after each layer to logit distributions in the vocabulary space. The logit distributions are indicative of what the model *believes* is the correct prediction after a layer. We perform this analysis on single token answers for multiple-choice question (MCQ) datasets. We analyze the distribution data to obtain tokens' probability trajectory through the layers, as well as the prediction depth—the layer at which the model commits to its answer.

### 2.1. Probability Trajectory Analysis

We analyze how the probabilities of possible predictions, namely single letter answer labels for the MCQs, change across layers. To clarify, we extract the probability distribution in the possible prediction space by performing softmax on just the logits for the label tokens. Then, we plot the probabilities for each label across layers to obtain the probability trajectories. We aggregate these probability trajectories across questions for all the datasets. Notably, we aggregate the trajectories for questions correctly and incorrectly answered separately. Specifically, for each question, we extract the trajectory for the label ranked first in terms of its probability at the final layer and obtain its average trajectory across questions. We condense the remaining lower ranked labels into a single trajectory. Hence, we obtain four average trajectories for each model, one pair of top and low rank trajectories for correctly answered questions and another for incorrectly answered questions. We compare the trajectories for correctly and incorrectly answered questions to see how inference dynamics differ.

### 2.2. Prediction Depth Analysis

The layer at which the model commits to its answer, or the prediction depth as coined by Baldock et al. (2021), would be the one at which a model's top prediction is different from the previous layer's and is maintained for subsequent layers. We calculate the PD for all questions in all datasets. We then plot the aggregate percentage distribution of prediction depths for all datasets across layers to see if distributions for when the model commits to its outputs diverge. Additionally, we calculate the Pearson correlation of answer incorrectness with PD for model and dataset pairs to see if uncertainty affects PD. Further, we analyze how the Kappa

value—accuracy adjusted for chance—on a dataset relates to the average PD difference between correct and incorrect outputs on that dataset. A greater PD difference would mean a greater degree of adaptive layer usage. Through this analysis, we can determine if sensitivity of inference dynamics to uncertainty changes according to task proficiency and, by implication, model competence. A positive correlation would mean that such abilities may emerge to greater degrees as a model's capability grows.

### 2.2.1. MODELS AND DATASETS

We use three models with Tuned Lenses provided by the authors of the Tuned Lens paper (Belrose et al., 2023): `Llama-3-8B`, `Llama-3-8B-Instruct` (Grattafiori et al., 2024), and `vicuna-13b-v1.1` (Chiang et al., 2023). Additionally, we train new Tuned Lenses on two models, `Mistral-7B-Instruct-v0.1` (Jiang et al., 2023) and `Mistral-Nemo-Instruct-2407` (Mistral AI, 2024).

We use datasets covering a range of tasks: ANLI (Nie et al., 2020), ARC (Clark et al., 2018), BoolQ (Clark et al., 2019), CommonsenseQA (Talmor et al., 2019), HellaSwag (Zellers et al., 2019), LogiQA (Liu et al., 2021), MMLU (Hendrycks et al., 2021), QASC (Khot et al., 2020), QuAIL (Rogers et al., 2020), RACE (Lai et al., 2017), and SciQ (Welbl et al., 2017).

## 3. Experimental Results

Figure 1 shows a consistent pattern across models in which the top prediction trajectory sharply increases and maintains its divergence with the lower rank trajectory. This is when the model seems to decide on its output. The trajectories for correct and incorrect predictions are strongly aligned, moving synchronously and abruptly increasing at the same layers. Albeit, the top prediction trajectory for the incorrect questions is consistently lower after the sharp increase, we cannot observe a definitive divergence between the two that would indicate different inference dynamics. We can also observe the alignment in the distribution of PD layers in Figure 2 where there are similar peaks, meaning the model committed to its final predictions at the same layers regardless of uncertainty.

PD correlations with answer incorrectness in Table 1 show weak positive correlations across models and datasets, with 80% below 0.300. However, the results are 97% positive. This consistency shows the potential for models to leverage adaptive inference techniques. To investigate this potential specifically in terms of emergence in competent models, we plot how the PD difference relates to the Kappa score for each dataset in Figure 3. The figure shows positive linear correlations between Kappa scores and the av-
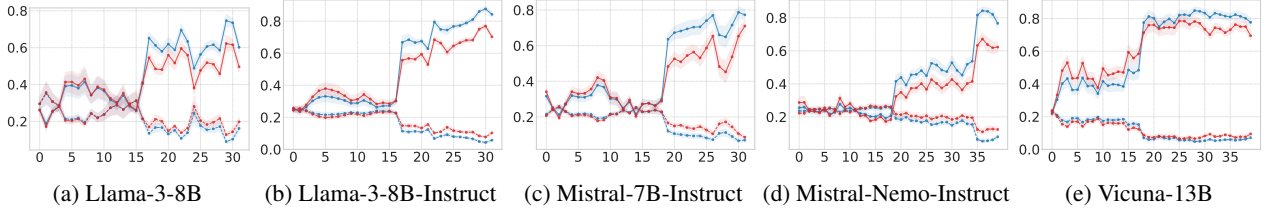
| (a) Llama-3-8B | (b) Llama-3-8B-Instruct | (c) Mistral-7B-Instruct | (d) Mistral-Nemo-Instruct | (e) Vicuna-13B |

*Figure 1.* Average probability trajectories. The x-axis denotes layer number and the y-axis denotes probability. The plots show how the final prediction token probability (*solid lines*) changes across layers for correct predictions (*in blue*) and incorrect predictions (*in red*). The probabilities for the other possible predictions in each case are condensed in the dashed line. The trajectories show strong alignment, moving synchronously across layers. The model seems to decide on its final prediction, as noted by the abrupt increase in probability, at similar layers as well.
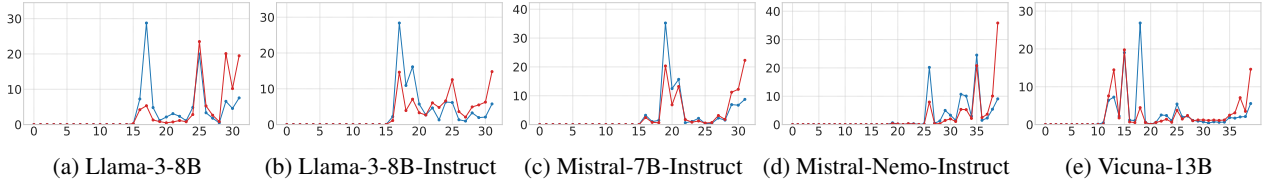


| (a) Llama-3-8B | (b) Llama-3-8B-Instruct | (c) Mistral-7B-Instruct | (d) Mistral-Nemo-Instruct | (e) Vicuna-13B |

*Figure 2.* Prediction depth distribution. The x-axis denotes layer number and the y-axis denotes percentage of questions. The plots show how the prediction depth, or the layer at which the model committed to its output, is distributed similarly for correct (*in blue*) and incorrect answers (*in red*). This alignment shows that the model reserves a similar number of layers for processing certain and uncertain outputs.
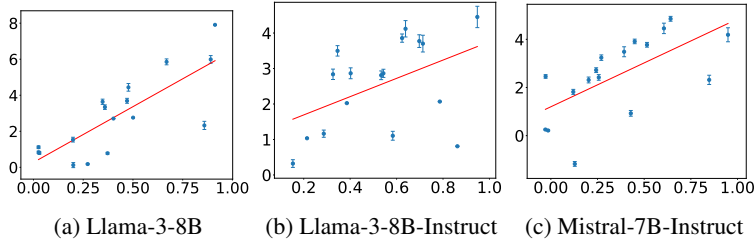


| (a) Llama-3-8B | (b) Llama-3-8B-Instruct | (c) Mistral-7B-Instruct |

*Figure 3.* Kappa vs. prediction depth difference. Each point in the plots represents a dataset; the y-axis denotes the average prediction depth difference, or difference in layers at which the model committed to answers, between correctly and incorrectly answered questions and the x-axis denotes Cohen's Kappa score, or accuracy adjusted for random chance. The positive correlations across the models show that as task competence increases, represented by the Kappa score, the degree to which uncertainty affects the inference dynamics increases, represented by the prediction depth difference. (`Llama-3-8B` and `Mistral-7B-Instruct` have $p < 0.05$ while `Llama-3-8B-Instruct` has $p = 0.06$. Other larger models do not observe significant trends.)

erage PD gap for datasets. Two out of the five models, `Llama-3-8B` and `Mistral-7B-Instruct`, show statistically significant positive correlations with $p < 0.05$ while `Llama-3-8B-Instruct` has $p = 0.06$. Insofar as we can interpret accuracy on a dataset as task proficiency, these results are preliminary evidence that stronger adaptive inference dynamics could emerge with greater competence.

## 4. Related Work

### 4.1. Logit Lens and Tuned Lens

To analyze the inference dynamics of models, we leverage Tuned Lens (Belrose et al., 2023), a variant of the method of the Logit Lens (nostalgebraist, 2020). The original Logit Lens is an early-stage technique that decodes intermediate activations into vocabulary space using the classification head. It has been leveraged to gain insights into how LLMs process and generate languages (Wang, 2025), including analysis on reasoning pathways (Li et al., 2024), multilingual representations (Wendler et al., 2024), and safety alignment (Golgoon et al., 2024). Building upon this, Belrose et al. (2023) proposed Tuned Lens, a refinement of Logit Lens, which trains an affine probe for each block in a frozen pre-trained model. This method offers more reliable interpretations of model predictions across layers.

*Table 1.* Pearson correlation between answer incorrectness and PD. Greater values would mean that incorrectness correlates with later prediction commitment layers. Correlations are generally weak, with 80% below 0.300, but with 97% being positive. This consistency shows LLMs' potential to leverage uncertainty in their inference dynamics. (Values in parentheses show standard error rates. Values above 0.300 are bold.)

| DATASET | LLAMA-3-8B | LLAMA-3-8B-INSTRUCT | VICUNA-13B | MISTRAL-7B-INSTRUCT | MISTRAL-NEMO-INSTRUCT |
|---|---|---|---|---|---|
| ANLI-R1 | 0.192*(0.015) | 0.266*(0.013) | 0.195*(0.014) | 0.118*(0.014) | **0.428**\*(0.012) |
| ANLI-R2 | 0.174*(0.016) | **0.321**\*(0.013) | 0.194*(0.014) | 0.146*(0.014) | **0.464**\*(0.011) |
| ANLI-R3 | 0.238*(0.015) | 0.258*(0.013) | 0.215*(0.013) | 0.284*(0.009) | 0.240*(0.013) |
| ARC-EASY | **0.449**\*(0.011) | **0.351**\*(0.019) | 0.069*(0.021) | **0.443**\*(0.017) | **0.332**\*(0.019) |
| ARC-CHALLENGE | **0.373**\*(0.017) | **0.363**\*(0.017) | 0.040(0.030) | **0.322**\*(0.018) | **0.357**\*(0.017) |
| BOOLQ | 0.061*(0.012) | 0.159*(0.014) | 0.058*(0.014) | -0.168*(0.014) | -0.083*(0.014) |
| BOOLQ W/ CONTEXT | 0.175*(0.011) | 0.132*(0.014) | 0.233*(0.013) | 0.107*(0.014) | -0.081*(0.014) |
| COMMONSENSEQA | **0.316**\*(0.011) | 0.255*(0.009) | -0.045*(0.010) | 0.274*(0.009) | 0.258*(0.010) |
| HELLASWAG | 0.010(0.012) | 0.047*(0.016) | 0.188*(0.014) | 0.261*(0.013) | **0.309**\*(0.013) |
| LOGIQA | 0.174*(0.014) | 0.120*(0.011) | 0.113*(0.014) | 0.216*(0.011) | 0.269*(0.015) |
| MMLU | **0.332**\*(0.013) | 0.217*(0.008) | 0.139*(0.014) | 0.265*(0.009) | **0.364**\*(0.013) |
| QASC | 0.287*(0.010) | 0.252*(0.010) | 0.129*(0.011) | 0.213*(0.011) | 0.248*(0.010) |
| QASC W/ CONTEXT | **0.374**\*(0.010) | 0.263*(0.010) | **0.408**\*(0.009) | 0.298*(0.010) | 0.164*(0.011) |
| QUAIL | **0.318**\*(0.013) | **0.337**\*(0.009) | 0.276*(0.009) | **0.381**\*(0.009) | **0.334**\*(0.013) |
| RACE | **0.315**\*(0.013) | **0.327**\*(0.013) | 0.247*(0.013) | **0.370**\*(0.009) | **0.396**\*(0.012) |
| SCIQ | **0.480**\*(0.007) | 0.213*(0.009) | 0.123*(0.014) | **0.484**\*(0.008) | 0.101*(0.014) |
| SCIQ W/ CONTEXT | 0.217*(0.014) | 0.083*(0.009) | 0.192*(0.014) | 0.286*(0.013) | 0.065*(0.014) |

$^*p < 0.05$

## 4.2. Prediction Depth

Prediction depth is the earliest layer in a deep learning model where the highest-probability prediction matches the final output layer's prediction, differs from the previous layer's prediction, and remains the same through subsequent layers. It can be interpreted as the commitment layer or the effective number of layers used for inference. Baldock et al. (2021) first defined and validated it as a measure of computational difficulty for individual examples in computer vision tasks. However, there is limited empirical research on PD in LLMs. While Belrose et al. (2023) has tested how effectively PD correlates with iteration learned on a single model, we analyze it with relation to epistemic uncertainty across a variety of models.

## 4.3. Inference Dynamics

The term "*inference dynamics*" as used in this work is an umbrella term that refers to the patterns that arise at inference, namely prediction probability trajectories and prediction refinement across layers. Studies in this field have found conflicting results for how a model produces its final output, with some results showing a gradual refinement across layers while others have found layers disproportionately contributing to the output and inducing sudden "*decisions*" (nostalgebraist, 2020; Kongmanee, 2025). While the research is inconclusive, methods for early-exiting still try to leverage the knowledge we have on inference dynamics to reduce computational cost by using predictions at intermediate layers (Xin et al., 2020; Schuster et al., 2022; Varshney et al., 2024). Our results related to PD and prediction trajec-

tories would help advance these endeavors.

## 5. Conclusion

In this paper, we analyze if epistemic uncertainty causes large language models (LLMs) to exhibit different inference dynamics. Using a Logit Lens-based method, we examine layer-wise token probability trajectories and prediction depth (PD) statistics of 5 LLMs across 11 datasets. We demonstrate that LLMs exhibit similar token probability dynamics and tend to commit to final decisions at specific layers largely independent of output uncertainty. This implies that inference is generally characterized by abrupt decision-making unaffected by levels of uncertainty. In addition, we observe that models exhibit the potential for greater adaptability in their inference dynamics as PD shows weak but consistent correlations with epistemic uncertainty. Further, models exhibit greater degrees of inference adaptability with respect to uncertainty at tasks they are more proficient at, implying greater degrees of adaptability could emerge in more competent models. These findings have implications for interpretability methods assessing uncertainty at inference as well as for adaptive inference techniques (*i.e.*, test-time scaling, early exiting) that leverage intermediate prediction confidence. Further testing should be done on more models to validate the results of this exploratory work, especially on more competent models. More rigorous tests that include analyses using varying levels of uncertainty and accounts for hallucination on known knowledge is needed as well. Additionally, similar experiments for aleatoric uncertainty may be done for a more thorough investigation.

## Acknowledgements

## References

Ahdritz, G., Qin, T., Vyas, N., Barak, B., and Edelman, B. L. Distinguishing the knowable from the unknowable with language models. In *Proceedings of the 41st International Conference on Machine Learning*, ICML'24. JMLR.org, 2024.

Baldock, R., Maennel, H., and Neyshabur, B. Deep learning through the lens of example difficulty. In Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 10876–10889. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/5a4b25aaed25c2ee1b74de72dc03c14e-Paper.pdf.

Belrose, N., Furman, Z., Smith, L., Halawi, D., Ostrovsky, I., McKinney, L., Biderman, S., and Steinhardt, J. Eliciting latent predictions from transformers with the tuned lens. *arXiv preprint arXiv:2303.08112*, 2023. URL https://arxiv.org/abs/2303.08112.

Chiang, W.-L., Li, Z., Lin, Z., Sheng, Y., Wu, Z., Zhang, H., Zheng, L., Zhuang, S., Zhuang, Y., Gonzalez, J. E., Stoica, I., and Xing, E. P. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL https://lmsys.org/blog/2023-03-30-vicuna/.

Clark, C., Lee, K., Chang, M.-W., Kwiatkowski, T., Collins, M., and Toutanova, K. BoolQ: Exploring the surprising difficulty of natural yes/no questions. In Burstein, J., Doran, C., and Solorio, T. (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 2924–2936, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1300. URL https://aclanthology.org/N19-1300/.

Clark, P., Cowhey, I., Etzioni, O., Khot, T., Sabharwal, A., Schoenick, C., and Tafjord, O. Think you have solved question answering? try arc, the ai2 reasoning challenge, 2018. URL https://arxiv.org/abs/1803.05457.

Ferrando, J., Obeso, O., Rajamanoharan, S., and Nanda, N. Do i know this entity? knowledge awareness and hallucinations in language models. *arXiv preprint arXiv:2411.14257*, 2024. URL https://arxiv.org/abs/2411.14257.

Golgoon, A., Filom, K., and Ravi Kannan, A. Mechanistic interpretability of large language models with applications to the financial services industry. In *Proceedings of the 5th ACM International Conference on AI in Finance*, ICAIF '24, pp. 660–668, New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400710810. doi: 10.1145/3677052.3698612. URL https://doi.org/10.1145/3677052.3698612.

Grattafiori, A., Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Vaughan, A., Yang, A., Fan, A., Goyal, A., Hartshorn, A., Yang, A., Mitra, A., Sravankumar, A., Korenev, A., Hinsvark, A., Rao, A., Zhang, A., Rodriguez, A., Gregerson, A., Spataru, A., Roziere, B., Biron, B., Tang, B., Chern, B., Caucheteux, C., Nayak, C., Bi, C., Marra, C., McConnell, C., Keller, C., Touret, C., Wu, C., Wong, C., Ferrer, C. C., Nikolaidis, C., Allonsius, D., Song, D., Pintz, D., Livshits, D., Wyatt, D., Esiobu, D., Choudhary, D., Mahajan, D., Garcia-Olano, D., Perino, D., Hupkes, D., Lakomkin, E., AlBadawy, E., Lobanova, E., Dinan, E., Smith, E. M., Radenovic, F., Guzmán, F., Zhang, F., Synnaeve, G., Lee, G., Anderson, G. L., Thattai, G., Nail, G., Mialon, G., Pang, G., Cucurell, G., Nguyen, H., Korevaar, H., Xu, H., Touvron, H., Zarov, I., Ibarra, I. A., Kloumann, I., Misra, I., Evtimov, I., Zhang, J., Copet, J., Lee, J., Geffert, J., Vranes, J., Park, J., Mahadeokar, J., Shah, J., van der Linde, J., Billock, J., Hong, J., Lee, J., Fu, J., Chi, J., Huang, J., Liu, J., Wang, J., Yu, J., Bitton, J., Spisak, J., Park, J., Rocca, J., Johnstun, J., Saxe, J., Jia, J., Alwala, K. V., Prasad, K., Upasani, K., Plawiak, K., Li, K., Heafield, K., Stone, K., El-Arini, K., Iyer, K., Malik, K., Chiu, K., Bhalla, K., Lakhotia, K., Rantala-Yeary, L., van der Maaten, L., Chen, L., Tan, L., Jenkins, L., Martin, L., Madaan, L., Malo, L., Blecher, L., Landzaat, L., de Oliveira, L., Muzzi, M., Pasupuleti, M., Singh, M., Paluri, M., Kardas, M., Tsimpoukelli, M., Oldham, M., Rita, M., Pavlova, M., Kambadur, M., Lewis, M., Si, M., Singh, M. K., Hassan, M., Goyal, N., Torabi, N., Bashlykov, N., Bogoychev, N., Chatterji, N., Zhang, N., Duchenne, O., Çelebi, O., Alrassy, P., Zhang, P., Li, P., Vasic, P., Weng, P., Bhargava, P., Dubal, P., Krishnan,

P., Koura, P. S., Xu, P., He, Q., Dong, Q., Srinivasan, R., Ganapathy, R., Calderer, R., Cabral, R. S., Stojnic, R., Raileanu, R., Maheswari, R., Girdhar, R., Patel, R., Sauvestre, R., Polidoro, R., Sumbaly, R., Taylor, R., Silva, R., Hou, R., Wang, R., Hosseini, S., Chennabasappa, S., Singh, S., Bell, S., Kim, S. S., Edunov, S., Nie, S., Narang, S., Raparthy, S., Shen, S., Wan, S., Bhosale, S., Zhang, S., Vandenhende, S., Batra, S., Whitman, S., Sootla, S., Collot, S., Gururangan, S., Borodinsky, S., Herman, T., Fowler, T., Sheasha, T., Georgiou, T., Scialom, T., Speck-bacher, T., Mihaylov, T., Xiao, T., Karn, U., Goswami, V., Gupta, V., Ramanathan, V., Kerkez, V., Gonguet, V., Do, V., Vogeti, V., Albiero, V., Petrovic, V., Chu, W., Xiong, W., Fu, W., Meers, W., Martinet, X., Wang, X., Wang, X., Tan, X. E., Xia, X., Xie, X., Jia, X., Wang, X., Goldschlag, Y., Gaur, Y., Babaei, Y., Wen, Y., Song, Y., Zhang, Y., Li, Y., Mao, Y., Coudert, Z. D., Yan, Z., Chen, Z., Papakipos, Z., Singh, A., Srivastava, A., Jain, A., Kelsey, A., Shajn-feld, A., Gangidi, A., Victoria, A., Goldstand, A., Menon, A., Sharma, A., Boesenberg, A., Baevski, A., Feinstein, A., Kallet, A., Sangani, A., Teo, A., Yunus, A., Lupu, A., Alvarado, A., Caples, A., Gu, A., Ho, A., Poulton, A., Ryan, A., Ramchandani, A., Dong, A., Franco, A., Goyal, A., Saraf, A., Chowdhury, A., Gabriel, A., Bharambe, A., Eisenman, A., Yazdan, A., James, B., Maurer, B., Leonhardi, B., Huang, B., Loyd, B., Paola, B. D., Paran-jape, B., Liu, B., Wu, B., Ni, B., Hancock, B., Wasti, B., Spence, B., Stojkovic, B., Gamido, B., Montalvo, B., Parker, C., Burton, C., Mejia, C., Liu, C., Wang, C., Kim, C., Zhou, C., Hu, C., Chu, C.-H., Cai, C., Tindal, C., Feichtenhofer, C., Gao, C., Civin, D., Beaty, D., Kreymer, D., Li, D., Adkins, D., Xu, D., Testuggine, D., David, D., Parikh, D., Liskovich, D., Foss, D., Wang, D., Le, D., Holland, D., Dowling, E., Jamil, E., Montgomery, E., Presani, E., Hahn, E., Wood, E., Le, E.-T., Brinkman, E., Arcaute, E., Dunbar, E., Smothers, E., Sun, F., Kreuk, F., Tian, F., Kokkinos, F., Ozgenel, F., Caggioni, F., Kanayet, F., Seide, F., Florez, G. M., Schwarz, G., Badeer, G., Swee, G., Halpern, G., Herman, G., Sizov, G., Guangyi, Zhang, Lakshminarayanan, G., Inan, H., Shojanazeri, H., Zou, H., Wang, H., Zha, H., Habeeb, H., Rudolph, H., Suk, H., Aspegren, H., Goldman, H., Zhan, H., Damlaj, I., Molybog, I., Tufanov, I., Leontiadis, I., Veliche, I.-E., Gat, I., Weissman, J., Geboski, J., Kohli, J., Lam, J., Asher, J., Gaya, J.-B., Marcus, J., Tang, J., Chan, J., Zhen, J., Reizenstein, J., Teboul, J., Zhong, J., Jin, J., Yang, J., Cummings, J., Carvill, J., Shepard, J., McPhie, J., Tor-res, J., Ginsburg, J., Wang, J., Wu, K., U, K. H., Saxena, K., Khandelwal, K., Zand, K., Matosich, K., Veeraragha-van, K., Michelena, K., Li, K., Jagadeesh, K., Huang, K., Chawla, K., Huang, K., Chen, L., Garg, L., A, L., Silva, L., Bell, L., Zhang, L., Guo, L., Yu, L., Moshkovich, L., Wehrstedt, L., Khabsa, M., Avalani, M., Bhatt, M., Mankus, M., Hasson, M., Lennie, M., Reso, M., Gro-shev, M., Naumov, M., Lathi, M., Keneally, M., Liu, M., Seltzer, M. L., Valko, M., Restrepo, M., Patel, M., Vy-atskov, M., Samvelyan, M., Clark, M., Macey, M., Wang, M., Hermoso, M. J., Metanat, M., Rastegari, M., Bansal, M., Santhanam, N., Parks, N., White, N., Bawa, N., Sing-hal, N., Egebo, N., Usunier, N., Mehta, N., Laptev, N. P., Dong, N., Cheng, N., Chernoguz, O., Hart, O., Salpekar, O., Kalinli, O., Kent, P., Parekh, P., Saab, P., Balaji, P., Rittner, P., Bontrager, P., Roux, P., Dollar, P., Zvyagina, P., Ratanchandani, P., Yuvraj, P., Liang, Q., Alao, R., Rodriguez, R., Ayub, R., Murthy, R., Nayani, R., Mitra, R., Parthasarathy, R., Li, R., Hogan, R., Battey, R., Wang, R., Howes, R., Rinott, R., Mehta, S., Siby, S., Bondu, S. J., Datta, S., Chugh, S., Hunt, S., Dhillon, S., Sidorov, S., Pan, S., Mahajan, S., Verma, S., Yamamoto, S., Ra-maswamy, S., Lindsay, S., Lindsay, S., Feng, S., Lin, S., Zha, S. C., Patil, S., Shankar, S., Zhang, S., Zhang, S., Wang, S., Agarwal, S., Sajuyigbe, S., Chintala, S., Max, S., Chen, S., Kehoe, S., Satterfield, S., Govindaprasad, S., Gupta, S., Deng, S., Cho, S., Virk, S., Subramanian, S., Choudhury, S., Goldman, S., Remez, T., Glaser, T., Best, T., Koehler, T., Robinson, T., Li, T., Zhang, T., Matthews, T., Chou, T., Shaked, T., Vontimitta, V., Ajayi, V., Montanez, V., Mohan, V., Kumar, V. S., Mangla, V., Ionescu, V., Poenaru, V., Mihailescu, V. T., Ivanov, V., Li, W., Wang, W., Jiang, W., Bouaziz, W., Constable, W., Tang, X., Wu, X., Wang, X., Wu, X., Gao, X., Kleinman, Y., Chen, Y., Hu, Y., Jia, Y., Qi, Y., Li, Y., Zhang, Y., Zhang, Y., Adi, Y., Nam, Y., Yu, Wang, Zhao, Y., Hao, Y., Qian, Y., Li, Y., He, Y., Rait, Z., DeVito, Z., Rosnbrick, Z., Wen, Z., Yang, Z., Zhao, Z., and Ma, Z. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. URL https://arxiv.org/abs/2407.21783.

Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. Measuring massive multitask language understanding. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=d7KBjmI3GmQ.

Hou, B., Liu, Y., Qian, K., Andreas, J., Chang, S., and Zhang, Y. Decomposing uncertainty for large language models through input clarification ensembling. In *Proceedings of the 41st International Conference on Machine Learning*, ICML'24. JMLR.org, 2024.

Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., de las Casas, D., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., Lavaud, L. R., Lachaux, M.-A., Stock, P., Scao, T. L., Lavril, T., Wang, T., Lacroix, T., and Sayed, W. E. Mistral 7b, 2023. URL https://arxiv.org/abs/2310.06825.

Jiang, C., Qi, B., Hong, X., Fu, D., Cheng, Y., Meng, F., Yu, M., Zhou, B., and Zhou, J. On large lan-guage models' hallucination with regard to known facts.

In Duh, K., Gomez, H., and Bethard, S. (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 1041–1053, Mexico City, Mexico, June 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.naacl-long.60. URL https://aclanthology.org/2024.naacl-long.60/.

Kadavath, S., Conerly, T., Askell, A., Henighan, T., Drain, D., Perez, E., Schiefer, N., Hatfield-Dodds, Z., Das-Sarma, N., Tran-Johnson, E., et al. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*, 2022. URL https://arxiv.org/abs/2207.05221.

Khot, T., Clark, P., Guerquin, M., Jansen, P., and Sabharwal, A. Qasc: A dataset for question answering via sentence composition. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05): 8082–8090, Apr. 2020. doi: 10.1609/aaai.v34i05.6319. URL https://ojs.aaai.org/index.php/AAAI/article/view/6319.

Kongmanee, J. An attempt to unraveling token prediction refinement and identifying essential layers of large language models, 2025. URL https://arxiv.org/abs/2501.15054.

Lai, G., Xie, Q., Liu, H., Yang, Y., and Hovy, E. RACE: Large-scale ReAding comprehension dataset from examinations. In Palmer, M., Hwa, R., and Riedel, S. (eds.), *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pp. 785–794, Copenhagen, Denmark, September 2017. Association for Computational Linguistics. doi: 10.18653/v1/D17-1082. URL https://aclanthology.org/D17-1082/.

Li, Z., Jiang, G., Xie, H., Song, L., Lian, D., and Wei, Y. Understanding and patching compositional reasoning in LLMs. In Ku, L.-W., Martins, A., and Srikumar, V. (eds.), *Findings of the Association for Computational Linguistics: ACL 2024*, pp. 9668–9688, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.576. URL https://aclanthology.org/2024.findings-acl.576/.

Liu, J., Cui, L., Liu, H., Huang, D., Wang, Y., and Zhang, Y. Logiqa: a challenge dataset for machine reading comprehension with logical reasoning. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence*, IJCAI'20, 2021. ISBN 9780999241165.

Merity, S., Xiong, C., Bradbury, J., and Socher, R. Pointer sentinel mixture models. In *International Conference on Learning Representations*, 2017. URL https://openreview.net/forum?id=Byj72udxe.

Mistral AI. Mistral-nemo-instruct-2407. https://huggingface.co/mistralai/Mistral-Nemo-Instruct-2407, 2024. Accessed: 2025-05-18.

Nie, Y., Williams, A., Dinan, E., Bansal, M., Weston, J., and Kiela, D. Adversarial NLI: A new benchmark for natural language understanding. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, 2020.

nostalgebraist. Interpreting gpt: the logit lens, 2020. URL https://www.lesswrong.com/posts/AcKRB8wDpdaN6v6ru/interpreting-gpt-the-logit-lens.

Rogers, A., Kovaleva, O., Downey, M., and Rumshisky, A. Getting closer to ai complete question answering: A set of prerequisite real tasks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34 (05):8722–8731, Apr. 2020. doi: 10.1609/aaai.v34i05.6398. URL https://ojs.aaai.org/index.php/AAAI/article/view/6398.

Schuster, T., Fisch, A., Gupta, J., Dehghani, M., Bahri, D., Tran, V., Tay, Y., and Metzler, D. Confident adaptive language modeling. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 17456–17472. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/6fac9e316a4ae75ea244ddcef1982c71-Paper-Conference.pdf.

Talmor, A., Herzig, J., Lourie, N., and Berant, J. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In Burstein, J., Doran, C., and Solorio, T. (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4149–4158, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1421. URL https://aclanthology.org/N19-1421/.

Varshney, N., Chatterjee, A., Parmar, M., and Baral, C. Investigating acceleration of LLaMA inference by enabling intermediate layer decoding via instruction tuning with 'LITE'. In Duh, K., Gomez, H., and Bethard, S. (eds.), *Findings of the Association for Computational Linguistics: NAACL 2024*, pp. 3656–3677, Mexico City, Mexico, June 2024. Association for Computa-

tional Linguistics. doi: 10.18653/v1/2024.findings-naacl.232. URL https://aclanthology.org/2024.findings-naacl.232/.

Wang, X., Zhang, Z., Li, Q., Chen, G., Hu, M., li, Z., Luo, B., Gao, H., Han, Z., and Wang, H. Ubench: Benchmarking uncertainty in large language models with multiple choice questions, 2024. URL https://arxiv.org/abs/2406.12784.

Wang, Z. Logitlens4llms: Extending logit lens analysis to modern large language models. *arXiv preprint arXiv:2503.11667*, 2025. URL https://arxiv.org/abs/2503.11667.

Welbl, J., Liu, N. F., and Gardner, M. Crowdsourcing multiple choice science questions. In Derczynski, L., Xu, W., Ritter, A., and Baldwin, T. (eds.), *Proceedings of the 3rd Workshop on Noisy User-generated Text*, pp. 94–106, Copenhagen, Denmark, September 2017. Association for Computational Linguistics. doi: 10.18653/v1/W17-4413. URL https://aclanthology.org/W17-4413/.

Wendler, C., Veselovsky, V., Monea, G., and West, R. Do llamas work in English? on the latent language of multilingual transformers. In Ku, L.-W., Martins, A., and Srikumar, V. (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15366–15394, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.820. URL https://aclanthology.org/2024.acl-long.820/.

Xin, J., Tang, R., Lee, J., Yu, Y., and Lin, J. DeeBERT: Dynamic early exiting for accelerating BERT inference. In Jurafsky, D., Chai, J., Schluter, N., and Tetreault, J. (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 2246–2251, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.204. URL https://aclanthology.org/2020.acl-main.204/.

Zellers, R., Holtzman, A., Bisk, Y., Farhadi, A., and Choi, Y. HellaSwag: Can a machine really finish your sentence? In Korhonen, A., Traum, D., and Màrquez, L. (eds.), *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4791–4800, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1472. URL https://aclanthology.org/P19-1472/.

# Appendix

## A. Broader Impacts Statement

This paper presents work whose goal is to advance the field of machine learning. There are many potential societal consequences of our work, none of which we feel must be specifically highlighted here.

## B. Reproducibility Statement

### B.1. Prompts

Question prompts from the original MCQ datasets were augmented with directions to answer with a single letter. If supplementary context that is not necessarily needed to answer the question is provided, a separate version of the questions with contexts were also tested. Note the use of brackets around the label in the directions and after the '*Answer:* ' to prompt the model to output a letter label. The format is as following:

> Answer the question with a single letter like [A].
> Mesophiles grow best in moderate temperature, typically between 25°C and 40°C (77°F and 104°F). Mesophiles are often found living in or on the bodies of humans or other animals. The optimal growth temperature of many pathogenic mesophiles is 37°C (98°F), the normal human body temperature. Mesophilic organisms have important uses in food preparation, including cheese, yogurt, beer and wine.
> What type of organism is commonly used in preparation of foods such as cheese and yogurt?
> A. viruses
> B. protozoa
> C. gymnosperms
> D. mesophilic organisms
> Answer: [

### B.2. Sample Size

The model did not produce relevant answer labels for every question. For example, it sometimes tried to answer with a sentence such as "*The answer is · · · .*" Hence, we consider only the questions that were answered with a single label token. The filtering results in sample sizes that are different across model and dataset pairs as is shown in Table 2.

Table 2. Number of sensical answers for model and dataset pairs.

| DATASET | LLAMA-3-8B | LLAMA-3-8B-INSTRUCT | VICUNA-13B | MISTRAL-7B-INSTRUCT | MISTRAL-NEMO-INSTRUCT |
|---|---|---|---|---|---|
| ANLI-R1 | 3909 | 5000 | 5000 | 5000 | 4996 |
| ANLI-R2 | 3787 | 5000 | 4998 | 5000 | 5000 |
| ANLI-R3 | 4037 | 5000 | 5000 | 10000 | 5000 |
| ARC-EASY | 4935 | 2241 | 2206 | 2151 | 2237 |
| ARC-CHALLENGE | 2556 | 2567 | 1116 | 2557 | 2531 |
| BOOLQ | 7479 | 5000 | 5000 | 4999 | 5000 |
| BOOLQ W/ CONTEXT | 7481 | 5000 | 5000 | 5000 | 5000 |
| COMMONSENSEQA | 6372 | 9737 | 9741 | 9706 | 9325 |
| HELLASWAG | 7394 | 3772 | 4988 | 5000 | 4720 |
| LOGIQA | 4687 | 7375 | 4996 | 7373 | 3745 |
| MMLU | 4969 | 14032 | 4991 | 9822 | 4742 |
| QASC | 7952 | 8122 | 8134 | 8112 | 8062 |
| QASC W/ CONTEXT | 7917 | 8134 | 8134 | 8132 | 7998 |
| QUAIL | 4996 | 9999 | 10246 | 10000 | 5000 |
| RACE | 4981 | 4998 | 5000 | 9999 | 4975 |
| SCIQ | 11679 | 11679 | 5000 | 10000 | 4996 |
| SCIQ W/ CONTEXT | 4972 | 11679 | 5000 | 4999 | 4992 |

## B.3. Tuned Lens Training

We trained Tuned Lens for the `Mistral-7B-Instruct-v0.1` and `Mistral-Nemo-Instruct-2407` using WikiText-103-dataset (Merity et al., 2017), following the original implementations on GitHub [1]. We use Intel(R) Xeon(R) Silver 4214R CPU @ 2.40GHz (24 cores), 188GiB System memory, and four 48GB NVIDIA Quadro RTX 8000 GPUs. Table 3 shows datasets and hyperparameters for training Tuned Lenses.

*Table 3.* Training details for Tuned Lens

| HYPERPARAMETERS | VALUES |
|---|---|
| OPTIMZIER | ADAM |
| LEARNING RATE | 0.001 |
| TRAINING STEPS | 1000 |
| WEIGHT DECAY | 0.01 |
| MOMENTUM | 0.9 |
| TOKENS PER STEP | $2^{18}$ |

---

[1] https://github.com/AlignmentResearch/tuned-lens