

---

# DO STOP ME NOW: DETECTING BOILERPLATE RESPONSES WITH A SINGLE ITERATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Large Language Models (LLMs) often expend significant computational resources generating boilerplate responses, such as refusals, simple acknowledgements and casual greetings, which adds unnecessary cost and latency. To address this inefficiency, we propose a simple yet highly effective method for detecting such responses after only a single generation step. We demonstrate that the log-probability distribution of the first generated token serves as a powerful signal for classifying the nature of the entire subsequent response. Our experiments, conducted across a diverse range of small, large, and reasoning-specialized models, show that the first-token log-probability vectors form distinctly separable clusters for different response types. Using a lightweight k-NN classifier, we achieve high accuracy in predicting whether a response will be a substantive answer or a form of boilerplate response, including user-specified refusals. The primary implication is a practical, computationally trivial technique, optimizing LLM inference by enabling early termination or redirection to a smaller model, thereby yielding significant savings in computational cost. This work presents a direct path toward more efficient and sustainable LLM deployment.

## 1 INTRODUCTION

Large Language Models (LLMs) have revolutionized artificial intelligence with their ability to understand and generate human-like text across diverse applications, from conversational agents to code assistants. This remarkable capability comes with a significant challenge: the high computational and financial costs associated with llm inference for each user query (Regmi & Pun, 2024). These costs are especially wasteful in real-world scenarios where LLMs generates unwanted or predictable outputs. For instance, OpenAI CEO Sam Altman has publicly stated that politeness expressions like "please" and "thank you" have cost the company tens of millions of dollars, due to the electricity consumption associated with generating these boilerplate responses (USA Today, 2025). This highlights a critical inefficiency: LLMs often produce unnecessary tokens that consume resources without contributing to the user's core intent. The ability to accurately and cost-effectively characterize LLM responses *prior to or early in their generation* is thus paramount for optimizing inference costs, reducing latency, and enhancing the overall sustainability of LLM-powered systems.

To address this inefficiency, recent research purposes several novel methods, specifically concerning refusal-to-comply (i.e. due to safeguards). For example, strategies like refusal tokens (Jain et al., 2025) involve prepending special tokens during training that the model learns to generate first when a refusal is appropriate. Other studies show that LLMs encode global attributes of their future responses in hidden representations even before any tokens are generated, enabling emergent response planning (Dong et al., 2025). Another study (Arditi et al., 2024), reveals that refusal behavior is classifiable by a one-dimensional subspace in LLM activations, leading to the development of a "refusal metric". This metric, derived from summing probabilities assigned to specific "refusal tokens" at early generation stages, serves as an efficient proxy for measuring refusal likelihood without full response generation, aligning closely with our objectives for cost-effective content filtering.

This work builds upon these advancements and shows that the **log-probabilities of the first generated token** are sufficient for accurate **prediction of multiple response types**, including user-specified refusals. By focusing on early prediction and detection of boilerplate content (such as refusals, gratitude acknowledgements and other non-task-solving elements), we aim to significantly

---

054 reduce the computational and environmental footprint of LLM inference, ultimately leading to more  
055 efficient and sustainable AI systems.

## 057 2 PRELIMINARIES

### 060 2.1 DEFINING BOILERPLATE RESPONSES

061 We define boilerplate responses as responses that are interchangeable within their class. For exam-  
062 ple, **refusal type response** (i.e. *"I'm sorry, I cannot help you with ----"* or *"You're welcome! I'm*  
063 *glad I could assist you with ----"*) are interchangeable up to context. Efficient and cost-effective  
064 identification of boilerplate enables optimizing inference costs in real-world applications. This in-  
065 volves distinguishing content such as refusals, gratitude, or simple acknowledgements from genu-  
066 inely meaningful, task-solving outputs. Recent advancements in LLM research offer several av-  
067 enues to address this challenge, often by leveraging the internal mechanisms of these models for  
068 early prediction or classification.

## 070 3 RELATED WORK

### 073 3.1 DIALOGUE ACT CLASSIFICATION FOR CONVERSATIONAL INTENT

074 In conversational settings, Dialogue Act Classification (DAC) identifies the communicative intent  
075 behind utterances (Zhangwenbo & Yuhan, 2025; Aljanaideh, 2025). LLMs exhibit zero-shot DAC  
076 capabilities, which can be refined iteratively with online feedback without requiring labeled data  
077 (Zhangwenbo & Yuhan, 2025). This allows for the early classification of dialogue acts, including  
078 those that signify conversational fillers, gratitude, or simple acknowledgements, thereby identifying  
079 non-meaningful conversational turns.

### 081 3.2 DISTINGUISHING MEANINGFUL CONTENT FROM BOILERPLATE ELEMENTS

083 A key challenge is the precise differentiation between essential, reasoning tokens and repetitive, non-  
084 critical boilerplate tokens (e.g., formatting, transitional phrases like *"Based on the user's request..."*)  
085 (Ye et al., 2025). The Shuffle-Aware Discriminator (SHAD) offers an automated and adaptive so-  
086 lution by exploiting predictability differences after shuffling input-output combinations: boilerplate  
087 tokens remain predictable, while reasoning tokens do not.

### 089 3.3 OPERATIONAL EFFICIENCY METHODS FOR COST REDUCTION

091 Beyond direct classification, other techniques aim to reduce LLM inference costs by avoiding un-  
092 necessary computation. Semantic caching mechanisms, such as GPT Semantic Cache, leverage  
093 query embeddings to identify semantically similar questions and retrieve pre-generated responses,  
094 significantly reducing API calls and improving response times, especially for repeated boilerplate  
095 queries (Regmi & Pun, 2024). Additionally, advancements in multi-token prediction enable LLMs  
096 to jointly predict several subsequent tokens in a single inference step (Orgad et al., 2025). While  
097 this primarily speeds up generation, it could potentially be adapted to quickly scan for and flag boil-  
098 erplate patterns, allowing for early termination. Another approach is to use model routing, sending  
099 simpler prompts (such as those leading to boilerplate responses) to a smaller, cheaper model (Ding  
100 et al., 2024).

### 101 3.4 EARLY PREDICTION OF OUTPUT CHARACTERISTICS

103 A foundational concept in this area is emergent response planning, where LLMs' internal hidden  
104 representations can encode global attributes of an entire future response before any tokens are gen-  
105 erated (Dong et al., 2025). By probing these pre-generation representations, it is possible to predict  
106 various characteristics of the upcoming output, such as response length, reasoning steps, answer  
107 confidence, or factual consistency. This capability is instrumental for pre-generation resource allo-  
cation optimization, allowing systems to anticipate the nature of a response and potentially avoid

---

108 costly full generations of boilerplate content. Complementary to this, methods like TRAIL (Sha-  
109 hout et al., 2025) leverage recycled LLM layer embeddings to dynamically predict remaining output  
110 length with low overhead and high accuracy, refining predictions at each token generation step. This  
111 approach, similar to others employing separate lightweight LLMs or BERT models (Devlin et al.,  
112 2019) for length prediction, aims to optimize scheduling and reduce latency, indirectly supporting  
113 early termination for short, non-meaningful responses.

### 114 115 3.5 REFUSAL DETECTION AT INFERENCE TIME 116

117 For specific boilerplate types like refusals, specialized mechanisms have been developed. The Re-  
118 fusal Tokens strategy proposes prepending a special [ `refuse` ] token (or category-specific tokens)  
119 to responses during training, allowing the model to learn to generate this token first when a refusal is  
120 appropriate (Jain et al., 2025). At test-time, the softmax probability of this refusal token quantifies  
121 the likelihood of a refusal, enabling calibrated control over refusal rates without retraining. This en-  
122 ables a "cheap sweep" by allowing identification of optimal refusal thresholds with a single forward  
123 pass, avoiding full response generation. Notably, this method allows for fine-grained control over  
124 various refusal types if multiple tokens are used.

125 Similarly, (Arditi et al., 2024) demonstrates that refusal behavior in LLMs is mediated by a one-  
126 dimensional subspace within their residual stream activations. This allows for the derivation of a  
127 "refusal metric" by summing the probabilities assigned to a predefined set of "refusal tokens" (e.g.,  
128 "I'm sorry", "I cannot") at the first token position of the prompt. This metric serves as an efficient  
129 proxy for estimating the likelihood of a model refusing an instruction without requiring full response  
130 generation. This approach is the most closely related to our work, as it focuses on leveraging early  
131 signals from the model's generation output to predict the nature of the upcoming response. Yet,  
132 unlike our work, it requires manual listing of these "refusal tokens".

133 Beyond explicit refusals, research into LLM transparency also explores their internal signals for  
134 other content characteristics. Studies on LLM hallucinations indicate that truthfulness information  
135 is concentrated in "exact answer tokens" within the generated response (Orgad et al., 2025). Probing  
136 classifiers trained on intermediate representations of these tokens can predict errors, suggesting that  
137 LLMs encode information about their own truthfulness. While focused on error detection, this  
138 highlights the general potential to probe internal states for the "meaningfulness" of a response.

## 139 140 4 OUR METHOD 141

142 When using LLMs to generate responses, they do so one token at a time. At each iteration, *all*  
143 *possible tokens* are assigned probabilities, but only one is selected. Thus, examining all token prob-  
144 abilities of a *single iteration* provides an overview of all the possible subsequent responses. We  
145 hypothesize that by using the (log-)probabilities of the first token generation, it is possible to clas-  
146 sify certain response types. To validate it, we measure the similarities between log-probabilities of  
147 the first token generated by different prompts, designed to induce boilerplate responses versus oth-  
148 ers that suppose to elicit more detailed responses related to the chat. We perform these validation  
149 over different models, both small and large language models, as the token-probability-space is not  
150 comparable between different models.

### 151 152 4.1 DATASET

153 We create a unique dataset of ~3k different chats of different lengths and different classes. We define  
154 several types of classes:

- 156 • **Refusal:** Chats or messages an assistant will refuse to answer due to internal learnt safe-  
157 guards.
- 158 • **Thanks:** Chats ending with the user thanking the assistant for its assistance. These usually  
159 make the assistant respond with common phrases like "You're welcome!" or "My plea-  
160 sure!"
- 161 • **Hello:** Chats starting with the user saying "Hello!", "Hi!" or similar texts.

- 
- **Chat:** All other chats that do not belong to above classes, where the user and the assistant are having a regular conversation.

The dataset was created in the following way:

1. We use the AdvBench dataset (Zou et al., 2023), containing harmful prompts, and classify these as Refusal.
2. We then sampled ~500 random prompts from the Alpaca dataset (Taori et al., 2023), and classify these as Chat. The input-prompts of the Alpaca dataset are split into two columns: *instruction* and *input*, where in most cases, the *input* column is empty. We explicitly used only the *instruction* column as the prompt, thus creating some cases of chats with missing context.
3. Per each of the prompts we now have, we prompted an LLM to respond, and recorded the response. We then asked an LLM to continue the conversation as the user would, and if the model refused to reply - steer the conversation to a legitimate follow up question. We combined the User-Assistant-User interactions as new examples, and labeled them as Chat.
4. For the original harmful prompts, we asked an LLM to hypothesize what was the legitimate prompt the user asked *prior* to the harmful prompt. We then asked an LLM to respond to the legitimate prompt, and added these User(safe)-Assistant-User(harmful) interactions as Refusal.
5. Next, we asked an LLM to come with 250 "benign" thank-you prompts a user might send to an assistant (i.e. "Thank you for your help!"). We then sampled 500 User-Assistant-User chats, replaced the last message with a random thank-you prompt, and labeled them as Thanks.
6. Finally, we created a list of ~30 "benign" prompts which can be used to initiate a conversation with an assistant without any additional request (such as "Hello!" or "Good morning!"). We labeled these as Hello.

The result is a dataset of ~3k chats, containing both single prompt chats and User-Assistant-User chats. The dataset is available at `ANON_URL`<sup>1</sup>.

## 4.2 RESPONSE TYPE CLUSTERING AND CLASSIFICATION

We prompt selected language models with the chats from our dataset, and then record the log-probabilities vectors of the first token generated as the LLM's response. We visualize the results using 2D t-SNE (van der Maaten & Hinton, 2008). To quantitatively evaluate the effectiveness of our approach, we trained k-Nearest Neighbors (k-NN) classifiers (Cover & Hart, 1967) on the first-token log-probability embeddings for each model category. The k-NN algorithm was chosen for its simplicity and interpretability, enabling direct measurement of the clusters separability. For each model, we performed 5-fold stratified cross-validation to ensure robust evaluation across all response types (Chat, Hello, Refusal, Thanks). We fixed  $k = 3$  for all models to enable direct comparison across different architectures. All reported metrics (accuracy, precision, recall, F1-score) are cross-validation results averaged across the 5 folds, providing more reliable estimates than single train-test splits. We report macro-averaged precision, recall, and F1-score to account for class imbalance, particularly for the Hello class which represents only ~1.4% of the dataset.

## 5 EXPERIMENTS

We verify our hypothesis on three scenarios: **Small Language Models**, **Reasoning Models** and **Large Language Models**.

### 5.1 SMALL LANGUAGE MODELS

We perform an evaluation of the first-token log-probabilities of the following models:

---

<sup>1</sup>Dataset is attached to anonymous submission, dataset link is hidden for anonymity

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269

- Llama 3.2 3B (Meta AI, 2024)
- Qwen 2.5 1.5B (Qwen Team, 2024; Yang et al., 2024)
- Gemma 3 1B (Gemma Team et al., 2025)

A 2D t-SNE plot of these models' log-probabilities is shown in Figure 1 and Table 1. We clearly see a separation between the classes.

### 5.1.1 REFUSAL DUE TO INCAPABILITY

When examining Chat samples located much closer to the Refusal class, we find mostly the chats with the missing context we created by omitting the *input* column from the Alpaca dataset. The following list provides a few examples of such incomplete chat messages found more closely to the Refusal class:

- *"Based on the provided input paragraph, provide a summary of its content."*
- *"Translate the sentence below into Japanese."*
- *"Describe the character of the protagonist in the given TV show."*
- *"Group the given list into 3 Groups."*

The assistants are now incapable of replying to these prompts, as they now miss critical context. They therefore trigger a refusal response from the assistants, as they try to explain to the user that they cannot assist - not due to safeguards, but due to lack of context.

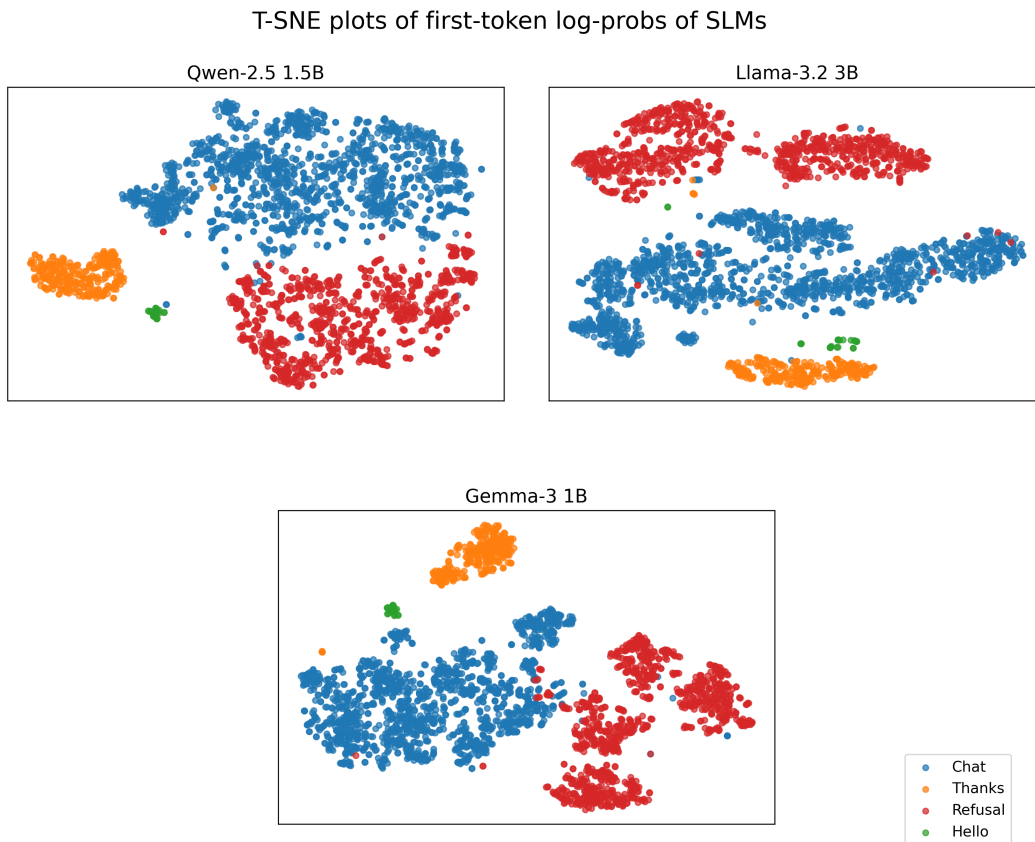


Figure 1: 2D T-SNE plot of first-token log-probabilities for Small Language Models (Llama 3.2 3B, Qwen 2.5 1.5B, Gemma 3 1B). Each point represents a chat, colored by class.

270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
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

Table 1: Small Language Models Performance on Type Classification (k=3)

Model	Accuracy	Precision	Recall	F1	Chat F1	Hello F1	Refusal F1	Thanks F1
Qwen2.5-1.5B	0.997	0.991	0.998	0.994	0.998	1.000	0.998	0.998
Llama-3.2-3B	0.995	0.996	0.984	0.990	0.998	1.000	0.998	0.996
Gemma-3-1B-IT	0.994	0.997	0.997	0.997	0.998	1.000	0.997	1.000

### 5.1.2 REFUSAL DUE TO SYSTEM PROMPT

So far, chats marked as Refusal are those AI assistants tend to refuse to reply to due to internal trained safeguards. We wish to check whether this method applies too when the refusal is due to an arbitrary system prompt, stating that the assistant cannot reply to certain scenarios.

In the experiment shown in Figure 2, we send the assistant a request for a recipe for a Black Forest Cake. We did so twice per model, once with a system prompt explicitly stating that the assistant *cannot* provide a recipe for a Black Forest Cake, and once without any additional restrictions. Figure 2 shows the same 2D T-SNE plot of the first-token log-probabilities of SLM, along with the log-probabilities of the first token of each of the Black Forest Cake prompt variations. We clearly see that when the model was instructed not to provide a recipe, the first token log-probabilities are closer to the Refusal class center-of-mass. While not shown here, other such experiments yielded similar results.

T-SNE plots of first-token log-probs of SLMs, including the Cake Recipe experiment

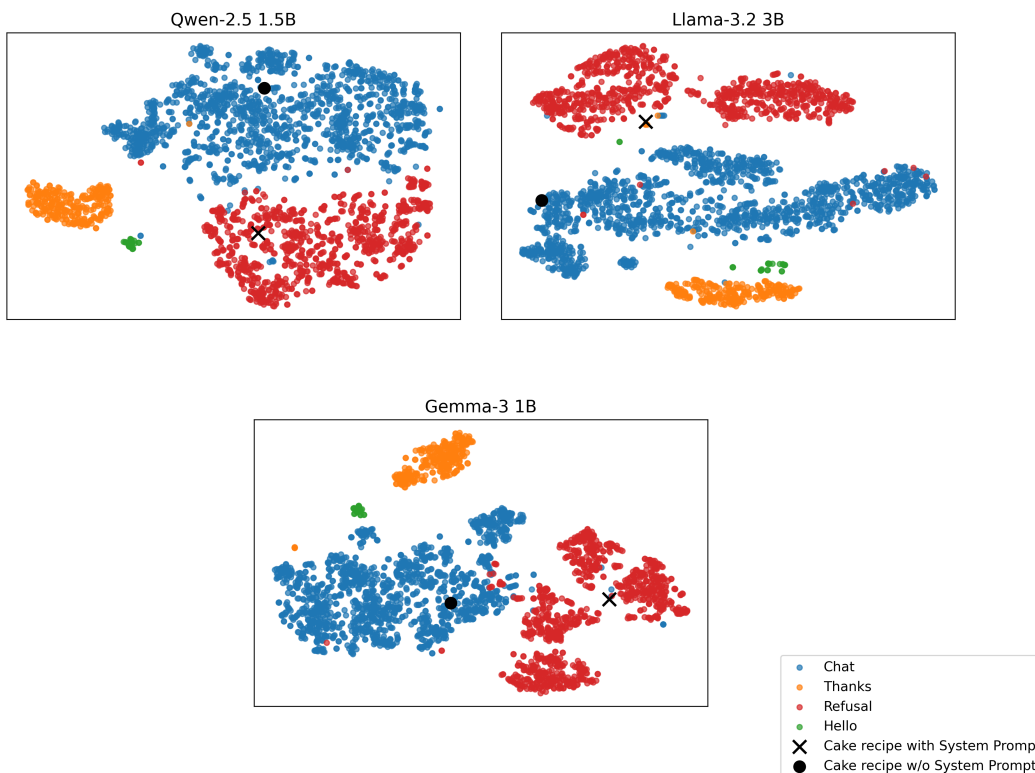


Figure 2: 2D T-SNE plot of first-token log-probabilities for Small Language Models, including the Cake Recipe experiment. The black circle represents a request for a recipe for a Black Forest Cake, and the black cross represents the same requests, but with a system prompt instructing the assistant not to provide the recipe.

## 5.2 REASONING MODELS

The examination of reasoning models adds a bit more complexity, as they tend to begin with a thinking phase which includes self-explaining and "self-talking" (i.e. "So, the user asked...", "I need to...", etc.). Since our classification applies to the first token of the **response**, we slightly adjusted the chat inputs by adding an empty-thinking phase to the assistant's message (i.e. "<think></think>"), and only then checked the log-probabilities of the first token.

We evaluated the following reasoning models:

- DeepSeek-R1 8B (Llama Distillation) (DeepSeek-AI, 2025)
- Phi-4 Reasoning Plus (Abdin et al., 2025)

Results are shown in Figure 3 and Table 2. Here too, we see a clear separation between the classes.

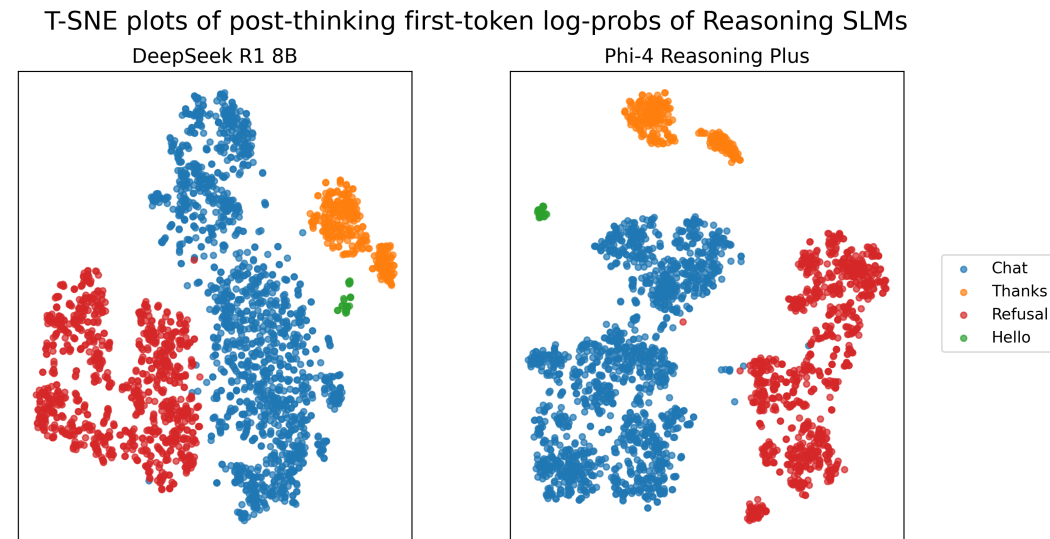


Figure 3: 2D T-SNE plot of post-empty thinking ("<think></think>") first-token partial log-probabilities for Reasoning Models (DeepSeek-R1 8B, Phi-4 Reasoning Plus). Each point represents a chat, colored by class.

Table 2: Reasoning Models Performance on Type Classification (k=3)

Model	Accuracy	Precision	Recall	F1	Chat F1	Hello F1	Refusal F1	Thanks F1
Phi-4-Reasoning+	0.998	0.999	0.999	0.999	0.999	1.000	0.999	1.000
DeepSeek-R1-8B	0.998	0.998	0.989	0.993	0.999	1.000	0.999	1.000

## 5.3 LARGE LANGUAGE MODELS

We perform an evaluation of the first-token log-probabilities of the following cloud-based models:

- OpenAI GPT-4o<sup>2</sup> (OpenAI et al., 2024)
- Gemini 2.0 Flash (Hassabis & Kavukcuoglu, 2024)

Unlike open-source models, these models do not provide the full log-probabilities of generated tokens, but only the top 20. We reconstructed the partial log-probabilities vectors, as the APIs of

<sup>2</sup>Running on Microsoft Azure

these models provide both the log-probability and token IDs (token indices) of the top 20 tokens.<sup>3</sup> Therefore, the vectors displayed in Figure 4 and Table 3 are *trimmed* log-probabilities vectors. Even with the trimmed log-probabilities, we still see clear separation between the classes.

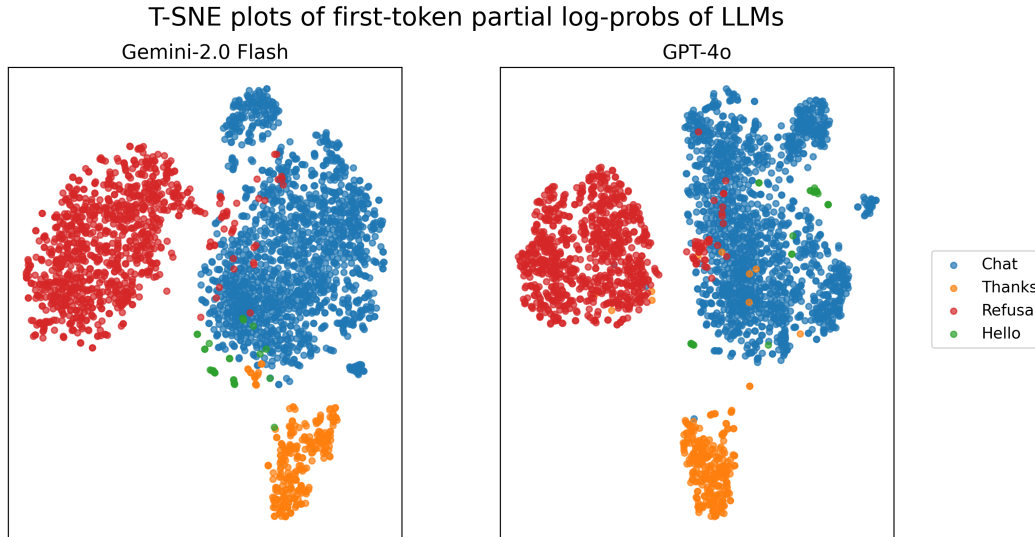


Figure 4: 2D T-SNE plot of first-token partial log-probabilities for Large Language Models (GPT-4o, Gemini-2.0 Flash). Each point represents a chat, colored by class.

Table 3: Large Language Models Performance on Type Classification (k=3)

Model	Accuracy	Precision	Recall	F1	Chat F1	Hello F1	Refusal F1	Thanks F1
Gemini-2.0-Flash	0.979	0.989	0.844	0.884	0.993	0.826	0.995	0.993
GPT-4o	0.974	0.983	0.914	0.941	0.992	0.941	0.987	0.982

## 6 CONCLUSION

In this work, we addressed the significant computational inefficiency of Large Language Models generating boilerplate responses. We introduced a simple yet highly effective method to predict the nature of an entire response by analyzing the log-probability distribution of just the first generated token.

Our comprehensive experiments across a diverse range of small, large, and reasoning-specialized models consistently demonstrated that the log-probabilities of the first token form distinct, separable clusters for different response types, such as substantive answers, refusals, simple acknowledgements or greetings. We showed that a lightweight k-NN classifier can leverage these clusters to achieve high accuracy in predicting the response category after a single generation step. Furthermore, our method successfully identifying refusals prompted by both inherent model safeguards and arbitrary, user-defined system prompts.

The primary implication of our findings is a practical and computationally trivial technique to optimize LLM inference. By enabling early termination of unwanted boilerplate generation, this approach offers substantial savings in computational cost and latency. This work presents a direct path toward more efficient, economical, and sustainable deployment of LLM systems, paving the way for more responsive and cost-effective applications.

<sup>3</sup>For GPT-4o, we used `tiktoken` to retrieve the token IDs from the token string-values: [github.com/openai/tiktoken](https://github.com/openai/tiktoken)

---

432 Future work could focus on applying this technique to a wider range of boilerplate categories, multi-  
433 language scenarios, and exploring its effectiveness in multi-modal contexts.  
434

## 435 REFERENCES 436

437 Marah Abdin, Sahaj Agarwal, Ahmed Awadallah, Vidhisha Balachandran, Harkirat Behl, Lingjiao  
438 Chen, Gustavo de Rosa, Suriya Gunasekar, Mojan Javaheripi, Neel Joshi, Piero Kauffmann, Yash  
439 Lara, Caio César Teodoro Mendes, Arindam Mitra, Besmira Nushi, Dimitris Papailiopoulos,  
440 Olli Saarikivi, Shital Shah, Vaishnavi Shrivastava, Vibhav Vineet, Yue Wu, Safoora Yousefi, and  
441 Guoqing Zheng. Phi-4-reasoning technical report, 2025. URL [https://arxiv.org/abs/  
442 2504.21318](https://arxiv.org/abs/2504.21318).

443 Ahmad Aljanaideh. Speech act patterns for improving generalizability of explainable politeness  
444 detection models. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher  
445 Pilehvar (eds.), *Findings of the Association for Computational Linguistics: ACL 2025*, pp. 18945–  
446 18954, Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-  
447 89176-256-5. doi: 10.18653/v1/2025.findings-acl.970. URL [https://aclanthology.  
448 org/2025.findings-acl.970/](https://aclanthology.org/2025.findings-acl.970/).

449 Andy Arditi, Oscar Balcells Obeso, Aaqib Syed, Daniel Paleka, Nina Rimsky, Wes Gurnee, and  
450 Neel Nanda. Refusal in language models is mediated by a single direction. In *The Thirty-  
451 eighth Annual Conference on Neural Information Processing Systems*, 2024. URL [https:  
452 //openreview.net/forum?id=pH3XAQME6c](https://openreview.net/forum?id=pH3XAQME6c).

453 T. Cover and P. Hart. Nearest neighbor pattern classification. *IEEE Transactions on Information  
454 Theory*, 13(1):21–27, 1967. doi: 10.1109/TIT.1967.1053964.  
455

456 DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning,  
457 2025. URL <https://arxiv.org/abs/2501.12948>.

458 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of  
459 deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and  
460 Tamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of  
461 the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long  
462 and Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Com-  
463 putational Linguistics. doi: 10.18653/v1/N19-1423. URL [https://aclanthology.org/  
464 N19-1423/](https://aclanthology.org/N19-1423/).

465 Dujian Ding, Ankur Mallick, Chi Wang, Robert Sim, Subhabrata Mukherjee, Victor Rühle, Laks  
466 V. S. Lakshmanan, and Ahmed Hassan Awadallah. Hybrid LLM: Cost-efficient and quality-aware  
467 query routing. In *The Twelfth International Conference on Learning Representations*, 2024. URL  
468 <https://openreview.net/forum?id=02f3mUtqnM>.

469 Zhichen Dong, Zhanhui Zhou, Zhixuan Liu, Chao Yang, and Chaochao Lu. Emergent response  
470 planning in llms, 2025. URL <https://arxiv.org/abs/2502.06258>.  
471

472 Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej,  
473 Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard, Thomas  
474 Mesnard, Geoffrey Cideron, Jean bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Cas-  
475 bon, Etienne Pot, Ivo Penchev, Gaël Liu, Francesco Visin, Kathleen Kenealy, Lucas Beyer, Xi-  
476 aohai Zhai, Anton Tsitsulin, Robert Busa-Fekete, Alex Feng, Noveen Sachdeva, Benjamin Cole-  
477 man, Yi Gao, Basil Mustafa, Iain Barr, Emilio Parisotto, David Tian, Matan Eyal, Colin Cherry,  
478 Jan-Thorsten Peter, Danila Sinopalnikov, Surya Bhupatiraju, Rishabh Agarwal, Mehran Kazemi,  
479 Dan Malkin, Ravin Kumar, David Vilar, Idan Brusilovsky, Jiaming Luo, Andreas Steiner, Abe  
480 Friesen, Abhanshu Sharma, Abheesht Sharma, Adi Mayrav Gilady, Adrian Goedeckemeyer, Alaa  
481 Saade, Alex Feng, Alexander Kolesnikov, Alexei Bendebury, Alvin Abdagic, Amit Vadi, Andrés  
482 György, André Susano Pinto, Anil Das, Ankur Bapna, Antoine Miech, Antoine Yang, Antonia  
483 Paterson, Ashish Shenoy, Ayan Chakrabarti, Bilal Piot, Bo Wu, Bobak Shahriari, Bryce Pettrini,  
484 Charlie Chen, Charline Le Lan, Christopher A. Choquette-Choo, CJ Carey, Cormac Brick, Daniel  
485 Deutsch, Danielle Eisenbud, Dee Cattle, Derek Cheng, Dimitris Pappas, Divyashree Shivaku-  
mar Sreepathihalli, Doug Reid, Dustin Tran, Dustin Zelle, Eric Noland, Erwin Huizenga, Eu-  
gene Kharitonov, Frederick Liu, Gagik Amirkhanyan, Glenn Cameron, Hadi Hashemi, Hanna

---

486 Klimczak-Plucińska, Harman Singh, Harsh Mehta, Harshal Tushar Lehri, Hussein Hazimeh, Ian  
487 Ballantyne, Idan Szpektor, Ivan Nardini, Jean Pouget-Abadie, Jetha Chan, Joe Stanton, John Wi-  
488 eting, Jonathan Lai, Jordi Orbay, Joseph Fernandez, Josh Newlan, Ju yeong Ji, Jyotinder Singh,  
489 Kat Black, Kathy Yu, Kevin Hui, Kiran Vodrahalli, Klaus Greff, Linhai Qiu, Marcella Valentine,  
490 Marina Coelho, Marvin Ritter, Matt Hoffman, Matthew Watson, Mayank Chaturvedi, Michael  
491 Moynihan, Min Ma, Nabila Babar, Natasha Noy, Nathan Byrd, Nick Roy, Nikola Momchev, Ni-  
492 lay Chauhan, Noveen Sachdeva, Oskar Bunyan, Pankil Botarda, Paul Caron, Paul Kishan Ruben-  
493 stein, Phil Culliton, Philipp Schmid, Pier Giuseppe Sessa, Pingmei Xu, Piotr Stanczyk, Pouya  
494 Tafti, Rakesh Shivanna, Renjie Wu, Renke Pan, Reza Rokni, Rob Willoughby, Rohith Vallu,  
495 Ryan Mullins, Sammy Jerome, Sara Smoot, Sertan Girgin, Shariq Iqbal, Shashir Reddy, Shruti  
496 Sheth, Siim Pöder, Sijal Bhatnagar, Sindhu Raghuram Panyam, Sivan Eiger, Susan Zhang, Tianqi  
497 Liu, Trevor Yacovone, Tyler Liechty, Uday Kalra, Utku Evci, Vedant Misra, Vincent Roseberry,  
498 Vlad Feinberg, Vlad Kolesnikov, Woohyun Han, Woosuk Kwon, Xi Chen, Yinlam Chow, Yuvein  
499 Zhu, Zichuan Wei, Zoltan Egyed, Victor Cotruta, Minh Giang, Phoebe Kirk, Anand Rao, Kat  
500 Black, Nabila Babar, Jessica Lo, Erica Moreira, Luiz Gustavo Martins, Omar Sanseviero, Lucas  
501 Gonzalez, Zach Gleicher, Tris Warkentin, Vahab Mirrokni, Evan Senter, Eli Collins, Joelle Bar-  
502 ral, Zoubin Ghahramani, Raia Hadsell, Yossi Matias, D. Sculley, Slav Petrov, Noah Fiedel, Noam  
503 Shazeer, Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Farabet, Elena  
504 Buchatskaya, Jean-Baptiste Alayrac, Rohan Anil, Dmitry, Lepikhin, Sebastian Borgeaud, Olivier  
505 Bachem, Armand Joulin, Alek Andreev, Cassidy Hardin, Robert Dadashi, and Léonard Hussenot.  
506 Gemma 3 technical report, 2025. URL <https://arxiv.org/abs/2503.19786>.

506 Demis Hassabis and Koray Kavukcuoglu. Introducing gemini 2.0: Our new ai model  
507 for the agentic era. [https://blog.google/technology/google-deepmind/  
508 google-gemini-ai-update-december-2024/](https://blog.google/technology/google-deepmind/google-gemini-ai-update-december-2024/), December 2024. Published December  
509 11, 2024 on The Keyword (Google Blog).

510 Neel Jain, Aditya Shrivastava, Chenyang Zhu, Daben Liu, Alfy Samuel, Ashwinee Panda, Anoop  
511 Kumar, Micah Goldblum, and Tom Goldstein. Refusal tokens: A simple way to calibrate re-  
512 fusals in large language models, 2025. URL [https://openreview.net/forum?id=  
513 QnylufReka](https://openreview.net/forum?id=QnylufReka).

514 Lingrui Mei, Shenghua Liu, Yiwei Wang, Baolong Bi, Ruibin Yuan, and Xueqi Cheng. Hid-  
515 danguard: Fine-grained safe generation with specialized representation router, 2025. URL  
516 <https://openreview.net/forum?id=NgCNMLTXx9>.

517

518 Meta AI. Llama 3.2 connect 2024: Vision, edge & mobile devices. [https://ai.meta.com/  
519 blog/llama-3-2-connect-2024-vision-edge-mobile-devices/](https://ai.meta.com/blog/llama-3-2-connect-2024-vision-edge-mobile-devices/), September  
520 2024. Meta AI blog post.

521

522 Sander Noels, Guillaume Bied, Maarten Buyl, Alexander Rogiers, Yousra Fettach, Jeffrey Lijffijt,  
523 and Tijn De Bie. What large language models do not talk about: An empirical study of moderation  
524 and censorship practices, 2025. URL <https://arxiv.org/abs/2504.03803>.

525

526 OpenAI, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan  
527 Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Madry, Alex Baker-  
528 Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol,  
529 Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Con-  
530 neau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoochian,  
531 Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein,  
532 Andrew Cann, Andrew Codispoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey  
533 Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia,  
534 Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben  
535 Leimberger, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake  
536 Samic, Bob McGrew, Bobby Spero, Bogo Giertler, Bowen Cheng, Brad Lightcap, Brandon  
537 Walkin, Brendan Quinn, Brian Guarraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo  
538 Lugaresi, Carroll Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson, Chak Li,  
539 Chan Jun Shern, Channing Conger, Charlotte Barette, Chelsea Voss, Chen Ding, Cheng Lu,  
Chong Zhang, Chris Beaumont, Chris Hallacy, Chris Koch, Christian Gibson, Christina Kim,  
Christine Choi, Christine McLeavey, Christopher Hesse, Claudia Fischer, Clemens Winter, Coley  
Czarnecki, Colin Jarvis, Colin Wei, Constantin Koumouzelis, Dane Sherburn, Daniel Kappler,

---

540 Daniel Levin, Daniel Levy, David Carr, David Farhi, David Mely, David Robinson, David Sasaki,  
541 Denny Jin, Dev Valladares, Dimitris Tsipras, Doug Li, Duc Phong Nguyen, Duncan Findlay,  
542 Edede Oiwoh, Edmund Wong, Ehsan Asdar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow,  
543 Eric Kramer, Eric Peterson, Eric Sigler, Eric Wallace, Eugene Brevdo, Evan Mays, Farzad Kho-  
544 rasani, Felipe Petroski Such, Filippo Raso, Francis Zhang, Fred von Lohmann, Freddie Sulit,  
545 Gabriel Goh, Gene Oden, Geoff Salmon, Giulio Starace, Greg Brockman, Hadi Salman, Haiming  
546 Bao, Haitang Hu, Hannah Wong, Haoyu Wang, Heather Schmidt, Heather Whitney, Heewoo Jun,  
547 Hendrik Kirchner, Henrique Ponde de Oliveira Pinto, Hongyu Ren, Huiwen Chang, Hyung Won  
548 Chung, Ian Kivlichan, Ian O’Connell, Ian O’Connell, Ian Osband, Ian Silber, Ian Sohl, Ibrahim  
549 Okuyucu, Ikai Lan, Ilya Kostrikov, Ilya Sutskever, Ingmar Kanitscheider, Ishaan Gulrajani, Ja-  
550 cob Coxon, Jacob Menick, Jakub Pachocki, James Aung, James Betker, James Crooks, James  
551 Lennon, Jamie Kiros, Jan Leike, Jane Park, Jason Kwon, Jason Phang, Jason Teplitz, Jason Wei,  
552 Jason Wolfe, Jay Chen, Jeff Harris, Jenia Varavva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui  
553 Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joanne Jang, Joaquin Quinero Candela, Joe Beutler, Joe  
554 Landers, Joel Parish, Johannes Heidecke, John Schulman, Jonathan Lachman, Jonathan McKay,  
555 Jonathan Uesato, Jonathan Ward, Jong Wook Kim, Joost Huizinga, Jordan Sitkin, Jos Kraaijeveld,  
556 Josh Gross, Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao, Joyce Lee, Juntang Zhuang,  
557 Justyn Harriman, Kai Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kavin Karthik, Kayla Wood,  
558 Kendra Rimbach, Kenny Hsu, Kenny Nguyen, Keren Gu-Lemberg, Kevin Button, Kevin Liu, Kiel  
559 Howe, Krithika Muthukumar, Kyle Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lauren Work-  
560 man, Leher Pathak, Leo Chen, Li Jing, Lia Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka,  
561 Lilian Weng, Lindsay McCallum, Lindsey Held, Long Ouyang, Louis Feuvrier, Lu Zhang, Lukas  
562 Kondraciuk, Lukasz Kaiser, Luke Hewitt, Luke Metz, Lyric Doshi, Mada Aflak, Maddie Simens,  
563 Madelaine Boyd, Madeleine Thompson, Marat Dukhan, Mark Chen, Mark Gray, Mark Hudnall,  
564 Marvin Zhang, Marwan Aljubeh, Mateusz Litwin, Matthew Zeng, Max Johnson, Maya Shetty,  
565 Mayank Gupta, Meghan Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao Zhong, Mia Glaese,  
566 Mianna Chen, Michael Janner, Michael Lampe, Michael Petrov, Michael Wu, Michele Wang,  
567 Michelle Fradin, Michelle Pokrass, Miguel Castro, Miguel Oom Temudo de Castro, Mikhail  
568 Pavlov, Miles Brundage, Miles Wang, Minal Khan, Mira Murati, Mo Bavarian, Molly Lin, Murat  
569 Yesildal, Nacho Soto, Natalia Gimelshein, Natalie Cone, Natalie Staudacher, Natalie Summers,  
570 Natan LaFontaine, Neil Chowdhury, Nick Ryder, Nick Stathas, Nick Turley, Nik Tezak, Niko Fel-  
571 ix, Nithanth Kudige, Nitish Keskar, Noah Deutsch, Noel Bundick, Nora Puckett, Ofir Nachum,  
572 Ola Okelola, Oleg Boiko, Oleg Murk, Oliver Jaffe, Olivia Watkins, Olivier Godement, Owen  
573 Campbell-Moore, Patrick Chao, Paul McMillan, Pavel Belov, Peng Su, Peter Bak, Peter Bakkum,  
574 Peter Deng, Peter Dolan, Peter Hoeschele, Peter Welinder, Phil Tillet, Philip Pronin, Philippe  
575 Tillet, Prafulla Dhariwal, Qiming Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Rajan Troll, Ran-  
576 dall Lin, Rapha Gontijo Lopes, Raul Puri, Reah Miyara, Reimar Leike, Renaud Gaubert, Reza  
577 Zamani, Ricky Wang, Rob Donnelly, Rob Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchan-  
578 dani, Romain Huet, Rory Carmichael, Rowan Zellers, Roy Chen, Ruby Chen, Ruslan Nigmat-  
579 ullin, Ryan Cheu, Saachi Jain, Sam Altman, Sam Schoenholz, Sam Toizer, Samuel Miserendino,  
580 Sandhini Agarwal, Sara Culver, Scott Ethersmith, Scott Gray, Sean Grove, Sean Metzger, Shamez  
581 Hermani, Shantanu Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shirong Wu, Shuaiqi, Xia,  
582 Sonia Phene, Spencer Papay, Srinivas Narayanan, Steve Coffey, Steve Lee, Stewart Hall, Suchir  
583 Balaji, Tal Broda, Tal Stramer, Tao Xu, Tarun Gogineni, Taya Christianson, Ted Sanders, Tejal  
584 Patwardhan, Thomas Cunningham, Thomas Degry, Thomas Dimson, Thomas Raoux, Thomas  
585 Shadwell, Tianhao Zheng, Todd Underwood, Todor Markov, Toki Sherbakov, Tom Rubin, Tom  
586 Stasi, Tomer Kaftan, Tristan Heywood, Troy Peterson, Tyce Walters, Tyna Eloundou, Valerie Qi,  
587 Veit Moeller, Vinnie Monaco, Vishal Kuo, Vlad Fomenko, Wayne Chang, Weiyi Zheng, Wenda  
588 Zhou, Wesam Manassra, Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian, Yongjik Kim,  
589 Youlong Cheng, Yu Zhang, Yuchen He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and Yury Malkov.  
590 Gpt-4o system card, 2024. URL <https://arxiv.org/abs/2410.21276>.

588 Hadas Orgad, Michael Toker, Zorik Gekhman, Roi Reichart, Idan Szpektor, Hadas Kotek, and  
589 Yonatan Belinkov. LLMs know more than they show: On the intrinsic representation of LLM  
590 hallucinations. In *The Thirteenth International Conference on Learning Representations, 2025*.  
591 URL <https://openreview.net/forum?id=KRnsX5Em3W>.

592

593 Qwen Team. Qwen2.5: A party of foundation models, September 2024. URL <https://qwenlm.github.io/blog/qwen2.5/>.

---

594 Sajal Regmi and Chetan Phakami Pun. Gpt semantic cache: Reducing llm costs and latency via  
595 semantic embedding caching, 2024. URL <https://arxiv.org/abs/2411.05276>.  
596

597 Rana Shahout, eran malach, Chunwei Liu, Weifan Jiang, Minlan Yu, and Michael Mitzenmacher.  
598 DON't STOP ME NOW: EMBEDDING BASED SCHEDULING FOR LLMS. In *The Thirteenth*  
599 *International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=7JhGdZvW4T>.  
600

601 Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy  
602 Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model.  
603 [https://github.com/tatsu-lab/stanford\\_alpaca](https://github.com/tatsu-lab/stanford_alpaca), 2023.  
604

605 USA Today. “please, thank you”: Chatgpt, openai energy  
606 costs. [https://www.usatoday.com/story/tech/2025/04/22/  
607 please-thank-you-chatgpt-openai-energy-costs/83207447007/](https://www.usatoday.com/story/tech/2025/04/22/please-thank-you-chatgpt-openai-energy-costs/83207447007/), April  
608 2025. Accessed: 2025-08-14.

609 Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Ma-*  
610 *chine Learning Research*, 9(86):2579–2605, 2008. URL [http://jmlr.org/papers/v9/  
611 vandermaaten08a.html](http://jmlr.org/papers/v9/vandermaaten08a.html).

612 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,  
613 Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang,  
614 Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai,  
615 Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng  
616 Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai  
617 Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan  
618 Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang  
619 Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. Qwen2  
620 technical report. *arXiv preprint arXiv:2407.10671*, 2024.

621 Ziang Ye, Zhenru Zhang, Yang Zhang, Jianxin Ma, Junyang Lin, and Fuli Feng. Disentangling  
622 reasoning tokens and boilerplate tokens for language model fine-tuning. In Wanxiang Che,  
623 Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Findings of the As-*  
624 *sociation for Computational Linguistics: ACL 2025*, pp. 20939–20957, Vienna, Austria, July  
625 2025. Association for Computational Linguistics. ISBN 979-8-89176-256-5. doi: 10.18653/v1/  
626 2025.findings-acl.1078. URL [https://aclanthology.org/2025.findings-acl.  
627 1078/](https://aclanthology.org/2025.findings-acl.1078/).

628 Zhangwenbo Zhangwenbo and Wang Yuhan. Act2P: LLM-driven online dialogue act classification  
629 for power analysis. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher  
630 Pilehvar (eds.), *Findings of the Association for Computational Linguistics: ACL 2025*, pp. 20494–  
631 20504, Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-  
632 89176-256-5. doi: 10.18653/v1/2025.findings-acl.1052. URL [https://aclanthology.  
633 org/2025.findings-acl.1052/](https://aclanthology.org/2025.findings-acl.1052/).

634 Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial  
635 attacks on aligned language models, 2023.  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647