DIFFERENTIALLY PRIVATE RETRIEVAL AUGMENTED GENERATION WITH RANDOM PROJECTION

Dixi Yao, Tian Li

Department of Computer Science University of Chicago Chicago, IL 60637, USA {dixi,litian}@uchicago.edu

Abstract

Large Language Models (LLMs) have gained widespread interest and driven advancements across various fields. Retrieval-Augmented Generation (RAG) enables LLMs to incorporate domain-specific knowledge without retraining. However, evidence shows that RAG poses significant privacy risks due to leakage of sensitive information stored in the retrieval database. In this work, we focus on the notion of differential privacy (DP) and propose a private randomized mechanism to project both the queries and the datastore into a lower-dimensional space using Gaussian matrices, while preserving the similarities for effective retrieval. Empirical evaluation on different RAG architectures demonstrates that our solution achieves strong privacy protection with $\epsilon \approx 5$, and negligible impact on generation performance and latency compared to prior methods.

1 INTRODUCTION

Large Language Models (LLMs) excel across diverse applications. Retrieval-Augmented Generation (RAG) (Khandelwal et al., 2019; Lewis et al., 2020; Min et al., 2023; Edge et al., 2024) enhances LLMs by enabling them to answer domain-specific questions using external knowledge without additional training, ensuring easy deployment. However, recent studies reveal that sensitive information in the knowledge database is vulnerable to leakage if attackers craft specific prompts (Huang et al., 2023; Zeng et al., 2024b; Koga et al., 2024). For example, a company's user data, such as emails and names, stored in RAG documents, can be exposed through API interactions.

Despite the growing adoption of RAG and the simplicity of these attacks, effective countermeasures are lacking. Attacks mimic ordinary user queries, and simple access control may block legitimate users. This challenge motivates us to leverage Differential Privacy (DP) (Dwork et al., 2006), allowing us to answer queries as usual while replace sensitive information like personal emails and phone numbers with other information hard to differentiate.

This paper addresses two popular RAG scenarios: KNN-LM (Khandelwal et al., 2019) and direct prompting with retrieval outputs (Min et al., 2023). KNN-LM enhances next-word prediction by combining LLM probabilities with those from retrieved nearest-neighbor entries. Direct prompting feeds the retrieved items together with the query into the LLM to generate responses.

We propose a novel differentially private solution for RAG using random projection. By projecting the private database onto a permuted space via the Johnson-Lindenstrauss transformation (Johnson et al., 1986), we ensure answers exclude sensitive information while retaining essential content. Our key insight is that random projection preserves pairwise similarity on low dimensionality, yet changes embedding values such that attackers extract alternative items instead of the original text. This allows users to receive accurate answers from the RAG while adversaries fail.

We provide formal proof demonstrating that our mechanism satisfies differential privacy requirements. Empirical results show that our method outperforms prior work and direct random projection Li & Li (2023), where data embeddings are projected by a matrix of IID random variables from a standard Gaussian distribution, which is then combined with another Gaussian random matrix. The second matrix ensures differential privacy. Our approach achieves better generation quality and datastore privacy under the same privacy budget, offering a promising direction for private RAG.

2 RELATED WORK

Retrieval Augmented Generation. Retrieval-Augmented Generation (RAG) enhances pre-trained LLMs with non-parametric retrieval of external knowledge. Retrieval methods based on *k*-nearest neighbors (Khandelwal et al., 2019) retrieve relevant database entries, combined with LLMs to produce the final answer. In KNN-LM, the next-word probability combines language model predictions with probabilities from retrieved items. Another RAG architecture (Lewis et al., 2020; Min et al., 2023) directly prompts LLMs with a combination of retrieved items and the original query to generate final responses. We call such methods direct-prompting RAG in this paper. GraphRAG (Edge et al., 2024) builds a graph-based text index and uses it to retrieve from the knowledge base.

Privacy of RAG. LLMs can pose privacy risks due to strong memorization of information. Previous works show that personal data (emails, phone numbers, URLs) and random datastore content can be extracted from KNN-LM models (Huang et al., 2023). Zeng et al. and Jiang et al. confirm similar vulnerabilities in RAG models via direct prompting. Efforts to defend such attacks include DP-based sampling and aggregation (Koga et al., 2024) over tokens and usage of synthetic data (Zeng et al., 2024a). However, these methods reduce efficiency, particularly with large datasets, requiring extra data, computation, and resources.

Random Projection. Random projection, based on the Johnson-Lindenstrauss lemma (Johnson et al., 1986), reduces dimensionality while preserving pairwise distances with high probability. By projecting query and datastore embeddings to lower dimensions, retrieval accuracy remains largely unaffected (Section 5). Additionally, the Johnson-Lindenstrauss transform inherently satisfies differential privacy (Blocki et al., 2012). Prior work applied DP random projection to aggregation (Li & Li, 2023) and image retrieval (Ibrahim et al., 2024). Instead of trivially applying random projection after embeddings in RAG, we propose an algorithm involving matrix projection to perturb data embeddings and prove its differential privacy property.

3 PRELIMINARY

Definition 1 (Renyi Differential Privacy (RDP) (Mironov, 2017)). A randomized mechanism \mathcal{M} : $\mathcal{D} \to S$ satisfies (α, ϵ) -differential privacy if for any neighboring datasets D_1 and D_2 differing by one record, it holds that

$$D_{\alpha}(\mathcal{M}(D_1) \| \mathcal{M}(D_2)) = \frac{1}{\alpha - 1} \log \left(\frac{\mathcal{M}(D_1)}{\mathcal{M}(D_2)} \right)^{\alpha} \le \epsilon$$

In our context, let the original dataset be D_1 , and define w_m is the word of the highest probability being generated by a given query. Removing w_m from D_1 yields a new dataset D_2 , with the only difference being the absence of w_m . The key idea is to design a RAG algorithm that identifies an alternative word w_r instead of w_m while maintaining correct outputs within the answer set S. Such an idea is first enlightened in Huang et al. (2023) and further confirmed in Koga et al. (2024).

4 Method

4.1 PROPOSED APPROACH

Our proposed algorithm is shown in Algorithm 1. In the first step, we calculate the variance of random Gaussian noise to be used in the random projection (detailed in Section 4.2). In Line 3, we only pick the first *n* entries of the document *D* so that the output size after applying DP mechanisms is fixed. A pre-trained language model $f(\cdot)$ serves as the encoder, encoding $D = \{w_1, \dots, w_n\}$ into $f(D) = \{f(w_1), \dots, f(w_n)\}$ (Line 3), which can be further finetuned. But we assume a fixed pre-trained LM throughout the paper for simplicity. Each embedding has a dimension of *d*.

We then randomly generate an IID Gaussian Matrix $R \in \mathbb{R}^{d \times k}$, where each cell is sampled from $\mathcal{N}(0, \sigma^2)$. Next, after normalizing (or clipping) D, we project f(D) to f(D)R (Line 4). To preserve

	Algorithm 1: Algorithm: Differential Privacy Guaranteed						
	Random Projection of RAG-LLM						
	Data: Datastore D, Embedding encoder $f(\cdot)$, $M(x, f(D))$						
1	. Parameters : Gaussian matrix magnitude σ , privacy budget ϵ ,						
	max document entry n.						
	Input: An arbitrary query x.						
	Output: Next work prediction <i>y</i> .						
2	Run once before taking user queries:						
3	Compute γ , Δ , and k using Theorem 1, then generate a $d \times k$ IID						
	random matrix from $\mathcal{N}(0, \sigma^2)$.						
4	Pick the first n entries of the document D , embed D with some						
	encoder $f(\cdot)$, and get $f(D)$.						
5	Normalize and clip $f(D)$ so that each element						
	$\gamma \leq \ f(w_i)\ _2 \leq \Delta, (1 \leq i \leq n).$ Project $f(D) \to f(D)R$.						
6	for each user query x: do						
7	$y \leftarrow M(xR, f(D)R)$	E:					
8	end	F1					



Figure 1: Working flow of KNN-LM and direct prompting RAG.

similarities between queries and datastores, we project the input x to the same dimension of k, using the same matrix R. Then, we calculate the distances using the projected embeddings of the query and dataset for the nearest neighbor search. The following steps are the same as in the common process reflected by calculation by $M(\cdot)$ in Line 6.

Our algorithm is compatible with both KNN-LM and direct prompting in RAG architectures. In KNN-LM, the next-word probability combines logits from a pre-trained LLM and softmax probabilities based on the distance of token candidates to the nearest neighbor. In direct prompting, retrieved texts, such as those from KNN, are combined with queries and input into another LLM to generate a final answer. As illustrated in Fig. 1, the green block denotes our random projection operations, corresponding to line 4 in Algorithm 1. After retrieving embeddings from the datastore, we can either directly output the prediction, yielding results from KNN-LM architectures, or combine the query with the retrieved information and input it into another pre-trained LLM. In the other architecture, we prompt the LLM with the retrieved texts alongside the query to generate the result.

4.2 DIFFERENTIAL PRIVACY GUARANTEE

Theorem 1. Algorithm 1 satisfies (α, ϵ) -differential privacy (Definition 1) if projection dimension k, and normalization bound γ and Δ along with privacy budgets α and ϵ meet following equation:

$$\frac{\Gamma^{1-\alpha}}{(1-\alpha)\Gamma+\alpha} < e^{\frac{2(\alpha-1)\epsilon}{k}}, \forall \Gamma \in \left[\frac{\gamma^2}{\Delta^2}, \frac{\Delta^2}{\gamma^2}\right]$$
(1)

We provide a proof sketch here for better understanding of our theorems and defer the complete proof to Appendix A. We pick any two datasets D_1 and D_2 including cases they have one line difference, D_2 has one line more, and one line less. The idea is that, with some $Z \in DR$, we have $\arg \max p(y|x, Z) \in S$. Back to the definition of DP in our question, the idea is to show that the probabilities of D_1R and D_2R equaling to Z are hard to differentiate (\Pr_{D_1} and \Pr_{D_2} are close).

We assume the distribution of each element in D_1R and D_2R follows i.i.d. normal distribution and they can be expressed as $\mathcal{N}(0, \sigma_1^2)$ and $\mathcal{N}(0, \sigma_2^2)$ respectively. We can convert σ_1^2 to $\sigma^2 ||D_{1_m}||_2^2$ and σ_2^2 to $\sigma^2 ||D_{2_m}||_2^2$. With the condition of normalization, we then further construct the relationship between privacy budgets and parameters to set: normalization bounds and projected dimensions k.

From Theorem 1, several implications arise. Smaller budgets ϵ and δ require a smaller k. Smaller k means a smaller projected dimension and higher compression rate of information. Since embedding vectors are typically sparse, a smaller k helps maintain generation performance while meeting privacy requirements. However, if k is too small, much information will lose. To adopt larger $\frac{\Delta^2}{\gamma^2}$ we need a larger privacy budget since usually α is bigger than 2, very small value of Γ can lead the left side formula very large.

Methods	Email	URL	Phone Number	Perp.	Perp. -Sens	Methods	Email	URL	Phone Number	Perp.	Perp. -Sens
kNN-LM	0	13	0	2.872	2.872	RAG	36	56	125	1.58	1.58
Ours	0	0	0	2.89	2.20	Ours	0	3	0	2.05	1.06
DP-RP-G	0	2	0	2.96	2.89	DP-RP-G	0	3	0	2.05	1.59

Table 1: 1	KNN-LM
------------	--------

Table 2: Direct-prompting RAG

Apart from k and normalization bounds, which are directly related to privacy budgets, another parameter can implicitly influence the tradeoff between resistance to extraction attacks and language generation performance. From Theorem 1, we observe that σ is independent of the privacy budget. Setting σ too high increases the input magnitude to the language model, potentially leading to inaccurate predictions. Conversely, if σ is too small, the robustness of LLMs to minor noise may prevent effective replacement of sensitive words, as a small perturbation in embeddings can still preserve the original semantic features from the model's perspective.

5 EVALUATION

5.1 METRICS AND SETUP

For content quality, we follow prior work (Huang et al., 2023), measuring perplexity. Lower perplexity indicates higher content quality. We also measure the perplexity where sensitive information is removed from the label. We call it perplexity-sens. The reason is that if the privacy-preserving solution works well, the generated content will not have sensitive information. Therefore, sensitive information should not be counted on the label for a fair comparison.

To evaluate defense against attacks, we conduct empirical extraction attacks (Huang et al., 2023) on the RAG model using specific prompts to extract details (e.g., emails, websites, phone numbers). We use the Enron Email dataset (Klimt & Yang, 2004), which does not overlap with those used to pre-train major LLMs like GPT-2 (Radford et al., 2019), ensuring the RAG data are unseen by the pre-trained model. Thus, private data is limited to RAG documents. We used GPT-2 (Radford et al., 2019) as the pre-trained LLM. In KNN-LM, λ in KNN-LM is 0.1 and K in KNN is 1024. We set $\gamma = 1$, k = 64, $\Delta = 2$, and $\sigma = 0.1$. For Enron email dataset, the datastore size n is 465026. For GPT-2, the model embedding size is 768.

5.2 Results

To verify our algorithm's effectiveness, we compare it with baselines of no privacy-preserving methods and directly applying DP random projections under different privacy budgets. Apart from performance and generation, we measure the latency for one RAG inference. We set a tight privacy bound as $\epsilon = 5.051$ and $\alpha = 99$. With the same setting, using GPT-2 without any RAG methods or datastore yields a perplexity of 3.31, which is far from satisfactory.

5.2.1 KNN-LM

The results show that even under a tight privacy budget, generation performance is largely maintained with minimal loss. Interestingly, performance improves when sensitive information is excluded, as random projection masks overly detailed data, allowing kNN searches to focus on meaningful content rather than specifics like phone numbers. This leads to more contextually appropriate word selection and better outcomes. Directly applying random projection (Li & Li, 2023) over query and datastore embeddings requires privacy preservation at the cost of performance degradation, with two items still leaked. Averaging over 50,000 inferences, each inference takes NN-LM 0.793s and our method 0.811s on Nvidia RTX 4090.

5.2.2 DIRECT-PROMPTING LM

Apart from KNN-LM, we apply our methods to modern RAG architecture. We first use KNN-LM to retrieve the text. We then put the retrieved content along with the query in a message template and input them into a pre-trained LLM for answers, following the paradigm of the prevalent RAG.

Compared to kNN-LM, modern RAG achieves better generation performance by combining retrieval results with queries for more aggressive generation. However, this increases privacy leakage, with over 100 personal phone numbers retrievable. Fortunately, our methods prevent such leaks while maintaining similar results and improving perplexity-sens. Perplexity worsens slightly as sensitive information is excluded from the generated results. In summary, our approach successfully defends against attacks without sacrificing generation performance. Each interface takes RAG 2.996s and our applied methods 2.975s.

6 DISCUSSION AND FUTURE WORK

This paper presents preliminary results through empirical study and theoretical analysis to show that randomly projecting RAG document embedding onto a lower space is differentially private and can maintain the performance of RAG. In future work, we aim to have a further study to verify how such an algorithm performs with advanced and more commonly used architectures and language tasks. For example, instead of GPT-2, modern LLMs such as Llama 3 Touvron et al. (2023), Gemini Team et al. (2023), and DeepSeek R1 Guo et al. (2025) should be taken into consideration. Datasets more commonly used in evaluating RAG tasks also need to be considered. Apart from that, more various RAG patterns such as GraphRAG Edge et al. (2024) also worth studying.

To further study the tradeoff between privacy and generation performance, an important direction is how we can choose the best privacy budget in recording of k and normalization bounds to achieve the best performance. Apart from that, the selection of σ can impact the performance as well. Hence, a future direction will be an investigation of the automatic way of selecting these parameters.

REFERENCES

- Jeremiah Blocki, Avrim Blum, Anupam Datta, and Or Sheffet. The johnson-lindenstrauss transform itself preserves differential privacy. In *Proceedings of the 2012 IEEE 53rd Annual Symposium on Foundations of Computer Science (FOCS)*, pp. 410–419. IEEE Computer Society, 2012.
- Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Proceedings of the Third Theory of Cryptography Conference, TCC 2006, New York, NY, USA, March 4-7, 2006.*, pp. 265–284. Springer, 2006.
- Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, and Jonathan Larson. From local to global: A graph rag approach to query-focused summarization. arXiv preprint arXiv:2404.16130, 2024.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-R1: Incentivizing reasoning capability in Ilms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Yangsibo Huang, Samyak Gupta, Zexuan Zhong, Kai Li, and Danqi Chen. Privacy implications of retrieval-based language models. In *Proceedings of the 2023 Conference on Empirical Methods* in Natural Language Processing (EMNLP), pp. 14887–14902, 2023.
- Alaa Mahmoud Ibrahim, Mohamed Farouk, and Mohamed Waleed Fakhr. Privacy preserving image retrieval using multi-key random projection encryption and machine learning decryption. *Journal of Advanced Research in Applied Sciences and Engineering Technology*, 42(2):155–174, 2024.
- Changyue Jiang, Xudong Pan, Geng Hong, Chenfu Bao, and Min Yang. RAG-thief: Scalable extraction of private data from retrieval-augmented generation applications with agent-based attacks. *arXiv preprint arXiv:2411.14110*, 2024.
- William B Johnson, Joram Lindenstrauss, and Gideon Schechtman. Extensions of lipschitz maps into banach spaces. *Israel Journal of Mathematics*, 54(2):129–138, 1986.

- Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. Generalization through memorization: Nearest neighbor language models. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.
- Bryan Klimt and Yiming Yang. The enron corpus: A new dataset for email classification research. In *Proceedings of the European Conference on Machine Learning (ECML)*, pp. 217–226. Springer, 2004.
- Tatsuki Koga, Ruihan Wu, and Kamalika Chaudhuri. Privacy-preserving retrieval augmented generation with differential privacy. arXiv preprint arXiv:2412.04697, 2024.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Proceedings of the Advances in Neural Information Processing Systems (NeuraIPS)*, volume 33, pp. 9459–9474, 2020.
- Ping Li and Xiaoyun Li. Differential privacy with random projections and sign random projections. arXiv preprint arXiv:2306.01751, 2023.
- Sewon Min, Weijia Shi, Mike Lewis, Xilun Chen, Wen-tau Yih, Hannaneh Hajishirzi, and Luke Zettlemoyer. Nonparametric masked language modeling. *Proceedings of the Findings of the Association for Computational Linguistics (ACl Findings)*, 2023.
- Ilya Mironov. Rényi differential privacy. In *Proceedings of the 2017 IEEE 30th computer security foundations symposium (CSF)*, pp. 263–275. IEEE, 2017.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- Shenglai Zeng, Jiankun Zhang, Pengfei He, Jie Ren, Tianqi Zheng, Hanqing Lu, Han Xu, Hui Liu, Yue Xing, and Jiliang Tang. Mitigating the privacy issues in retrieval-augmented generation (rag) via pure synthetic data. *arXiv preprint arXiv:2406.14773*, 2024a.
- Shenglai Zeng, Jiankun Zhang, Pengfei He, Yue Xing, Yiding Liu, Han Xu, Jie Ren, Shuaiqiang Wang, Dawei Yin, Yi Chang, et al. The good and the bad: Exploring privacy issues in retrieval-augmented generation (RAG). *arXiv preprint arXiv:2402.16893*, 2024b.

A COMPLETE PROOF OF THEOREM 1

Proof. We consider any adjacent datasets $D_1, D_2 \in \mathcal{D}$ where they differ by one record (row) and assume $D_1, D_2 \in \mathbb{R}^{n \times d}$. We normalize each token vector (row) in D_1 and D_2 to Δ as the upper bound and δ as the lower bound ($\gamma \leq ||D_{1_j}|| \leq \Delta$), for each embedding in D_1 and D_2 . Next, we use Johnson Lindenstrauss Leema to calculate that when we want to achieve a (ϵ, α) -DP guarantee, we need to project a vector of dimension d to dimension k. We then let R be an IID Gaussian Matrix in the dimension of $d \times k$ whose entries are IID samples from $\mathcal{N}(0, \sigma^2)$. Next, we project c_i to $c_i R$. After the random projection, $P = \Pr[\arg \max p(y|x, D_1 R) \in S]$.

We note that the D_1R and D_2R distributions follow normal distributions. Hence, we are going to calculate the probability of matrix D_1R equaling to Z. Except the *m*th row with the difference, the other rows in D_1R and D_2R are the same. So, the probabilities that they are the same as the vector on the corresponding row in Z are the same. We only now to consider the probability of sampling $[D_1R]_m$ and $[D_2R]_m$.

We first calculate the expectations of $[D_1R]_{mj}$ and $[D_2R]_{mj}$. Because $\mathbb{E}[[D_1R]_{mj}] = \mathbb{E}[\Sigma_l D_{1_{ml}} R_{lj}] = \Sigma_l(\mathbb{E}[D_{1_{ml}} R_{lj}]) = \Sigma_l(\mathbb{E}[D_{1_{ml}}]\mathbb{E}[R_{lj}]) = 0$. The next step is to calculate the variance of D_1R_m and D_2R_m . We assume that the variance are σ_1 and σ_2 . We know that R is an i.i.d Gaussian matrix. Hence, we can view $D_{1_m}R$ as a random Gaussian vector in the dimension of k. For each element in $D_{1_m}R$, the variance is σ_1 . We pick a random element $[D_1R]_{mj}$ and can have $Var([D_1R]_{mj}) = Var(\Sigma_{l=1}^d D_{1_{ml}}R_{lj} = \Sigma_{l=1}^d D_{1_{ml}}^2 Var(R_{lj}) = \sigma_1^2 ||D_{1_m}||_2^2$

Back to our target, to prove privacy guarantees, we need to have that

$$D_{\alpha}(\Pr_{D_1R_m}(Z_m) \| \Pr_{D_2R_m}(Z_m)) \le \epsilon$$
⁽²⁾

As Z_m should follow the multi-variant Gaussian distribution, and according to the definition of Ren-yi divergence over the continuous random variable, we can have:

$$\frac{1}{\alpha - 1} \ln \int \Pr_{D_1 R_m}(Z_m)^{\alpha} \Pr_{D_2 R_m}(Z_m))^{1 - \alpha} \le \epsilon$$
(3)

We also have that

$$\Pr_{D_1 R_m}(Z_m) = \frac{1}{(2\pi)^{k/2} |\Sigma_1|^{1/2}} e^{-\frac{1}{2} Z_m^T \Sigma_1^{-1} Z_m}$$
(4)

$$\Pr_{D_2 R_m}(Z_m) = \frac{1}{(2\pi)^{k/2} |\Sigma_2|^{1/2}} e^{-\frac{1}{2} Z_m^T \Sigma_2^{-1} Z_m}$$
(5)

We know that $\Sigma_1 = \sigma_1^2 I_k$ and $\Sigma_2 = \sigma_2^2 I_k$. So, we have

$$\Pr_{D_1 R_m}(Z_m) = \frac{1}{(2\pi)^{k/2} \sigma_1^k} e^{-\frac{\|Z_m\|_2^2}{2\sigma_1^2}}$$
(6)

$$\Pr_{D_2 R_m}(Z_m) = \frac{1}{(2\pi)^{k/2} \sigma_2^k} e^{-\frac{\|Z_m\|_2^2}{2\sigma_2^2}}$$
(7)

Hence, LHS in Eq. (3) becomes

$$\frac{1}{\alpha-1}\ln\left[\frac{1}{(2\pi)^{k/2}\sigma_1^{\alpha k}\sigma_2^{(1-\alpha)k}}\int_{\mathbb{R}^k}e^{-\frac{1}{2}\left(\frac{\alpha}{\sigma_1^2}+\frac{1-\alpha}{\sigma_2^2}\right)Z_m^T I_k Z_m}d_{Z_m}\right] \tag{8}$$

$$= \frac{1}{\alpha - 1} \ln \left[\frac{1}{(2\pi)^{k/2} \sigma_1^{\alpha k} \sigma_2^{(1-\alpha)k}} (2\pi)^{k/2} \left(\frac{\alpha}{\sigma_1^2} + \frac{1-\alpha}{\sigma_2^2} \right)^{-k/2} \right]$$
(9)

$$= \frac{1}{\alpha - 1} \ln \left[\sigma_1^{-\alpha k} \sigma_2^{-(1-\alpha)k} \left(\frac{\sigma_1^2 \sigma_2^2}{\alpha \sigma_2^2 + (1-\alpha)\sigma_1^2} \right)^{k/2} \right]$$
(10)

$$= \frac{k}{2(\alpha - 1)} \ln \frac{(\sigma_1^2)^{1 - \alpha} (\sigma_2^2)^{\alpha}}{(1 - \alpha)\sigma_1^2 + \alpha\sigma_2^2}$$
(11)

$$=\frac{k}{2(\alpha-1)}\ln\left[\frac{\left(\frac{\sigma_1^2}{\sigma_2^2}\right)^{1-\alpha}}{\left(1-\alpha\right)\left(\frac{\sigma_1^2}{\sigma_2^2}\right)^2+\alpha}\right]$$
(12)

Now, we let $\frac{\sigma_1^2}{\sigma_2^2} = \frac{\sigma^2 \|D_{1_m}\|_2^2}{\sigma^2 \|D_{2_m}\|_2^2} = \Gamma.$

As a result, we need to pick Δ and γ to satisfy the following equation:

$$\frac{\Gamma^{1-\alpha}}{(1-\alpha)\Gamma+\alpha} < e^{\frac{2(\alpha-1)\epsilon}{k}}$$
(13)

For this equation, we can see that, in the case where $\Gamma \ge 1$, the left equation is always smaller than 1. In this case, the in-equation is always satisfied. Hence, we just need to choose proper $\frac{\gamma^2}{\Delta^2}$ so that for any values of $\Gamma \ge \frac{\gamma^2}{\Delta^2}$, Eq. (13) is satisified.

B FROM RDP TO DP

Further, we here show how we can convert the (α, ϵ) RDP guarantee to (ϵ, δ) -DP guarantee.

Definition 2 (Differential Privacy). A randomized mechanism $\mathcal{M} : \mathcal{D} \to \mathcal{S}$ satisfies (ϵ, δ) differential privacy if for any neighboring datasets D_1 and D_2 differing by one record, it holds that

$$\Pr[\mathcal{M}(D_1) \in \mathcal{S}] \le e^{\epsilon} \Pr[\mathcal{M}(D_2) \in \mathcal{S}] + \delta.$$

According to the proposition 3 in Mironov (2017), if \mathcal{M} is a (α, ϵ) -RDP mechanism, it also satisfies $(\epsilon + \frac{\ln 1/\delta}{\alpha - 1}, \delta)$ -differential privacy for any $0 < \delta < 1$. We let $\epsilon_{RDP} = \epsilon_{DP} - \frac{\ln 1/\delta}{\alpha - 1}$. In (ϵ, δ) -DP, we need to satisfy that

$$\frac{\Gamma^{1-\alpha}}{(1-\alpha)\Gamma+\alpha} < e^{\frac{2(\alpha-1)(\epsilon_{DP} - \frac{\ln 1/\delta}{\alpha-1})}{k}}$$
(14)