
Efficient Prompt Caching for Large Language Model Inference via Embedding Similarity

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Large language models (LLMs) have achieved huge success in numerous natu-
2 ral language process (NLP) tasks. However, it faces the challenge of significant
3 resource consumption during inference. In this paper, we aim to improve the
4 inference efficiency of LLMs by prompt caching, i.e., if the current prompt can be
5 answered by the same response of a previous prompt, one can directly utilize that
6 response without calling the LLM. Specifically, we focus on the prediction accuracy
7 of prompt caching for single-round question-answering tasks via embedding simi-
8 larity. The existing embeddings of prompts mostly focus on whether two prompts
9 are semantically similar, which is not necessarily equivalent to whether the same
10 response can answer them. Therefore, we propose a distillation-based method to
11 fine-tune the existing embeddings for better caching prediction. Theoretically, we
12 provide finite-sample guarantees for the convergence of our method under different
13 types of loss functions. Empirically, we construct a dataset based on Kwiatkowski
14 et al. [2019] and fine-tune the embedding from Wang et al. [2022], which improves
15 the AUC of caching prediction from 0.85 to 0.92 within 10 minutes of training.

16 1 Introductions

17 The recent development of large language models (LLMs) and foundation models has notably
18 enhanced the potential of AI systems [Ziegler et al., 2019, Wei et al., 2022, Chowdhery et al., 2022,
19 Ouyang et al., 2022, Bubeck et al., 2023, Nori et al., 2023, OpenAI, 2023, Beeching et al., 2023,
20 Anil et al., 2023] However, due to the large scale of those models, it causes significant resource
21 consumptions not only during the training process, but also in the inference stage [Sharir et al., 2020,
22 Patterson et al., 2021, Bommasani et al., 2022]. Moreover, the latency of LLMs during inference is not
23 negligible since the model only generates one token at a time due to its auto-regressive nature, which
24 makes it unfavorable to be applied to systems desiring high throughput, such as search engines [Zhu
25 et al., 2023]. Therefore, it would be appealing to reduce the resource consumption and latency
26 without degrading the performance of LLMs.

27 A natural idea to reduce resource consumption and latency is to reduce the number of calls to LLMs,
28 which can be implemented by caching, a technique that has a long history of being studied and
29 applied to important areas such as computer architecture and web retrieval [Smith, 1982, Wang,
30 1999, Kumar and Singh, 2016]. Zhu et al. [2023] studies prompt (or query) caching for LLMs, i.e.,
31 some of the previous prompt-response pairs are stored in a cache with limited size, and whenever a
32 new prompt arrives, one can search in the cache whether a prompt has the same semantic meaning as
33 the current prompt, and can directly reuse the response of the previous prompt without calling LLMs
34 if there is a hit (see Figure 1 for a figurative illustration).

35 Zhu et al. [2023] focuses on caching algorithm design and directly assumes a semantic search
36 oracle. Although previous literature studies semantic search or embedding-based methods [Bast et al.,

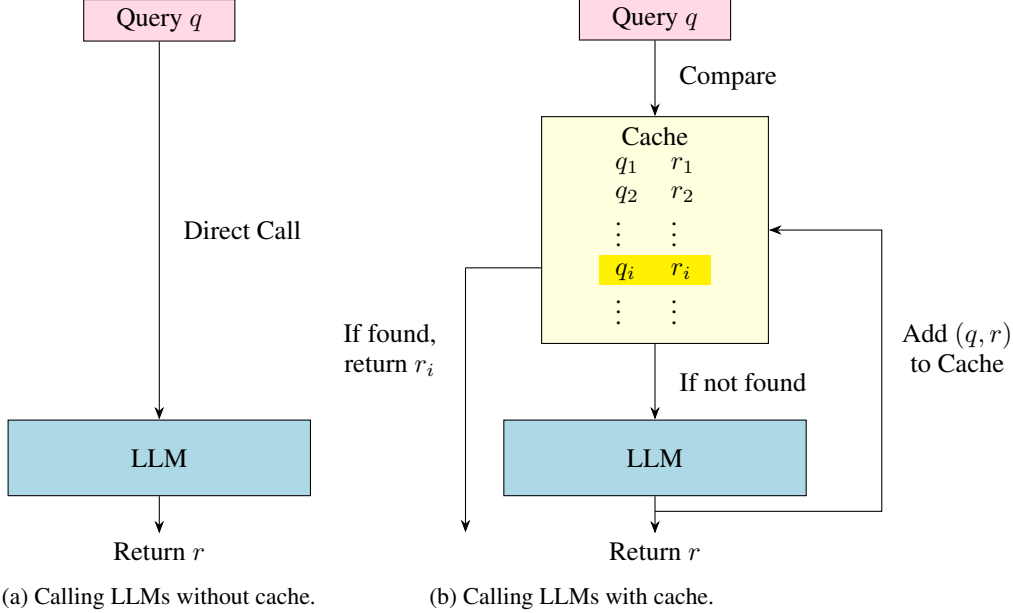


Figure 1: The procedure of calling LLMs with or without cache. When a cache is available, one can store some of the previous prompt-response pairs in the cache, and for a new prompt, one can search in the cache whether a prompt has the same semantic meaning as the current prompt. If there is a hit, one can directly reuse the response of the previous prompt without calling LLMs.

2016, Chang et al., 2020, Kamaloo et al., 2023], which could serve as solutions to the caching hit problem [zilliztech, 2023]¹, it is challenging to obtain a good embedding that can accurately represent the semantic meaning of a prompt. Moreover, a semantically similar prompt pair cannot necessarily be answered by the same response, which implies that we need a different vector embedding specifically for the caching hitting problem that can be used to search a similar prompt more efficiently and better predict the probability that a pair of prompts can be answered by the same response.

In this paper, we aim to learn a good vector embedding such that the similarity of embeddings of a prompt pair could encode the information of whether the pair of prompts can be answered by the same response, i.e., to better predict the probability that they can be answered by the same response. We propose a distillation-based method, which aims to learn the ground-truth probability of whether a prompt pair can be answered by the response via cosine similarity of the embeddings of the prompt pair, to fine-tune an existing semantic vector embedding from Wang et al. [2022]. Theoretically, we provide finite sample guarantees for the learning error under mild assumptions using cross entropy and squared log difference errors, respectively (Section 3). Empirically, we construct a dataset based on Kwiatkowski et al. [2019] and fine-tune the embedding from Wang et al. [2022], which improves the AUC of caching prediction from 0.85 to 0.92 within 10 minutes of training using cross entropy error (Section 4).

2 Preliminaries

We introduce in this section some basic notations, definitions, and assumptions. Let \mathcal{Q} denote the set of all possible prompts (queries). For any prompt pair $(q_1, q_2) \in \mathcal{Q} \times \mathcal{Q}$, we denote the ground-truth probability that (q_1, q_2) can be answered by the same response by $P^*(q_1 = q_2)$. Assume there exists an underlying distribution μ of prompt pairs (q_1, q_2) . Note that we do not have direct access to the μ and instead are given a dataset $\mathcal{D} = \{(q_{i,1}, q_{i,2}, p_i)\}_{i=1}^N$, where $(q_{i,1}, q_{i,2}) \stackrel{\text{i.i.d.}}{\sim} \mu$ and $p_i \in [0, 1]$ with

¹i.e., searching whether there exists a prompt in the cache s.t. the current prompt can be answered by the same response.

60 $p_i \sim \mathcal{P}(\cdot | q_{i,1}, q_{i,2})^2$ and $\mathbb{E}[p_i | q_{i,1}, q_{i,2}] = P^*(q_{i,1} = q_{i,2})$. Below, we define the vector embedding
61 of prompts and define probability via embedding similarity.

62 **Definition 2.1** (Embedding of prompts). For any prompt q , let $v_\theta(q) \in \mathbb{R}^d$ denote its vector
63 embedding where v can be viewed as the mapping of prompts to a specific layer of a language model,
64 and $\theta \in \Theta$ is the parameters of that model.

65 **Definition 2.2** (Probability via embedding similarity). For any two prompts q_1, q_2 , we denote the
66 induced probability via embedding similarity that q_1, q_2 can be answered by the same response by

$$P_{\theta, \lambda, c}(q_1 = q_2) \triangleq \sigma(\text{sim}(v_\theta(q_1), v_\theta(q_2)) / \lambda - c),$$

67 where $\text{sim}(\cdot, \cdot)$ denotes the cosine similarity, i.e., $\text{sim}(x, y) = \frac{\langle x, y \rangle}{\|x\| \cdot \|y\|}$ for two vectors $x, y \in \mathbb{R}^d$,
68 $\sigma(x) = \frac{1}{1 + \exp(-x)}$ for $x \in \mathbb{R}$, and $\lambda \in \Lambda \subset \mathbb{R}_+$, $c \in \mathcal{C} \subset \mathbb{R}$ are two real-valued parameters.

69 We make the following assumptions for theoretical analysis.

70 **Assumption 2.3** (Realizability). Assume there exists $\theta^* \in \Theta$, $\lambda^* \in \Lambda$, $c^* \in \mathcal{C}$, s.t. for any prompt
71 pairs $(q_1, q_2) \in \mathcal{Q} \times \mathcal{Q}$, it holds that

$$P_{\theta^*, \lambda^*, c^*}(q_1 = q_2) = P^*(q_1 = q_2).$$

72 **Assumption 2.4** (Boundedness). Assume there exist constants $L_\lambda, B_c > 0$, s.t.

$$\lambda \geq L_\lambda, |c| \leq B_c, \quad \forall \lambda \in \Lambda, c \in \mathcal{C}.$$

73 For convenience, for any $p \in [0, 1]$, we denote $\bar{p} = 1 - p$. Also, for any prompt pairs q_1, q_2 , we
74 denote $\bar{P}^*(q_1 = q_2) = 1 - P^*(q_1 = q_2)$ and $\bar{P}_{\theta, \lambda, c}(q_1 = q_2) = 1 - P_{\theta, \lambda, c}(q_1 = q_2)$.

75 3 Theoretical Results

76 In this section, we provide finite sample guarantees for the convergence of the learning error. We
77 compare two different loss functions, i.e., binary cross entropy loss (Section 3.1) and squared log
78 difference loss (Section 3.2).

79 3.1 Convergence guarantee for cross-entropy loss

80 For any $\theta \in \Theta$, $\lambda \in \Lambda$, $c \in \mathcal{C}$, we denote the (binary) cross-entropy loss function as

$$\mathcal{L}_\mu^{\text{BCE}}(\theta, \lambda, c) = -\mathbb{E}_{(q_1, q_2) \sim \mu} [P^*(q_1 = q_2) \log P_{\theta, \lambda, c}(q_1 = q_2) + \bar{P}^*(q_1 = q_2) \log \bar{P}_{\theta, \lambda, c}(q_1 = q_2)]. \quad (1)$$

81 To recover the ground-truth parameter $(\theta^*, \lambda^*, c^*)$, one only needs to solve the optimization problem

$$\min_{\theta \in \Theta, \lambda \in \Lambda, c \in \mathcal{C}} \mathcal{L}_\mu^{\text{BCE}}(\theta, \lambda, c). \quad (2)$$

82 One may observe that (2) is equivalent to minimizing the expected KL divergence between the ground
83 truth probability P^* and $P_{\theta, \lambda, c}$.

84 Our algorithm minimizes the empirical version of the loss function $\mathcal{L}_\mathcal{D}^{\text{BCE}}(\theta, \lambda, c)$, where

$$\mathcal{L}_\mathcal{D}^{\text{BCE}}(\theta, \lambda, c) = \frac{-1}{N} \sum_{(q_{i,1}, q_{i,2}, p_i) \in \mathcal{D}} (p_i \log P_{\theta, \lambda, c}(q_{i,1} = q_{i,2}) + \bar{p}_i \log \bar{P}_{\theta, \lambda, c}(q_{i,1} = q_{i,2})).$$

85 Let $(\hat{\theta}, \hat{\lambda}, \hat{c})$ denote the minimizer of the empirical loss, i.e.,

$$(\hat{\theta}, \hat{\lambda}, \hat{c}) \in \arg \min_{\theta \in \Theta, \lambda \in \Lambda, c \in \mathcal{C}} \mathcal{L}_\mathcal{D}^{\text{BCE}}(\theta, \lambda, c).$$

86 The following theorem provides a finite sample guarantee of the convergence rate of the empirical
87 minimizer:

² $\mathcal{P}(\cdot | q_{i,1}, q_{i,2}) \in \Delta([0, 1])$ can be any distribution on $[0, 1]$.

88 **Theorem 3.1** (Convergence rate of the main algorithm, BCE loss). *Under Assumptions 2.3 and 2.4,*
 89 *for any $\delta > 0$, with probability at least $1 - \delta$, it holds that*

$$\mathbb{E}_{(q_1, q_2) \sim \mu} \left[\left| P^*(q_1 = q_2) - P_{\hat{\theta}, \hat{\lambda}, \hat{c}}(q_1 = q_2) \right| \right] \leq O \left(\frac{\sqrt{L_\lambda^{-1} + B_c} \cdot (\log(1/\delta))^{1/4}}{N^{1/4}} \right).$$

90 The proof of Theorem 3.1 is deferred to Appendix B.1.

91 3.2 Convergence guarantee for squared log difference loss

92 In this section, we analyze the convergence rate for another loss function. Define the squared log
 93 difference loss function as follows:

$$\mathcal{L}_\mu^{\text{sld}}(\theta, \lambda, c) = \mathbb{E}_{(q_1, q_2) \sim \mu} \left[(\log P^*(q_1 = q_2) - \log P_{\theta, \lambda, c}(q_1 = q_2))^2 \right]. \quad (3)$$

94 Similarly, we also define the empirical squared log difference loss as

$$\mathcal{L}_\mu^{\text{sld}}(\theta, \lambda, c) = \frac{1}{N} \sum_{(q_{i,1}, q_{i,2}, p_i) \in \mathcal{D}} (\log p_i - \log P_{\theta, \lambda, c}(q_{i,1} = q_{i,2}))^2. \quad (4)$$

95 Since $\log 0$ is not well-defined, for theoretical analysis, we assume that each p_i in the dataset \mathcal{D}
 96 satisfies $p_i = P^*(q_{i,1} = q_{i,2})$, i.e., the label is exact the ground-truth probability.

97 Now, we provide a finite sample convergence guarantee for the squared log difference loss:

98 **Theorem 3.2** (Convergence rate of the main algorithm, squared log difference loss). *Under Assump-*
 99 *tions 2.3 and 2.4, for any $\delta > 0$, with probability at least $1 - \delta$, it holds that*

$$\mathbb{E}_{(q_1, q_2) \sim \mu} \left[\left| P^*(q_1 = q_2) - P_{\hat{\theta}, \hat{\lambda}, \hat{c}}(q_1 = q_2) \right| \right] \leq O \left(\frac{(L_\lambda^{-1} + B_c) \cdot (\log(1/\delta))^{1/4}}{N^{1/4}} \right).$$

100 The proof of Theorem 3.2 is deferred to Appendix B.2.

101 **Remark 3.3.** *Compared to the bound for BCE loss, there is an additional $\sqrt{L_\lambda^{-1} + B_c}$ factor in the*
 102 *bound for squared log difference loss.*

103 4 Experiments

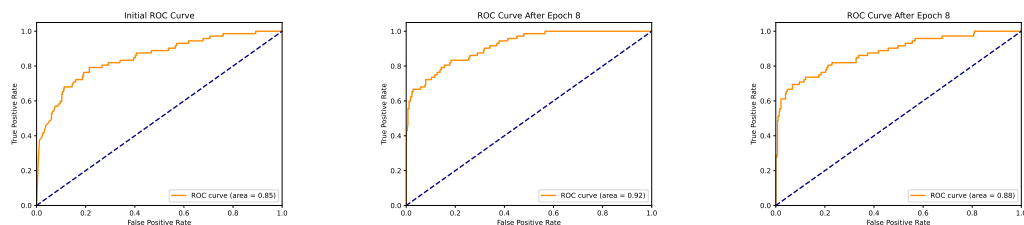
104 In this section, we show experimental results that our distillation-based method can indeed improve
 105 the accuracy of caching prediction.

106 **Construction of the dataset.** We first extract all prompts (queries) from the natural_questions
 107 dataset [Kwiatkowski et al., 2019] and compute a vector embedding for each prompt using the last
 108 layer of the intfloat/e5-large-v2 model [Wang et al., 2022]. After deleting repeated prompts, for each
 109 prompt, we search the five nearest neighbors using FAISS [Johnson et al., 2019]. We sample 1999
 110 prompts uniformly at random, and for each prompt, we choose the farthest three prompts³ among
 111 the five nearest neighbors to form three prompt pairs. Therefore, we get 5997 prompt pairs in total,
 112 and we use GPT-4 [OpenAI, 2023] to label whether each prompt pair can be answered by the same
 113 response (0 or 1). We split the dataset into a training set of size 5497 and a validation set of size 500.

114 **Fine-tuning of embeddings.** We fine-tune using cross-entropy loss or squared log difference loss
 115 from the embedding of Wang et al. [2022]. Slightly different from the theoretical version, we view λ
 116 and c as hyper-parameters and set them to $\lambda = 0.01$, $c = 80$. For squared log difference loss, we clip
 117 the label to $[10^{-10}, 1]$ to avoid calculating $\log 0$, which is not well-defined. We set the learning rate
 118 to be 10^{-5} and present the ROC curve as well as AUC on the validation set of the initial embedding
 119 and embeddings fine-tuned for eight epochs using two loss functions respectively in Figure 2.

120 As Figure 2 shows, fine-tuning on our constructed dataset using either loss function can help to
 121 improve the AUC, while the cross-entropy loss function shows a better performance than the squared
 122 log difference loss function, which is consistent with our theoretical results in Section 3.

³We choose the farthest three to construct a more “difficult” dataset. Note that for each prompt, the farthest three might still be close to it but cannot be answered by the same response.



(a) ROC before fine-tuning (AUC = 0.85)

(b) ROC after Epoch 8 using cross-entropy loss (AUC = 0.92)

(c) ROC after Epoch 8 using squared loss (AUC = 0.88)

Figure 2: Comparison of ROC curves. Both loss functions can help to improve the AUC, while the cross-entropy loss function shows a better performance than the squared log difference loss function, which is consistent with our theoretical results in Section 3.

123 5 Conclusions

124 In this paper, we study efficient prompt caching for LLMs by modeling the ground-truth probability
 125 of whether a prompt pair can be answered by the same response via embedding similarity, and fine-
 126 tuning existing semantic embeddings on our newly constructed dataset. We provide both theoretical
 127 guarantee and empirical evidence that our proposed distillation-based method can improve the
 128 accuracy of caching prediction. Interesting future directions include improving the $O(1/N^{1/4})$ rate
 129 and simulating the caching procedure using our fine-tuned embeddings.

130 References

131 Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos,
 132 Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, Eric Chu, Jonathan H. Clark,
 133 Laurent El Shafey, Yanping Huang, Kathy Meier-Hellstern, Gaurav Mishra, Erica Moreira, Mark
 134 Omernick, Kevin Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao, Yuanzhong Xu, Yujing Zhang,
 135 Gustavo Hernandez Abrego, Junwhan Ahn, Jacob Austin, Paul Barham, Jan Botha, James Bradbury,
 136 Siddhartha Brahma, Kevin Brooks, Michele Catasta, Yong Cheng, Colin Cherry, Christopher A.
 137 Choquette-Choo, Aakanksha Chowdhery, Cl ement Crepy, Shachi Dave, Mostafa Dehghani, Sunipa
 138 Dev, Jacob Devlin, Mark Diaz, Nan Du, Ethan Dyer, Vlad Feinberg, Fangxiaoyu Feng, Vlad
 139 Fienber, Markus Freitag, Xavier Garcia, Sebastian Gehrmann, Lucas Gonzalez, Guy Gur-Ari,
 140 Steven Hand, Hadi Hashemi, Le Hou, Joshua Howland, Andrea Hu, Jeffrey Hui, Jeremy Hurwitz,
 141 Michael Isard, Abe Ittycheriah, Matthew Jagielski, Wenhao Jia, Kathleen Kenealy, Maxim Krikun,
 142 Sneha Kudugunta, Chang Lan, Katherine Lee, Benjamin Lee, Eric Li, Music Li, Wei Li, YaGuang
 143 Li, Jian Li, Hyeontaek Lim, Hanzhao Lin, Zhongtao Liu, Frederick Liu, Marcello Maggioni,
 144 Aroma Mahendru, Joshua Maynez, Vedant Misra, Maysam Moussalem, Zachary Nado, John
 145 Nham, Eric Ni, Andrew Nystrom, Alicia Parrish, Marie Pellat, Martin Polacek, Alex Polozov,
 146 Reiner Pope, Siyuan Qiao, Emily Reif, Bryan Richter, Parker Riley, Alex Castro Ros, Aurko Roy,
 147 Brennan Saeta, Rajkumar Samuel, Renee Shelby, Ambrose Slone, Daniel Smilkov, David R. So,
 148 Daniel Sohn, Simon Tokumine, Dasha Valter, Vijay Vasudevan, Kiran Vodrahalli, Xuezhi Wang,
 149 Pidong Wang, Zirui Wang, Tao Wang, John Wieting, Yuhuai Wu, Kelvin Xu, Yunhan Xu, Linting
 150 Xue, Pengcheng Yin, Jiahui Yu, Qiao Zhang, Steven Zheng, Ce Zheng, Weikang Zhou, Denny
 151 Zhou, Slav Petrov, and Yonghui Wu. Palm 2 technical report, 2023.

152 Hannah Bast, Bj orn Buchhold, Elmar Haussmann, et al. Semantic search on text and knowledge
 153 bases. *Foundations and Trends  in Information Retrieval*, 10(2-3):119–271, 2016.

154 Edward Beeching, Younes Belkada, Kashif Rasul, Lewis Tunstall, Leandro von Werra, Nazneen
 155 Rajani, and Nathan Lambert. Stackllama: An rl fine-tuned llama model for stack exchange question
 156 and answering, 2023. URL <https://huggingface.co/blog/stackllama>.

157 Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx,
 158 Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson,
 159 Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri Chatterji, Annie Chen, Kathleen Creel,
 160 Jared Quincy Davis, Dora Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano

- 161 Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren
162 Gillespie, Karan Goel, Noah Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter
163 Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil
164 Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, Omar
165 Khattab, Pang Wei Koh, Mark Krass, Ranjay Krishna, Rohith Kuditipudi, Ananya Kumar, Faisal
166 Ladhak, Mina Lee, Tony Lee, Jure Leskovec, Isabelle Levent, Xiang Lisa Li, Xuechen Li, Tengyu
167 Ma, Ali Malik, Christopher D. Manning, Suvir Mirchandani, Eric Mitchell, Zanele Munyikwa,
168 Suraj Nair, Avaniika Narayan, Deepak Narayanan, Ben Newman, Allen Nie, Juan Carlos Niebles,
169 Hamed Nilforoshan, Julian Nyarko, Giray Ogut, Laurel Orr, Isabel Papadimitriou, Joon Sung
170 Park, Chris Piech, Eva Portelance, Christopher Potts, Aditi Raghunathan, Rob Reich, Hongyu
171 Ren, Frieda Rong, Yusuf Roohani, Camilo Ruiz, Jack Ryan, Christopher Ré, Dorsa Sadigh,
172 Shiori Sagawa, Keshav Santhanam, Andy Shih, Krishnan Srinivasan, Alex Tamkin, Rohan Taori,
173 Armin W. Thomas, Florian Tramèr, Rose E. Wang, William Wang, Bohan Wu, Jiajun Wu, Yuhuai
174 Wu, Sang Michael Xie, Michihiro Yasunaga, Jiaxuan You, Matei Zaharia, Michael Zhang, Tianyi
175 Zhang, Xikun Zhang, Yuhui Zhang, Lucia Zheng, Kaitlyn Zhou, and Percy Liang. On the
176 opportunities and risks of foundation models, 2022.
- 177 Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican,
178 George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al.
179 Improving language models by retrieving from trillions of tokens. In *International conference on*
180 *machine learning*, pages 2206–2240. PMLR, 2022.
- 181 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar,
182 Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence:
183 Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- 184 Archana Bura, Desik Rengarajan, Dileep Kalathil, Srinivas Shakkottai, and Jean-Francois Chamber-
185 land. Learning to cache and caching to learn: Regret analysis of caching algorithms. *IEEE/ACM*
186 *Transactions on Networking*, 30(1):18–31, 2021.
- 187 Wei-Cheng Chang, Felix X Yu, Yin-Wen Chang, Yiming Yang, and Sanjiv Kumar. Pre-training tasks
188 for embedding-based large-scale retrieval. *arXiv preprint arXiv:2002.03932*, 2020.
- 189 Zheng Chang, Lei Lei, Zhenyu Zhou, Shiwen Mao, and Tapani Ristaniemi. Learn to cache: Machine
190 learning for network edge caching in the big data era. *IEEE Wireless Communications*, 25(3):
191 28–35, 2018.
- 192 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam
193 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm:
194 Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- 195 Fathima Zarin Faizal, Priya Singh, Nikhil Karamchandani, and Sharayu Moharir. Regret-optimal
196 online caching for adversarial and stochastic arrivals. In *EAI International Conference on Perfor-*
197 *mance Evaluation Methodologies and Tools*, pages 147–163. Springer, 2022.
- 198 Edouard Grave, Armand Joulin, and Nicolas Usunier. Improving neural language models with a
199 continuous cache. *arXiv preprint arXiv:1612.04426*, 2016.
- 200 Edouard Grave, Moustapha M Cisse, and Armand Joulin. Unbounded cache model for online
201 language modeling with open vocabulary. *Advances in neural information processing systems*, 30,
202 2017.
- 203 Ying He, Zheng Zhang, F Richard Yu, Nan Zhao, Hongxi Yin, Victor CM Leung, and Yanhua Zhang.
204 Deep-reinforcement-learning-based optimization for cache-enabled opportunistic interference
205 alignment wireless networks. *IEEE Transactions on Vehicular Technology*, 66(11):10433–10445,
206 2017.
- 207 Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane
208 Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. Few-shot learning with
209 retrieval augmented language models. *arXiv preprint arXiv:2208.03299*, 2022.
- 210 Wei Jiang, Gang Feng, Shuang Qin, Tak Shing Peter Yum, and Guohong Cao. Multi-agent reinforce-
211 ment learning for efficient content caching in mobile d2d networks. *IEEE Transactions on Wireless*
212 *Communications*, 18(3):1610–1622, 2019.

- 213 Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with GPUs. *IEEE*
214 *Transactions on Big Data*, 7(3):535–547, 2019.
- 215 Ehsan Kamaloo, Xinyu Zhang, Odunayo Ogundepo, Nandan Thakur, David Alfonso-Hermelo,
216 Mehdi Rezagholizadeh, and Jimmy Lin. Evaluating embedding apis for information retrieval.
217 *arXiv preprint arXiv:2305.06300*, 2023.
- 218 Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. Generalization
219 through memorization: Nearest neighbor language models. *arXiv preprint arXiv:1911.00172*,
220 2019.
- 221 Swadhesh Kumar and PK Singh. An overview of modern cache memory and performance analysis
222 of replacement policies. In *2016 IEEE International Conference on Engineering and Technology*
223 *(ICETECH)*, pages 210–214. IEEE, 2016.
- 224 Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris
225 Alberti, Danielle Epstein, Illia Polosukhin, Matthew Kelcey, Jacob Devlin, Kenton Lee, Kristina N.
226 Toutanova, Llion Jones, Ming-Wei Chang, Andrew Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov.
227 Natural questions: a benchmark for question answering research. *Transactions of the Association*
228 *of Computational Linguistics*, 2019.
- 229 Donghee Lee, Jongmoo Choi, Jong-Hun Kim, Sam H Noh, Sang Lyul Min, Yookun Cho, and
230 Chong Sang Kim. Lrfu: A spectrum of policies that subsumes the least recently used and least
231 frequently used policies. *IEEE transactions on Computers*, 50(12):1352–1361, 2001.
- 232 Sewon Min, Weijia Shi, Mike Lewis, Xilun Chen, Wen-tau Yih, Hannaneh Hajishirzi, and Luke
233 Zettlemoyer. Nonparametric masked language modeling. *arXiv preprint arXiv:2212.01349*, 2022.
- 234 Samrat Mukhopadhyay and Abhishek Sinha. Online caching with optimal switching regret. In *2021*
235 *IEEE International Symposium on Information Theory (ISIT)*, pages 1546–1551. IEEE, 2021.
- 236 Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. Capabilities
237 of gpt-4 on medical challenge problems. *arXiv preprint arXiv:2303.13375*, 2023.
- 238 OpenAI. Gpt-4 technical report, 2023.
- 239 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
240 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow
241 instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:
242 27730–27744, 2022.
- 243 David Patterson, Joseph Gonzalez, Quoc Le, Chen Liang, Lluís-Miquel Munguia, Daniel Rothchild,
244 David So, Maud Texier, and Jeff Dean. Carbon emissions and large neural network training. *arXiv*
245 *preprint arXiv:2104.10350*, 2021.
- 246 Or Sharir, Barak Peleg, and Yoav Shoham. The cost of training nlp models: A concise overview.
247 *arXiv preprint arXiv:2004.08900*, 2020.
- 248 Junaid Shuja, Kashif Bilal, Waleed Alasmay, Hassan Sinky, and Eisa Alanazi. Applying machine
249 learning techniques for caching in next-generation edge networks: A comprehensive survey.
250 *Journal of Network and Computer Applications*, 181:103005, 2021.
- 251 Alan Jay Smith. Cache memories. *ACM Computing Surveys (CSUR)*, 14(3):473–530, 1982.
- 252 William Stallings. *Operating systems: internals and design principles*. Prentice Hall Press, 2011.
- 253 Jia Wang. A survey of web caching schemes for the internet. *ACM SIGCOMM Computer Communi-*
254 *cation Review*, 29(5):36–46, 1999.
- 255 Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder,
256 and Furu Wei. Text embeddings by weakly-supervised contrastive pre-training. *arXiv preprint*
257 *arXiv:2212.03533*, 2022.

- 258 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama,
259 Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models.
260 *arXiv preprint arXiv:2206.07682*, 2022.
- 261 Wenhao Zhan, Baihe Huang, Audrey Huang, Nan Jiang, and Jason Lee. Offline reinforcement
262 learning with realizability and single-policy concentrability. In *Conference on Learning Theory*,
263 pages 2730–2775. PMLR, 2022.
- 264 Zexuan Zhong, Tao Lei, and Danqi Chen. Training language models with memory augmentation.
265 *arXiv preprint arXiv:2205.12674*, 2022.
- 266 Banghua Zhu, Ying Sheng, Lianmin Zheng, Clark Barrett, Michael I Jordan, and Jiantao Jiao. On op-
267 timal caching and model multiplexing for large model inference. *arXiv preprint arXiv:2306.02003*,
268 2023.
- 269 Hanlin Zhu and Amy Zhang. Provably efficient offline goal-conditioned reinforcement learning with
270 general function approximation and single-policy concentrability. *arXiv preprint arXiv:2302.03770*,
271 2023.
- 272 Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul
273 Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv*
274 *preprint arXiv:1909.08593*, 2019.
- 275 zilliztech. Gptcache: Semantic cache for llms. fully integrated with langchain and llama_index.,
276 2023. URL <https://github.com/zilliztech/GPTCache>.

277 **A Related Works**

278 **Caching.** Caching algorithms are important to computer architecture and systems and have long
 279 been explored [Lee et al., 2001, Stallings, 2011, Bura et al., 2021]. In recent years, caching has also
 280 been applied to online learning analysis and machine learning advice [He et al., 2017, Chang et al.,
 281 2018, Jiang et al., 2019, Shuja et al., 2021, Mukhopadhyay and Sinha, 2021, Faizal et al., 2022]. Zhu
 282 et al. [2023] is the most related work and studies optimal caching algorithm for prompt in both online
 283 and offline learning settings. Instead of studying caching policy, we aim to study how to efficiently
 284 search and accurately predict whether there is a caching hit.

285 **Retrieval-based LLMs.** A line of work studies augmenting a language model by retrieval-based
 286 method [Grave et al., 2016, 2017, Khandelwal et al., 2019, Borgeaud et al., 2022, Izacard et al.,
 287 2022, Zhong et al., 2022, Min et al., 2022]. For example, the kNN-LM model [Khandelwal et al.,
 288 2019] interpolates a distribution obtained by the vector embedding of k nearest neighbors with the
 289 distribution of language models. Our formulation of probability via embedding similarity is inspired
 290 by these works.

291 **B Missing Proofs**

292 **B.1 Proof of Theorem 3.1**

293 To prove Theorem 3.1, we first present and prove Lemmas B.1 and B.2. Our proof strategy is similar
 294 to that of Zhan et al. [2022], Zhu and Zhang [2023].

295 **Lemma B.1.** *Under Assumptions 2.3 and 2.4, for any $\delta > 0$, with probability at least $1 - \delta$, it holds*
 296 *that*

$$|\mathcal{L}_{\mathcal{D}}^{\text{BCE}}(\theta, \lambda, c) - \mathcal{L}_{\mu}^{\text{BCE}}(\theta, \lambda, c)| \leq O\left((L_{\lambda}^{-1} + B_c)\sqrt{\frac{\log(|\Theta||\Lambda||\mathcal{C}|/\delta)}{N}}\right) \triangleq \epsilon_{\text{stat}}^{\text{BCE}}$$

297 for any $\theta \in \Theta, \lambda \in \Lambda, c \in \mathcal{C}$.

298 *Proof.* We first consider any fixed $\theta \in \Theta, \lambda \in \Lambda, c \in \mathcal{C}$. One can observe that $\mathcal{L}_{\mathcal{D}}^{\text{BCE}}(\theta, \lambda, c)$ is an
 299 unbiased estimator of $\mathcal{L}_{\mu}^{\text{BCE}}(\theta, \lambda, c)$ since

$$\begin{aligned} & \mathbb{E}[\mathcal{L}_{\mathcal{D}}^{\text{BCE}}(\theta, \lambda, c)] \\ &= -\mathbb{E}_{(q_1, q_2) \sim \mu, p \sim \mathcal{P}(\cdot|q_1, q_2)} [p \log P_{\theta, \lambda, c}(q_1 = q_2) + \bar{p} \log \bar{P}_{\theta, \lambda, c}(q_1 = q_2)] \\ &= -\mathbb{E}_{(q_1, q_2) \sim \mu} [\mathbb{E}_{p \sim \mathcal{P}(\cdot|q_1, q_2)} [p \log P_{\theta, \lambda, c}(q_1 = q_2) + \bar{p} \log \bar{P}_{\theta, \lambda, c}(q_1 = q_2) | q_1, q_2]] \\ &= -\mathbb{E}_{(q_1, q_2) \sim \mu} [P^*(q_1 = q_2) \log P_{\theta, \lambda, c}(q_1 = q_2) + \bar{P}^*(q_1 = q_2) \log \bar{P}_{\theta, \lambda, c}(q_1 = q_2)] \\ &= \mathcal{L}_{\mu}^{\text{BCE}}(\theta, \lambda, c), \end{aligned}$$

300 where the first equality holds since the data points in the dataset are i.i.d. distributed, the second
 301 equality holds due to tower property, and the third equality holds by the linearity of expectation.

302 Also, we note that the empirical loss for each data point can be upper bounded by

$$\begin{aligned} & |p_i \log P_{\theta, \lambda, c}(q_{i,1} = q_{i,2}) + \bar{p}_i \log \bar{P}_{\theta, \lambda, c}(q_{i,1} = q_{i,2})| \\ & \leq \max\{|\log P_{\theta, \lambda, c}(q_{i,1} = q_{i,2})|, |\log \bar{P}_{\theta, \lambda, c}(q_{i,1} = q_{i,2})|\} \\ & = \log\left(\max\left\{\frac{1}{P_{\theta, \lambda, c}(q_{i,1} = q_{i,2})}, \frac{1}{1 - P_{\theta, \lambda, c}(q_{i,1} = q_{i,2})}\right\}\right). \end{aligned}$$

303 Since $\sigma(-x) = 1 - \sigma(x)$ and $\text{sim}(v_{\theta}(q_1), v_{\theta}(q_2))/\lambda - c \in [-L_{\lambda}^{-1} - B_c, L_{\lambda}^{-1} + B_c]$ by Assump-
 304 tion 2.4, we can obtain that

$$\frac{1}{P_{\theta, \lambda, c}(q_1 = q_2)} = \frac{1}{\sigma(\text{sim}(v_{\theta}(q_1), v_{\theta}(q_2))/\lambda - c)} \leq 1 + \exp(L_{\lambda}^{-1} + B_c).$$

305 By the symmetry of $\sigma(\cdot)$ and the range of $\text{sim}(v_{\theta}(q_1), v_{\theta}(q_2))/\lambda - c$, it also holds that

$$\frac{1}{1 - P_{\theta, \lambda, c}(q_1 = q_2)} \leq 1 + \exp(L_{\lambda}^{-1} + B_c).$$

306 Therefore,

$$\begin{aligned} & |p_i \log P_{\theta, \lambda, c}(q_{i,1} = q_{i,2}) + \bar{p}_i \log \bar{P}_{\theta, \lambda, c}(q_{i,1} = q_{i,2})| \\ & \leq \log(1 + \exp(L_\lambda^{-1} + B_c)) \leq O(L_\lambda^{-1} + B_c). \end{aligned}$$

307 By Hoeffding's inequality, we have with probability at least $1 - \delta$, it holds that

$$|\mathcal{L}_D^{\text{BCE}}(\theta, \lambda, c) - \mathcal{L}_\mu^{\text{BCE}}(\theta, \lambda, c)| \leq O\left((L_\lambda^{-1} + B_c) \sqrt{\frac{\log(1/\delta)}{N}}\right).$$

308 Applying a union bound over all $\theta \in \Theta, \lambda \in \Lambda, c \in \mathcal{C}$ concludes the result. \square

309 **Lemma B.2.** *Under Assumptions 2.3 and 2.4, with probability at least $1 - \delta$, it holds that*

$$\mathcal{L}_\mu^{\text{BCE}}(\hat{\theta}, \hat{\lambda}, \hat{c}) - \mathcal{L}_\mu^{\text{BCE}}(\theta^*, \lambda^*, c^*) \leq 2\epsilon_{\text{stat}}^{\text{BCE}}.$$

310 where $\epsilon_{\text{stat}}^{\text{BCE}}$ is defined in Lemma B.1 and

$$(\hat{\theta}, \hat{\lambda}, \hat{c}) \in \arg \min_{\theta \in \Theta, \lambda \in \Lambda, c \in \mathcal{C}} \mathcal{L}_D^{\text{BCE}}(\theta, \lambda, c).$$

311 *Proof.* We condition on the high probability event in Lemma B.1. Note that

$$\begin{aligned} & \mathcal{L}_\mu^{\text{BCE}}(\hat{\theta}, \hat{\lambda}, \hat{c}) - \mathcal{L}_\mu^{\text{BCE}}(\theta^*, \lambda^*, c^*) \\ & = \underbrace{\mathcal{L}_\mu^{\text{BCE}}(\hat{\theta}, \hat{\lambda}, \hat{c}) - \mathcal{L}_D^{\text{BCE}}(\hat{\theta}, \hat{\lambda}, \hat{c})}_{(1)} + \underbrace{\mathcal{L}_D^{\text{BCE}}(\hat{\theta}, \hat{\lambda}, \hat{c}) - \mathcal{L}_D^{\text{BCE}}(\theta^*, \lambda^*, c^*)}_{(2)} \\ & \quad + \underbrace{\mathcal{L}_D^{\text{BCE}}(\theta^*, \lambda^*, c^*) - \mathcal{L}_\mu^{\text{BCE}}(\theta^*, \lambda^*, c^*)}_{(3)}. \end{aligned}$$

312 (1), (3) $\leq \epsilon_{\text{stat}}^{\text{BCE}}$ by Lemma B.1 and (2) ≤ 0 by the optimality of $(\hat{\theta}, \hat{\lambda}, \hat{c})$, which completes the
313 proof. \square

314 Equipped with Lemmas B.1 and B.2, we are now able to prove Theorem 3.1.

315 *Proof of Theorem 3.1.* We condition on the high probability event in Lemma B.2. Note that by the
316 realizability of the ground-truth probability (Assumption 2.3) and the property of binary cross-entropy,
317 we have

$$(\theta^*, \lambda^*, c^*) \in \arg \min_{\theta \in \Theta, \lambda \in \Lambda, c \in \mathcal{C}} \mathcal{L}_\mu^{\text{BCE}}(\theta, \lambda, c).$$

318 For any $\theta \in \Theta, \lambda \in \Lambda, c \in \mathcal{C}$, we map (θ, λ, c) to a function $f_{\theta, \lambda, c}(\cdot, \cdot) : \mathcal{Q} \times \mathcal{Q} \rightarrow [0, 1]$, where

$$f_{\theta, \lambda, c}(q_1, q_2) = P_{\theta, \lambda, c}(q_1 = q_2), \quad \forall q_1, q_2 \in \mathcal{Q}.$$

319 Moreover, we define functional h s.t. for any function $f : \mathcal{Q} \times \mathcal{Q} \rightarrow [0, 1]$,

$$h(f) = -\mathbb{E}_{(q_1, q_2) \sim \mu} [P^*(q_1 = q_2) \log f(q_1, q_2) + \bar{P}^*(q_1 = q_2) \log(1 - f(q_1, q_2))].$$

320 By the definition of BCE loss, $h(f_{\theta, \lambda, c}) = \mathcal{L}_\mu^{\text{BCE}}(\theta, \lambda, c)$. Therefore, Lemma B.2 translates to

$$h(f_{\hat{\theta}, \hat{\lambda}, \hat{c}}) - h(f_{\theta^*, \lambda^*, c^*}) \leq 2\epsilon_{\text{stat}}^{\text{BCE}}. \quad (5)$$

321 Note that $f_{\theta^*, \lambda^*, c^*}$ is still the minimizer of $h(f)$ even if f cannot be induced by some θ, λ, c .

322 We also observe that $h(f)$ is 1-strongly convex w.r.t. f in $\|\cdot\|_{2, \mu}$ norm. To see why this is the case,
323 one can calculate the second-order derivative of h w.r.t. $f(q_1, q_2)$ for any $(q_1, q_2) \in \mathcal{Q} \times \mathcal{Q}$, which is

$$\frac{\partial^2 h}{\partial f^2}(q_1, q_2) = \frac{P^*(q_1 = q_2)}{f^2(q_1, q_2)} + \frac{\bar{P}^*(q_1 = q_2)}{(1 - f(q_1, q_2))^2} \geq P^*(q_1 = q_2) + \bar{P}^*(q_1 = q_2) = 1,$$

324 which demonstrates the strong convexity. Therefore, by strong convexity and the optimality of
 325 $f_{\theta^*, \lambda^*, c^*}$, we can obtain that

$$h(f_{\hat{\theta}, \hat{\lambda}, \hat{c}}) \geq h(f_{\theta^*, \lambda^*, c^*}) + \frac{1}{2} \|f_{\theta^*, \lambda^*, c^*} - f_{\hat{\theta}, \hat{\lambda}, \hat{c}}\|_{2, \mu}^2.$$

326 Combining (5), we have

$$\|f_{\theta^*, \lambda^*, c^*} - f_{\hat{\theta}, \hat{\lambda}, \hat{c}}\|_{2, \mu} \leq 2\sqrt{\epsilon_{\text{stat}}^{\text{BCE}}}.$$

327 Finally, by Cauchy–Schwarz inequality, we can conclude

$$\begin{aligned} & \mathbb{E}_{(q_1, q_2) \sim \mu} \left[\left| P^*(q_1 = q_2) - P_{\hat{\theta}, \hat{\lambda}, \hat{c}}(q_1 = q_2) \right| \right] \\ &= \|f_{\theta^*, \lambda^*, c^*} - f_{\hat{\theta}, \hat{\lambda}, \hat{c}}\|_{1, \mu} \leq \|f_{\theta^*, \lambda^*, c^*} - f_{\hat{\theta}, \hat{\lambda}, \hat{c}}\|_{2, \mu} \leq 2\sqrt{\epsilon_{\text{stat}}^{\text{BCE}}}. \end{aligned}$$

328

□

329 B.2 Proof of Theorem 3.2

330 *Proof of Theorem 3.2.* The proof is similar to the proof of Theorem 3.1. First, it is easy to see
 331 that the empirical loss $\mathcal{L}_{\mathcal{D}}^{\text{sld}}(\theta, \lambda, c)$ is an unbiased estimator of $\mathcal{L}_{\mu}^{\text{sld}}(\theta, \lambda, c)$ by definition since
 332 $p_i = P^*(q_{i,1} = q_{i,2})$. Also,

$$(\log p_i - \log P_{\theta, \lambda, c}(q_{i,1} = q_{i,2}))^2 \leq (\log(1 + \exp(L_{\lambda}^{-1} + B_c)))^2 = O((L_{\lambda}^{-1} + B_c)^2).$$

333 Therefore, by Hoeffding’s inequality and union bound, we have that with probability at least $1 - \delta$, it
 334 holds that

$$|\mathcal{L}_{\mathcal{D}}^{\text{sld}}(\theta, \lambda, c) - \mathcal{L}_{\mu}^{\text{sld}}(\theta, \lambda, c)| \leq O\left((L_{\lambda}^{-1} + B_c)^2 \sqrt{\frac{\log(|\Theta||\Lambda||\mathcal{C}|/\delta)}{N}}\right) \triangleq \epsilon_{\text{stat}}^{\text{sld}}.$$

335 Similar to Lemma B.2, we can obtain that

$$\mathcal{L}_{\mu}^{\text{sld}}(\hat{\theta}, \hat{\lambda}, \hat{c}) - \mathcal{L}_{\mu}^{\text{sld}}(\theta^*, \lambda^*, c^*) \leq 2\epsilon_{\text{stat}}^{\text{sld}}. \quad (6)$$

336 Now, for any $\theta \in \Theta, \lambda \in \Lambda, c \in \mathcal{C}$, we map (θ, λ, c) to a function $f_{\theta, \lambda, c}(\cdot, \cdot) : \mathcal{Q} \times \mathcal{Q} \rightarrow [0, +\infty)$,
 337 where

$$f_{\theta, \lambda, c}(q_1, q_2) = -\log P_{\theta, \lambda, c}(q_1 = q_2), \quad \forall q_1, q_2 \in \mathcal{Q}.$$

338 Moreover, we define functional h s.t. for any function $f : \mathcal{Q} \times \mathcal{Q} \rightarrow [0, 1]$,

$$h(f) = \mathbb{E}_{(q_1, q_2) \sim \mu} [(f(q_1, q_2) + \log P^*(q_1 = q_2))^2].$$

339 By the definition of squared log difference loss, $h(f_{\theta, \lambda, c}) = \mathcal{L}_{\mu}^{\text{sld}}(\theta, \lambda, c)$. Therefore, (6) translates to

$$h(f_{\hat{\theta}, \hat{\lambda}, \hat{c}}) - h(f_{\theta^*, \lambda^*, c^*}) \leq 2\epsilon_{\text{stat}}^{\text{sld}}. \quad (7)$$

340 Note that $f_{\theta^*, \lambda^*, c^*}$ is still the minimizer of $h(f)$ even if f cannot be induced by some θ, λ, c .

341 We also observe that $h(f)$ is 2-strongly convex w.r.t. f in $\|\cdot\|_{2, \mu}$ norm by calculating the second-order
 342 derivative of h w.r.t. f . Therefore, by strong convexity and the optimality of $f_{\theta^*, \lambda^*, c^*}$, we can obtain
 343 that

$$h(f_{\hat{\theta}, \hat{\lambda}, \hat{c}}) \geq h(f_{\theta^*, \lambda^*, c^*}) + \|f_{\theta^*, \lambda^*, c^*} - f_{\hat{\theta}, \hat{\lambda}, \hat{c}}\|_{2, \mu}^2.$$

344 Combining (7), we can obtain

$$\mathbb{E}_{(q_1, q_2) \sim \mu} \left[(\log P^*(q_1 = q_2) - \log P_{\theta, \lambda, c}(q_1 = q_2))^2 \right] = \|f_{\theta^*, \lambda^*, c^*} - f_{\hat{\theta}, \hat{\lambda}, \hat{c}}\|_{2, \mu}^2 \leq 2\epsilon_{\text{stat}}^{\text{sld}}.$$

345 Note that by mean value theorem, for any $0 < x < y < 1$, $\log x - \log y = (x - y)/z$ for some
 346 $z \in (x, y)$. Therefore, $(\log x - \log y)^2 = (x - y)^2/z^2 > (x - y)^2$. This implies

$$\begin{aligned} & \mathbb{E}_{(q_1, q_2) \sim \mu} \left[(P^*(q_1 = q_2) - P_{\theta, \lambda, c}(q_1 = q_2))^2 \right] \\ & \leq \mathbb{E}_{(q_1, q_2) \sim \mu} \left[(\log P^*(q_1 = q_2) - \log P_{\theta, \lambda, c}(q_1 = q_2))^2 \right] = 2\epsilon_{\text{stat}}^{\text{sld}}. \end{aligned}$$

347 By Cauchy–Schwarz inequality, we can conclude

$$\begin{aligned} & \mathbb{E}_{(q_1, q_2) \sim \mu} \left[\left| P^*(q_1 = q_2) - P_{\hat{\theta}, \hat{\lambda}, \hat{c}}(q_1 = q_2) \right| \right] \\ & \leq \sqrt{\mathbb{E}_{(q_1, q_2) \sim \mu} \left[(P^*(q_1 = q_2) - P_{\theta, \lambda, c}(q_1 = q_2))^2 \right]} \leq \sqrt{2\epsilon_{\text{stat}}^{\text{slid}}}. \end{aligned}$$

348

□