
Evaluating Privacy Risks in Synthetic Clinical Text Generation in Spanish

Luis Miranda

Pontificia Universidad Católica de Chile
Av. Vicuña Mackenna 4860, Santiago de Chile
lmirandn@uc.cl

Jocelyn Dunstan

Pontificia Universidad Católica de Chile
Av. Vicuña Mackenna 4860, Santiago de Chile
jdunstan@uc.cl

Matias Toro

Universidad de Chile
Beaucheff 851, Santiago de Chile
mtoro@dcc.uchile.cl

Federico Olmedo

Universidad de Chile
Beaucheff 851, Santiago de Chile
federico.olmedo@dcc.uchile.cl

Félix Melo

Universidad de Chile
Beaucheff 851, Santiago de Chile
felix.melo@ug.uchile.cl

Abstract

1 Leveraging medical data for Deep Learning models holds great potential, but
2 ensuring the protection of sensitive patient information is paramount in the clinical
3 domain. A widely used approach to balance data utility and privacy is the generation
4 of synthetic text with Large Language Models (LLMs) under the framework of
5 differential privacy (DP). Techniques like Differentially Private Stochastic Gradient
6 Descent (DP-SGD) are typically considered to provide privacy guarantees, but
7 they rely on specific conditions. This research demonstrates how memorization in
8 LLMs can deteriorate when these privacy safeguards are not fully met, increasing
9 the risk of personal and sensitive information being leaked in synthetic clinical
10 reports. Addressing these vulnerabilities could enhance the reliability of DP in
11 protecting clinical text data while maintaining its utility.

12 1 Introduction

13 The utilization of Electronic Health Records (EHRs) for Natural Language Processing (NLP) offers
14 numerous benefits, particularly in enhancing healthcare research and outcomes [7]. However, pro-
15 tecting the privacy of the patients in these records is crucial. Privacy is recognized as a core human
16 right in the Universal Declaration of Human Rights, placing the control individuals have over their
17 personal information on par with the authority exercised by corporations and governments [14].

18 According to the 2021 Annual Report of the United Nations High Commissioner, privacy reflects
19 human dignity and plays a critical role in safeguarding individual autonomy and identity. In today's
20 digital age, privacy concerns are even more pronounced as personal data—often considered a valuable
21 commodity—can be collected, sold, and potentially misused. This is particularly concerning when
22 sensitive health data is involved (e.g., apps that collect reproductive information, or dating apps that
23 ask for HIV status) [6]. The mishandling of such data not only threatens privacy but can also foster
24 discrimination and erode human dignity.

25 There are several techniques to protect patient privacy in EHRs, such as Named Entity Recognition
26 (NER) for de-identification or pseudo-anonymization [3, 17, 16]. However, synthetic text genera-
27 tion with Differential Privacy (DP) is often preferred due to its formal privacy guarantees and its
28 widespread use [20, 10, 19, 2].

29 Synthetic text refers to artificially generated text that mimics human language and content. One way
30 to create it is by using Large Language Models (LLMs), which generate text through "next-token
31 prediction." This process involves predicting the next word in a sentence based on the previous ones,
32 allowing the model to generate coherent text. In this context, the goal is to create realistic synthetic
33 Electronic Health Records (EHRs) that are similar to original EHRs, making them useful for research
34 and other purposes. To achieve this, an LLM can be trained using real EHR data.

35 Training an LLM involves exposing the model to a dataset and adjusting its parameters based on the
36 patterns it learns. However, during this process, the model might memorize personal information and
37 reproduce it [4], which is critical when dealing with clinical data. To prevent this, DP can be applied.
38 DP, in essence, ensures that individual data points within a dataset do not significantly influence the
39 outcome of an algorithm, protecting information quantified by a level of privacy ϵ [9]. A common
40 technique used for training an LLM with DP is Differentially Private Stochastic Gradient Descent
41 (DP-SGD), which adds noise during training to prevent memorization, ensuring both privacy and
42 utility [1, 11].

43 However, the mere use of DP-SGD often leads to an assumption of privacy guarantees, but in practice,
44 is frequently overlooked. DP-SGD provides "sample-level" privacy [18, 11], meaning it protects
45 individual data points as long as the same individual does not appear in multiple samples. In clinical
46 datasets, this assumption is unfeasible, as the same individual may be represented in multiple samples.
47 This raises serious concerns about the true effectiveness of DP in such contexts.

48 To address potential privacy concerns, it is important to evaluate privacy beyond standard guarantees,
49 such as by assessing the level of memorization. Previous research has primarily focused on measuring
50 model memorization and the leakage of sensitive information in synthetic data, particularly the
51 leakage of isolated pieces of Personally Identifiable Information (PII) [20, 5]. Building on these
52 studies, this work introduces a novel method for analyzing the memorization of LLMs and the risk
53 of information leakage in synthetic EHRs generated in Spanish. This presents unique challenges
54 specific to the language (e.g. the more frequent use of gendered terms throughout sentences).

55 2 Experimental Setup

56 In this study we used the MEDDOCAN dataset [12] (available here), which consists of 1,000
57 manually crafted Spanish clinical reports enriched with personal information and annotated with
58 NER for PII and sensitive data. For computing limitations, the final dataset used consisted of 1000
59 reports, divided into 750 documents for training and 250 for validation. These documents are used to
60 analyze information leakage at the document level. We conducted the experiments using the LLM
61 meta-llama/Meta-Llama-3.1-8B-Instruct [8].

62 The training was executed with the following parameters: 7 epochs, batch size of 2, gradient
63 accumulation steps of 1, LoRA dimension of 4, LoRA alpha, and a learning rate of $3e-4$. Additionally,
64 the training of the LLMs was performed on 2 NVIDIA RTX 6000 Ada Generation GPUs.

65 3 Methodology

66 The training used DP-SGD, which adds noise to gradients during the training process to safeguard the
67 original data's privacy [1]. We trained the models using identical parameters across different dataset
68 versions, each with varying levels of differential privacy. ϵ , a key parameter in differential privacy,
69 measures privacy loss, with lower values providing stronger protection. The used values are $\epsilon = 8$,
70 16, and ∞ (no privacy).

71 After training, 500 synthetic documents were generated with each model. These documents were
72 analyzed to assess memorization and evaluate the quality and utility of the generated text. The
73 generation process was standardized putting the same training parameters to ensure comparable
74 results across models. Finally, we applied various metrics to examine the privacy-utility trade-off and
75 the extent of memorization.

76 **3.1 Utility Metrics**

77 The utility of the synthetic documents generated by each model was evaluated using key metrics
 78 such as MAUVE and perplexity (PPL). MAUVE [15] measures the quality and diversity of generated
 79 text using divergence frontiers, reflecting how closely the synthetic data aligns with the distribution
 80 of real text. PPL assesses how well a model predicts a sample, with lower values indicating better
 81 performance [13]. These metrics were used to evaluate the impact of differential privacy on the
 82 quality and coherence of the generated EHRs.

83 **3.2 Leakage of Sensitive Information**

84 To evaluate the impact of synthetic text generation with DP-SGD when private patient information
 85 is repeated across documents, we adapted the “canary” experiment [5]. This involved injecting a
 86 “canary” sentence containing a single piece of PII repeated across documents, allowing us to track
 87 how often it appeared in generated samples. In our version, two pieces of information—a reference to
 88 positive HIV as sensitive data and the name “Lopez Perez” to link it to personal information—were
 89 embedded into 0, 50, and 200 documents. We then counted how often this information appeared in
 90 the generated samples. In this way, we assess the memorization of links between sensitive data and
 91 individuals rather than the memorization of individual data points, which is crucial in the context
 92 of sensitive clinical data, as the ability to link sensitive information (such as an illness or medical
 93 history) to an individual must be protected.

94 **4 Results and Discussion**

Inj. Can.	ϵ	MAUVE		PPL		Leaked Can.	
		Model 1	Model 2	Model 1	Model 2	Model 1	Model 2
0	8	0.48	0.84	7.84±0.42	8.27±0.43	0	0
0	16	0.55	0.88	7.56±0.21	8.44±0.39	0	0
0	∞	0.83	0.89	6.02±0.29	4.73±0.24	0	0
50	8	0.47	0.76	7.76±0.44	8.34±0.37	0	1
50	16	0.59	0.80	7.57±0.23	8.76±0.32	2	2
50	∞	0.82	0.87	6.06±0.27	4.39±0.07	76	120
200	8	0.41	0.81	8.05±0.35	8.32±0.39	1	1
200	16	0.55	0.85	7.75±0.30	8.45±0.19	3	8
200	∞	0.84	0.95	5.72±0.32	4.91±0.64	103	331

Table 1: Privacy-utility evaluation results for Model 1 : mistralai/Mistral-7B-v0.1 and Model 2 : meta-llama/Meta-Llama-3.1-8B-Instruct. The models were evaluated across varying privacy levels ($\epsilon = 8, 16, \infty$) and different quantities of injected canaries (Inj. Can.). The evaluation metrics include MAUVE, Perplexity (PPL), and the number of leaked canaries (Leaked Can.) in the 500 synthetic generated data.

95 Table 1 shows the results of synthetically generated texts evaluated by models trained with different
 96 privacy levels ($\epsilon = 8, 16, \infty$) and varying numbers of injected canaries (0, 50, 200). The utility
 97 metrics, MAUVE and PPL, reveal that as privacy increases (lower ϵ), MAUVE decreases and PPL
 98 rises, indicating lower text quality and diversity due to the added noise from DP-SGD. Additionally,
 99 Model 1 displays lower PPL but also a lower MAUVE than Model 2, suggesting that while the
 100 text generated by Model 1 is more predictable, it is less natural and diverse—consistent with the
 101 definitions of MAUVE and PPL. Except in the case where there is no privacy ($\epsilon = \infty$), where Model
 102 1 shows both lower MAUVE and higher PPL than Model 2.

103 Regarding canary leakage, the more frequently a canary (e.g., name and disease) is injected into
 104 the training data, the more it appears in the generated texts, with over 15% of the text containing
 105 personal information in some cases. However, when differential privacy is applied, this percentage
 106 drops to less than 2%. Despite this reduction, conditions for privacy guarantees are still violated, as
 107 differential privacy requires that no individual appear in more than one sample. Consequently, the
 108 generated text would be leaking that the individual with the surname “Lopez Perez” is HIV positive.

109 5 Conclusions and Limitations

110 While DP-SGD is widely believed to provide strong privacy guarantees, our findings reveal that
111 memorization in LLMs occurs when those privacy guarantees are compromised, particularly in
112 cases where the same individual appears across multiple samples—an aspect rarely considered
113 when applying these methods. This was done by injecting the same linked personal and sensitive
114 information multiple times in the training data of an LLM and then quantifying the leakage of
115 this information in synthetic generated data by the model, offering a more comprehensive view
116 of information leakage across entire documents, rather than focusing on individual PII entities.
117 This raises concerns about the effectiveness of DP in clinical datasets, where privacy protection is
118 paramount. Despite these challenges, DP can still serve as a valuable tool for safeguarding individuals
119 if its conditions are properly fulfilled.

120 It is important to highlight a limitation of this work: while the goal is to analyze the level of
121 memorization in an entire clinical report, this study only focuses on whether a name paired with a
122 disease appears in any generated text. Although this approach is closer to identifying more than just
123 a name, it still falls short of capturing a complete clinical record. Clinical records (for the dataset
124 used in this study) typically include more detailed information such as medical history, addresses,
125 phone numbers, prescriptions, medications, diagnoses, and more. Therefore, a more robust method
126 for analyzing memorization in clinical reports that takes into account all this additional information
127 is needed.

128 As future work, we propose employing feature extraction and Named Entity Recognition (NER)
129 algorithms to detect personal and sensitive information in each synthetically generated document.
130 Once extracted, this information can be used to analyze memorization across different differentially
131 private algorithms for generating synthetic clinical data, comparing how these techniques perform in
132 both the original and synthetic texts.

133 References

- 134 [1] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang. Deep
135 learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on*
136 *Computer and Communications Security, CCS '16*, page 308–318, New York, NY, USA, 2016.
137 Association for Computing Machinery.
- 138 [2] N. C. Abay, Y. Zhou, M. Kantarcioglu, B. Thuraisingham, and L. Sweeney. Privacy preserving
139 synthetic data release using deep learning. In *Machine Learning and Knowledge Discovery*
140 *in Databases: European Conference, ECML PKDD 2018, Dublin, Ireland, September 10–14,*
141 *2018, Proceedings, Part I 18*, pages 510–526. Springer, 2019.
- 142 [3] C. Aracena, L. Miranda, T. Vakili, F. Villena, T. Quiroga, F. Núñez-Torres, V. Rocco, and
143 J. Dunstan. A privacy-preserving corpus for occupational health in Spanish: Evaluation for
144 NER and classification tasks. In T. Naumann, A. Ben Abacha, S. Bethard, K. Roberts, and
145 D. Bitterman, editors, *Proceedings of the 6th Clinical Natural Language Processing Workshop*,
146 pages 111–121, Mexico City, Mexico, June 2024. Association for Computational Linguistics.
- 147 [4] E. M. Bender, T. Gebru, A. McMillan-Major, and S. Shmitchell. On the dangers of stochastic
148 parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on*
149 *Fairness, Accountability, and Transparency, FAccT '21*, page 610–623, New York, NY, USA,
150 2021. Association for Computing Machinery.
- 151 [5] N. Carlini, C. Liu, Ú. Erlingsson, J. Kos, and D. Song. The secret sharer: Evaluating and testing
152 unintended memorization in neural networks. In *28th USENIX Security Symposium (USENIX*
153 *Security 19)*, pages 267–284, Santa Clara, CA, Aug. 2019. USENIX Association.
- 154 [6] D. Citron. *The Fight for Privacy: Protecting Dignity, Identity and Love in the Digital Age*.
155 Random House, 2022.
- 156 [7] H. Dalianis. *Clinical Text Mining: Secondary Use of Electronic Patient Records*. Springer
157 International Publishing, 2018.
- 158 [8] A. Dubey, A. Jauhri, A. Pandey, A. Kadian, and (...). The llama 3 herd of models, 2024.

- 159 [9] C. Dwork. Differential privacy. In M. Bugliesi, B. Preneel, V. Sassone, and I. Wegener, editors,
160 *Automata, Languages and Programming*, pages 1–12, Berlin, Heidelberg, 2006. Springer Berlin
161 Heidelberg.
- 162 [10] J. Flemings and M. Annavam. Differentially private knowledge distillation via synthetic text
163 generation. In L.-W. Ku, A. Martins, and V. Srikumar, editors, *Findings of the Association
164 for Computational Linguistics ACL 2024*, pages 12957–12968, Bangkok, Thailand and virtual
165 meeting, Aug. 2024. Association for Computational Linguistics.
- 166 [11] O. Klymenko, S. Meisenbacher, and F. Matthes. Differential privacy in natural language process-
167 ing the story so far. In O. Feyisetan, S. Ghanavati, P. Thaine, I. Habernal, and F. Mireshghallah,
168 editors, *Proceedings of the Fourth Workshop on Privacy in Natural Language Processing*, pages
169 1–11, Seattle, United States, July 2022. Association for Computational Linguistics.
- 170 [12] M. Marimon, A. Gonzalez-Agirre, A. Intxaurreondo, H. Rodriguez, J. L. Martin, M. Villegas,
171 and M. Krallinger. Automatic De-identification of Medical Texts in Spanish: the MEDDOCAN
172 Track, Corpus, Guidelines, Methods and Evaluation of Results. In *IberLEF@ SEPLN*, pages
173 618–638, 2019.
- 174 [13] A. Miaschi, D. Brunato, F. Dell’Orletta, and G. Venturi. What makes my model perplexed? a
175 linguistic investigation on neural language models perplexity. In E. Agirre, M. Apidianaki, and
176 I. Vulić, editors, *Proceedings of Deep Learning Inside Out (DeeLIO): The 2nd Workshop on
177 Knowledge Extraction and Integration for Deep Learning Architectures*, pages 40–47, Online,
178 June 2021. Association for Computational Linguistics.
- 179 [14] Z. Nampewo, J. H. Mike, and J. Wolff. Respecting, protecting and fulfilling the human right to
180 health. *International Journal for Equity in Health*, 21(1), Mar. 2022.
- 181 [15] K. Pillutla, S. Swayamdipta, R. Zellers, J. Thickstun, S. Welleck, Y. Choi, and Z. Harchaoui.
182 Mauve: Measuring the gap between neural text and human text using divergence frontiers. In
183 M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, editors, *Advances in
184 Neural Information Processing Systems*, volume 34, pages 4816–4828. Curran Associates, Inc.,
185 2021.
- 186 [16] T. Vakili, A. Henriksson, and H. Dalianis. End-to-End Pseudonymization of Fine-Tuned Clinical
187 BERT Models, Sept. 2023.
- 188 [17] S. Verkijk and P. Vossen. Efficiently and Thoroughly Anonymizing a Transformer Language
189 Model for Dutch Electronic Health Records: a Two-Step Method. In *Proceedings of the
190 Thirteenth Language Resources and Evaluation Conference*, pages 1098–1103, Marseille,
191 France, June 2022. European Language Resources Association.
- 192 [18] Y. Wang, Q. Wang, L. Zhao, and C. Wang. Differential privacy in deep learning: Privacy and
193 beyond. *Future Generation Computer Systems*, 148:408–424, 2023.
- 194 [19] B. Xin, Y. Geng, T. Hu, S. Chen, W. Yang, S. Wang, and L. Huang. Federated synthetic data
195 generation with differential privacy. *Neurocomputing*, 468:1–10, 2022.
- 196 [20] X. Yue, H. Inan, X. Li, G. Kumar, J. McAnallen, H. Shajari, H. Sun, D. Levitan, and R. Sim.
197 Synthetic text generation with differential privacy: A simple and practical recipe. In A. Rogers,
198 J. Boyd-Graber, and N. Okazaki, editors, *Proceedings of the 61st Annual Meeting of the
199 Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1321–1342, Toronto,
200 Canada, July 2023. Association for Computational Linguistics.