













# Your Finetuned Large Language Model is Already a Powerful OOD Detector



Andi Zhang, Tim Z. Xiao, Weiyang Liu, Robert Bamler, Damon Wischik

#### **Takeaway**



- $\triangleright$   $\theta$ : parameters of the pre-trained LLM
- $\triangleright$   $\theta'$ : parameters of the fine-tuned LLM on a specific domain

 $\frac{p_{\theta}(x)}{p_{\theta'}(x)}$  : powerful OOD detector

## **Takeaway**



- $\triangleright$   $\theta$ : parameters of the pre-trained LLM
- $\triangleright$   $\theta'$ : parameters of the fine-tuned LLM on a specific domain

```
\frac{p_{\theta}(x)}{p_{\theta'}(x)} : powerful OOD detector
```

#### Implementation in 3 lines:

```
pre_log_likelihood = -pretrained_model(x).loss * len(x)
fine_log_likelihood = -finetuned_model(x).loss * len(x)
log_ratio = pre_log_likelihood - fine_log_likelihood
```

Note: loss of autoregressive LLMs = average NLL per token (perplexity).







► Data come from a true distribution p<sub>true</sub>



- Data come from a true distribution p<sub>true</sub>
- ightharpoonup We only observe in-distribution data ightharpoonup assume they follow  $p_{\rm in}$



- Data come from a true distribution p<sub>true</sub>
- We only observe in-distribution data  $\Rightarrow$  assume they follow  $p_{in}$
- ► Extra assumption (not typical): OOD data also follow some distribution p<sub>out</sub>



- Data come from a true distribution p<sub>true</sub>
- We only observe in-distribution data  $\Rightarrow$  assume they follow  $p_{in}$
- Extra assumption (not typical): OOD data also follow some distribution p<sub>out</sub>

Neyman–Pearson lemma: Given x, the optimal test is

$$\frac{p_{\text{out}}(x)}{p_{\text{in}}(x)}$$

## **OOD Proxy**



► The true  $p_{OOD}$  is unknown

<sup>&</sup>lt;sup>1</sup>Andi Zhang and Damon Wischik. "Falsehoods that ML researchers believe about OOD detection". In: arXiv preprint arXiv:2210.12767 (2022).

<sup>&</sup>lt;sup>2</sup>Christopher M Bishop. "Novelty detection and neural network validation". In: *IEE Proceedings-Vision, Image and Signal processing* 141.4 (1994), pp. 217–222.

# **OOD Proxy**



- ightharpoonup The true  $p_{OOD}$  is unknown
- ► Zhang & Wischik¹: use a **proxy distribution**  $p_{OOD}^{proxy}$  to approximate the OOD distribution

 $\frac{p_{\text{OOD}}^{\text{proxy}}(x)}{p_{\text{in}}(x)}$ 

<sup>&</sup>lt;sup>1</sup>Andi Zhang and Damon Wischik. "Falsehoods that ML researchers believe about OOD detection". In: arXiv preprint arXiv:2210.12767 (2022).

<sup>&</sup>lt;sup>2</sup>Christopher M Bishop. "Novelty detection and neural network validation". In: *IEE Proceedings-Vision, Image and Signal processing* 141.4 (1994), pp. 217–222.

## **OOD Proxy**



- ightharpoonup The true  $p_{OOD}$  is unknown
- Zhang & Wischik<sup>1</sup>: use a proxy distribution p<sub>OOD</sub><sup>proxy</sup> to approximate the OOD distribution

$$\frac{p_{\text{OOD}}^{\text{proxy}}(x)}{p_{\text{in}}(x)}$$

**Special case:** If  $p_{OOD}^{proxy}$  is uniform  $\Rightarrow$  criterion reduces to in-distribution likelihood<sup>2</sup>

<sup>&</sup>lt;sup>1</sup>Andi Zhang and Damon Wischik. "Falsehoods that ML researchers believe about OOD detection". In: arXiv preprint arXiv:2210.12767 (2022).

<sup>&</sup>lt;sup>2</sup>Christopher M Bishop. "Novelty detection and neural network validation". In: *IEE Proceedings-Vision, Image and Signal processing* 141.4 (1994), pp. 217–222.

## Failure of the Uniform OOD Proxy



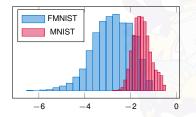
Nalisnick et al.<sup>3</sup> find: using uniform OOD proxy ( $p_{in}(x)$  as the criterion) can fail.

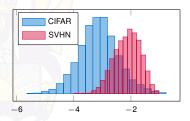
<sup>&</sup>lt;sup>3</sup>Eric Nalisnick et al. "Do deep generative models know what they don't know?" In: *International Conference on Learning Representations* (2019).

## Failure of the Uniform OOD Proxy



- Nalisnick et al.<sup>3</sup> find: using uniform OOD proxy ( $p_{in}(x)$  as the criterion) can fail.
- ► In some cases, OOD samples get higher likelihood than in-distribution samples.





Trained on Fashion-MNIST

Trained on CIFAR-10

Minus BPD  $\propto$  likelihood. OOD likelihood > in-distribution likelihood  $\Rightarrow$  uniform OOD proxy is **misleading**.

<sup>&</sup>lt;sup>3</sup>Eric Nalisnick et al. "Do deep generative models know what they don't know?" In: International Conference on Learning Representations (2019).



▶ We need some smarter OOD proxy



- We need some smarter OOD proxy
- In practice, OOD ≠ random characters from uniform distribution which can be filted by low level methods

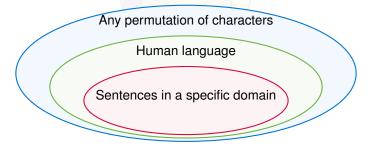


- We need some smarter OOD proxy
- In practice, OOD ≠ random characters from uniform distribution which can be filted by low level methods
- OOD of interest: human language outside the domain



- We need some smarter OOD proxy
- In practice, OOD ≠ random characters from uniform distribution which can be filted by low level methods
- ► OOD of interest: human language outside the domain
- Modern LLMs are trained on massive human-language corpora ⇒ can serve as an OOD proxy:

$$\frac{p_{\text{OOD}}^{\text{proxy}}(x)}{p_{\text{in}}(x)} = \frac{p_{\theta}(x)}{p_{\theta'}(x)} = \frac{\text{pretrained}}{\text{finetuned}}$$





► A finetuned LLM is already a powerful OOD detector



- A finetuned LLM is already a powerful OOD detector
- Just calculate the likelihood ratio:





- A finetuned LLM is already a powerful OOD detector
- Just calculate the likelihood ratio:

$$\frac{p_{\theta}(x)}{p_{\theta'}(x)}$$

Can be implemented by 3 lines of code



- A finetuned LLM is already a powerful OOD detector
- Just calculate the likelihood ratio:

$$\frac{p_{\theta}(x)}{p_{\theta'}(x)}$$

- Can be implemented by 3 lines of code
- Gives you an OOD criterion for free if you already have a finetuned LLM.

Q & A



