Personal Information Parroting in Language Models

Nishant Subramani¹ Kshitish Ghate¹ Mona T. Diab¹

Abstract

Modern language models (LM) are trained on large scrapes of the Web, containing millions of personal information (PI) instances, many of which LMs memorize, increasing privacy risks. In this work, we develop the regexes and rules (R&R) detector suite to detect email addresses, phone numbers, and IP addresses, which outperforms the best regex-based PI detectors. On a manually curated set of 483 instances of PI, we measure memorization: finding that 13.6% are parroted verbatim by the Pythia-6.9b model, *i.e.*, when the model is prompted with the tokens that precede the PI in the original document, greedy decoding generates the entire PI span exactly. We expand this analysis to study models of varying sizes (160M-6.9B) and timesteps of pretraining (70k-143k iterations) on the Pythia model suite and find that both model size and amount of pretraining are positively correlated with memorization. Even the smallest model, Pythia-160m, parrots 2.7% of the instances exactly. Consequently, we strongly recommend that pretraining datasets be aggressively filtered and anonymized to minimize PI parroting. The code for our detectors can be found at https://github. com/nishantsubramani/rr_pi_detectors/.

1. Introduction

Large language models (LLMs) are trained on trillions of tokens scraped from the Web, containing millions of instances of personal information (PI) (Elazar et al., 2023; Subramani et al., 2023; Soldaini et al., 2024). We use the term PI because it encapsulates the US definition of personally identifiable information (PII), the UN definition of personal data, and other definitions in other countries (Subramani et al., 2023). For many pretraining datasets, however, documentation of PI is absent. Furthermore, LLMs memorize training examples and can be prodded to extract PI via prompt-based methods (Carlini et al., 2019; 2021; 2022). LLMs can also be steered to generate exact sequences primarily via steering vectors and prompting, unrelated to privacy (Subramani et al., 2019; 2022; Subramani & Suresh, 2020; Shin et al., 2020; Li & Liang, 2021). This indicates a serious problem: LLMs memorize and generate PI and a malicious actor can gain access to these without complex extraction attacks. Better filtering can improve PI memorization for both filtered examples and for examples that were not caught by the filter (Borkar et al., 2025). However, PI filtering and anonymization has largely been ignored when curating pretraining datasets; those that do tend to resort to regularexpression (regex) based approaches because model-based approaches are computationally infeasible (Subramani et al., 2023; Elazar et al., 2023; Soldaini et al., 2024).

To address these limitations, we focus on character-based PI, which are among the highest risk PI types (Subramani et al., 2023). Accordingly, our work contributes the following:

- We develop the regexes and rules detector suite (R&R) containing four new PI detectors for email addresses, IP addresses, US phone numbers, and US phone numbers with the country code and show that our suite outperforms the strongest regex-based PI detectors (Elazar et al., 2023).
- 2. Using the Pythia suite, we measure the degree of PI parroting and analyze the effect of model size, pretraining timesteps, and prefix length on memorization.

2. R&R Detection Suite

PI Detection & Baselines Following Subramani et al. (2023), we focus on character-based PI since they are one of the highest risk categories for risk exposure as they can often uniquely identify a person. Model-based tools like Presidio (Microsoft, 2021) outperform regular-expression based systems slightly, but are infeasibly slow on large datasets (*e.g.*, pretraining corpora). As a result, we compare with the WIMBD detectors (Elazar et al., 2023), the strongest efficient baselines to our knowledge, which contain detectors for three PI types: email addresses, US phone numbers, and IP addresses.

¹Carnegie Mellon University, LTI. Correspondence to: Nishant Subramani <nishant2@cs.cmu.edu>.

Published at MemFM Workshop at ICML 2025, Vancouver, Canada. Copyright 2025 by the author(s).



Figure 1: Results on manually annotating 1750 detections, 250 per PI type per system. R&R (blue) outperforms WIMBD (orange) in 17 of 20 cases. WIMBD cannot detect US phone numbers when it has the country code (+1), yielding a precision of 0%. R&R significantly outperforms WIMBD for email addresses and both phone numbers and is not significantly better for IP addresses.

Dataset We choose the Pile (Gao et al., 2020), one of the largest open-source pretraining datasets containing 383 billion tokens. The dataset consists of 22 smaller subsets including OpenWebText2, arXiv, StackExchange, and YoutubeSubtitles with text from many genres including research papers, patents, subtitles, law, science, and mathematics. These are grouped into 5 categories: academic, dialogue, internet, prose, and misc. The Pile was also used to train the Pythia model suite, a suite of auto-regressive decoder-only language models ranging from 70 million parameters to 12 billion parameters with 143 intermediate model checkpoints (Biderman et al., 2023). We use the Pythia model suite in our memorization experiments in §4.

2.1. Regexes and Rules Detectors (R&R)

We develop our own regex patterns, improving upon the patterns in WIMBD and add rules to increase precision across PI types. R&R features four detectors for email addresses, IP addresses, US/Canada phone numbers, and US/Canada

PI type	total detections	expected PI counts	
email addresses	16.389,977	12,623,478	
IP addresses	7,801,628	4,411,309	
phone numbers	1,275,862	278,332	
phone numbers (global)	172,326	28,987	

Table 1: Total detections and expected PI counts across the Pile dataset using R&R for each PI type.

phone numbers with the +1 country code respectively.

Regexes For *email addresses*, we allow for a wider range of characters in the username, including all alphanumeric characters, valid special characters, and periods. This modification helps us detect non-traditional domain names through literals. For *IP addresses*, we add an additional pattern to detect IPv6 addresses because WIMBD only considers IPv4. For *phone numbers*, we ensure that any detected phone number is not followed by a digit, whereby removing false positives resulting from numbers with more than 10 digits. Since WIMBD, fails to detect phone numbers with a country code, we develop an additional regex to target phone numbers with a "+1" country code prefix, which reflects the common format for US and Canadian phone numbers. See Appendix C for details on the specific regular expressions used.

Rules To complement the updated regexes, we introduce new post-processing rules, particularly addressing contextual subtleties that a regex cannot easily capture. For phone numbers and IP addresses, we add filtering rules to check whether part numbers, ISBN numbers, and similar identification numbers are in the context. Additionally, we remove placeholder examples such as 123-456-7890. For phone numbers, we add an area code validator and a central office code validator ensuring compliance with the North American Numbering Plan (NANP). See Appendix C for details on the specific rules we use. Taken together these updates improve detection efficacy for all PI types.

2.2. R&R vs. WIMBD

Figure 1 shows the results of our manual audit of both the WIMBD and R&R detection suites. ¹ We run both detectors on the entire Pile dataset and take a random sample of 1750 detections. To improve coverage, we leverage stratified sampling from the detectors across 5 different subcategories

¹We focus solely on precision. This mirrors the annotation process of prior work, where the authors manually annotated detections (Subramani et al., 2023; Elazar et al., 2023). Annotating pretraining documents would be infeasible without a large pool of annotators, however using that pool would reveal PI publicly, drastically increasing privacy risks.

(academic, dialogue, internet, prose, and misc) of the Pile. We sample and annotate the data selected across all PI types. Overall, we find that R&R has a total of 483 true positives, with 99% of them having a perfect span. We use this gold set to quantify memorization.

For all four PI types (including phone numbers with +1) and for 17 of the 20 categories in Figure 1, R&R outperforms WIMBD. The improvement on phone numbers is especially notable: WIMBD has a precision of nearly 0, whereas R&R has a precision of 0.3 on average. Using both the total detected counts and the precision values calculated from the manual annotation, we compute the expected PI counts across the Pile dataset. Both total counts and expected counts are shown in Table 1. Email addresses and IP addresses are orders of magnitude more prevalent than phone numbers in the Pile. We run permutation tests with 10,000 resamples to test whether R&R is significantly better than WIMBD across all four types of PI. We find that R&R is significantly better for email addresses and both sets of phone numbers (p-value < 0.05), but not for IP addresses.

3. Quantifying Memorization and Risk

Measuring Memorization and Associated Model Parroting of Personal Information To quantify memorization and its associated model parroting, we use the definition of *p*-memorization (Carlini et al., 2022). A string *s* is said to be *p*-memorized if a model \mathcal{M} , when prompted with a string *s'* of length *p* generates *s* verbatim with greedy decoding and the concatenation, [s||s'], is in the training data of \mathcal{M} . This definition gives us a framework to quantify the extent to which a model has memorized a training instance in a deterministic fashion.²

Metrics Since character-based PI instances are strings, we use Levensthein distance between a candidate instance of PI and its ground-truth to compute a similarity score, which we call PARROTSCORE. More formally for a candidate string s_1 and a ground-truth string s_2 :

$$PARROTSCORE(s_1, s_2) = 1 - \frac{d_{levensthein}(s_1, s_2)}{|s_2|} \quad (1)$$

In practice, if $|s_1| > |s_2|$, we take every substring of s_1 of size $|s_2|$ and choose the one that has the maximum PAR-ROTSCORE with s_2 . Since PARROTSCORE $\in [0, 1]$, a score of 1 signifies verbatim parroting, while a score of 0 indicates that there is not a single character that overlaps between the two strings. A low score such as 0.1 can still pose privacy risks; just three characters at the end of an email address can



Figure 2: Here, we show the results of prompting each Pythia model with the prefix that occurs before each instance of PI and measuring the PARROTSCORE between the ground-truth PI and the model's greedily decoded generation across all PI types.

expose geographic information like the country a person lives in (*e.g.*, .nl).

4. Experiments

To measure PI parroting and memorization, we use the manually annotated and curated set of detections from our R&R detection suite which contains 483 instances of PI. We experiment with 6 models from the Pythia suite with 160m, 410m, 1b, 1.4b, 2.8b, 6.9b parameters respectively. For each instance of PI, we find its associated prefix in the Pile and truncate this to a maximum of 80 tokens. Using this potentially truncated prefix as a prompt, we generate from the LM using greedy decoding and evaluate whether it parrots the ground truth PI instance using PAR-ROTSCORE in equation (1). ³

5. Results & Analysis

Model size vs. memorization: Figure 2 shows how different models with varying model sizes parrot PI across all of the PI types considered. Email addresses are the most parroted, with an average PARROTSCORE of greater than 0.3 for all models. This is closely followed by IP addresses and then by phone numbers. Model size and PARROTSCORE are positively correlated, but even the smallest models have high PARROTSCORE indicating that even small models are prone to memorization and verbatim parroting. In fact, the 410m parameter model, the second smallest model we

²A model \mathcal{M} verbatim parroting an instance when prompted with a prefix of size p is similar to saying it has been p-memorized, if s' is the ground-truth PI.

³This is similar to measuring *p*-memorization for p = 80.

model sizes	email	ip	phone num	phone num +1
160m	2.34%	4.72%	1.64%	1.23%
410m	8.41%	7.09%	4.92%	2.47%
1b	12.62%	7.87%	3.28%	1.23%
1.4b	16.36%	11.81%	3.28%	1.23%
2.8b	19.63%	14.17%	1.64%	1.23%
6.9b	19.63%	14.17%	3.28%	4.94%

Table 2: Percent of total instances that are verbatim parroted. Note that this corresponds to the percent of instances that achieve a PARROTSCORE of 1. We find that verbatim parroting increases with model size and email addresses are the most likely to be parroted exactly.

tested, has a similar PARROTSCORE to the 2.8b model for both phone number types and has only a 0.1 lower PAR-ROTSCORE on email and IP addresses.

Model size vs. verbatim parroting: Table 2 presents the percent of total instances that are verbatim parroted by each model. Verbatim parroting and model size are positively correlated, email addresses and IP addresses contribute mostly to this trend. Both PI types are increasingly parroted: nearly 20% of all detected email addresses and over 14% of IP addresses are exactly parroted by the two largest models. Phone numbers have a much lower verbatim parrot rate, which is uncorrelated with model size, indicating that phone numbers can be challenging for LMs to memorize.

Pretraining steps vs. memorization: The top plot of Figure 3 shows that even from only half of the pretraining (70,000 steps) to fully pretrained (143,000 steps), PAR-ROTSCORE remains consistent, indicating that parroting exists for undertrained models and persists, even as models improve. Figure 3 shows results for the Pythia-6.9b model.

Prefix length vs. memorization: Recall that we are measuring *p*-memorization, which depends heavily on the prefix length *p*. In all preceding experiments, we set this number to be at most 80 tokens. To measure how prefix length affects PARROTSCORE, we experimented with reducing *p* to 40, 20, and 10 tokens. This is the maximum size of the prefix that precedes the target PI in the original document. The bottom part of Figure 3 indicates that PARROTSCORE is positively correlated with prefix length, but, even with as little as a 10 token prefix, the 6.9b model can parrot, achieving an average PARROTSCORE of 0.34. This indicates that models memorize PI rampantly and can be prompted with short prompts to parrot PI.

Memorization of constituent parts: We measure how each constituent part of a type of PI is verbatim parroted by each model in Tables 3–5. We split email addresses into

two groups via the '@' symbol: usernames and domains, IP addresses into four groups via each of the three '.' symbols, and phone numbers into two groups: area code and rest of the numbers. For email addresses, both usernames and domains are verbatim parroted often, while for IP addresses each of the four constituents are parroted less often than the preceding one. For phone numbers, area codes are five times more likely to be parroted than the full number, increasing privacy risks as area codes can be tied closely to location. For more details see Appendix E.

6. Related Work

Documentation and PI in Data: The community prioritized documentation, especially personal information, copyright, and autonomy more strongly before the current LLM wave when datasets were smaller (McEnery, 2019). Subramani et al. (2023) analyzed both C4 and the Pile for character-based PI including emails and phone numbers. Concurrently, a preliminary version of these filters was used during the creation of the BigScience ROOTS corpus used to train the BLOOM suite of models (Scao et al., 2022; Laurenccon et al., 2023). Elazar et al. (2023) built on top of this work to develop better regular expressions in the WIMBD suite. Our work improves upon WIMBD by developing better detectors for all PI types and annotating a larger set for better coverage.

Model Memorization and Privacy: Carlini et al. (2021) explore how language models like GPT-2 tend to memorize specific training examples, including PI, and that this can correlate with data frequency and model size. Other work investigate model forgetting, especially tailored to memorized examples throughout training (Jang et al., 2022; Jagielski et al., 2022; Carlini et al., 2022). Our work builds upon these: we quantify character-based PI parroting for the first time and analyze how model size, steps of pretraining, and prefix length affect it, further substantiating the claim that larger, better trained models tend to memorize and parrot more heavily.

7. Conclusion

We develop the regexes and rules (R&R) detection suite for email addresses, US/Canada phone numbers, and IP addresses, improving upon the WIMBD detectors across all PI types. We measure memorization of PI and find that verbatim parroting is rampant, especially as models get larger. This phenomenon is not isolated to larger models; even the smallest models pose privacy risks by parroting personal information verbatim. Consequently, we encourage the community to both develop better PI detectors and carefully filter and anonymize pretraining data when building language models.

References

- Biderman, S., Schoelkopf, H., Anthony, Q. G., Bradley, H., O'Brien, K., Hallahan, E., Khan, M. A., Purohit, S., Prashanth, U. S., Raff, E., et al. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pp. 2397–2430. PMLR, 2023.
- Borkar, J., Jagielski, M., Lee, K., Mireshghallah, N., Smith, D. A., and Choquette-Choo, C. A. Privacy ripple effects from adding or removing personal information in language model training. *ArXiv*, abs/2502.15680, 2025. URL https://api.semanticscholar.org/ CorpusID:276557913.
- Carlini, N., Liu, C., Erlingsson, Ú., Kos, J., and Song, D. The secret sharer: Evaluating and testing unintended memorization in neural networks. In 28th {USENIX} Security Symposium ({USENIX} Security 19), pp. 267– 284, 2019.
- Carlini, N., Tramer, F., Wallace, E., Jagielski, M., Herbert-Voss, A., Lee, K., Roberts, A., Brown, T., Song, D., Erlingsson, U., et al. Extracting training data from large language models. In *30th USENIX Security Symposium* (USENIX Security 21), pp. 2633–2650, 2021.
- Carlini, N., Ippolito, D., Jagielski, M., Lee, K., Tramer, F., and Zhang, C. Quantifying memorization across neural language models. *arXiv preprint arXiv:2202.07646*, 2022.
- Elazar, Y., Bhagia, A., Magnusson, I. H., Ravichander, A., Schwenk, D., Suhr, A., Walsh, E. P., Groeneveld, D., Soldaini, L., Singh, S., Hajishirzi, H., Smith, N. A., and Dodge, J. What's in my big data? In *The Twelfth International Conference on Learning Representations*, 2023.
- Gao, L., Biderman, S., Black, S., Golding, L., Hoppe, T., Foster, C., Phang, J., He, H., Thite, A., Nabeshima, N., et al. The pile: An 800gb dataset of diverse text for language modeling. arXiv preprint arXiv:2101.00027, 2020.
- Groeneveld, D., Beltagy, I., Walsh, E., Bhagia, A., Kinney, R., Tafjord, O., Jha, A., Ivison, H., Magnusson, I., Wang, Y., Arora, S., Atkinson, D., Authur, R., Chandu, K., Cohan, A., Dumas, J., Elazar, Y., Gu, Y., Hessel, J., Khot, T., Merrill, W., Morrison, J., Muennighoff, N., Naik, A., Nam, C., Peters, M., Pyatkin, V., Ravichander, A., Schwenk, D., Shah, S., Smith, W., Strubell, E., Subramani, N., Wortsman, M., Dasigi, P., Lambert, N., Richardson, K., Zettlemoyer, L., Dodge, J., Lo, K., Soldaini, L., Smith, N., and Hajishirzi, H. OLMo: Accelerating the science of language models. In Ku, L.-W., Martins, A., and Srikumar, V. (eds.), *Proceedings of the 62nd Annual*

Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 15789–15809, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.841. URL https://aclanthology.org/2024.acl-long.841/.

- Jagielski, M., Thakkar, O., Tramer, F., Ippolito, D., Lee, K., Carlini, N., Wallace, E., Song, S., Thakurta, A., Papernot, N., et al. Measuring forgetting of memorized training examples. arXiv preprint arXiv:2207.00099, 2022.
- Jang, J., Yoon, D., Yang, S., Cha, S., Lee, M., Logeswaran, L., and Seo, M. Knowledge unlearning for mitigating privacy risks in language models. *arXiv preprint arXiv:2210.01504*, 2022.
- Laurenccon, H., Saulnier, L., Wang, T., Akiki, C., del Moral, A. V., Scao, T. L., von Werra, L., Mou, C., Ponferrada, E. G., Nguyen, H., Frohberg, J., vSavsko, M., Lhoest, O., McMillan-Major, A., Dupont, G., Biderman, S., Rogers, A., Allal, L. B., Toni, F. D., Pistilli, G., Nguyen, O., Nikpoor, S., Masoud, M., Colombo, P., de la Rosa, J., Villegas, P., Thrush, T., Longpre, S., Nagel, S., Weber, L., Muñoz, M. S., Zhu, J., van Strien, D. A., Alyafeai, Z., Almubarak, K., Vu, M. C., Gonzalez-Dios, I., Etxabe, A. S., Lo, K., Dey, M., Suarez, P. O., Gokaslan, A., Bose, S., Adelani, D. I., Phan, L., Tran, H. T., Yu, I., Pai, S., Chim, J., Lepercq, V., Ilic, S., Mitchell, M., Luccioni, S., and Jernite, Y. The bigscience roots corpus: A 1.6tb composite multilingual dataset. ArXiv, abs/2303.03915, 2023. URL https: //api.semanticscholar.org/CorpusID:257378329.
- Li, X. L. and Liang, P. Prefix-tuning: Optimizing continuous prompts for generation. In Zong, C., Xia, F., Li, W., and Navigli, R. (eds.), Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pp. 4582–4597, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long. 353. URL https://aclanthology.org/2021.acl-long. 353/.
- McEnery, T. *Corpus linguistics*. Edinburgh University Press, 2019.
- Microsoft. Presidio data protection and anonymization api, 2021. URL https://github.com/microsoft/presidio. [Release Version 2.2.23, released on Nov 16, 2021].
- Scao, T. L., Fan, A., Akiki, C., Pavlick, E., Ili'c, S., Hesslow, D., Castagn'e, R., Luccioni, A. S., Yvon, F., Gallé, M., Tow, J., Rush, A. M., Biderman, S., Webson, A., Ammanamanchi, P. S., Wang, T., Sagot, B., Muennighoff, N., del Moral, A. V., Ruwase, O., Bawden, R., Bekman, S., McMillan-Major, A., Beltagy, I., Nguyen, H., Saulnier,

L., Tan, S., Suarez, P. O., Sanh, V., Laurenccon, H., Jernite, Y., Launay, J., Mitchell, M., Raffel, C., Gokaslan, A., Simhi, A., Etxabe, A. S., Aji, A. F., Alfassy, A., Rogers, A., Nitzav, A. K., Xu, C., Mou, C., Emezue, C. C., Klamm, C., Leong, C., van Strien, D. A., Adelani, D. I., Radev, D. R., Ponferrada, E. G., Levkovizh, E., Kim, E., Natan, E., Toni, F. D., Dupont, G., Kruszewski, G., Pistilli, G., ElSahar, H., Benyamina, H., Tran, H. T., Yu, I., Abdulmumin, I., Johnson, I., Gonzalez-Dios, I., de la Rosa, J., Chim, J., Dodge, J., Zhu, J., Chang, J., Frohberg, J., Tobing, J., Bhattacharjee, J., Almubarak, K., Chen, K., Lo, K., von Werra, L., Weber, L., Phan, L., Allal, L. B., Tanguy, L., Dey, M., Muñoz, M. R., Masoud, M., Grandury, M., vSavsko, M., Huang, M., Coavoux, M., Singh, M., Jiang, M. T.-J., Vu, M. C., mad A. Jauhar, M., Ghaleb, M., Subramani, N., Kassner, N., Khamis, N., Nguyen, O., Espejel, O., de Gibert, O., Villegas, P., Henderson, P., Colombo, P., Amuok, P., Lhoest, O., Harliman, R., Bommasani, R., L'opez, R., Ribeiro, R., Osei, S., Pyysalo, S., Nagel, S., Bose, S., Muhammad, S. H., Sharma, S. S., Longpre, S., Nikpoor, S., Silberberg, S., Pai, S., Zink, S., Torrent, T. T., Schick, T., Thrush, T., Danchev, V., Nikoulina, V., Laippala, V., Lepercq, V., Prabhu, V., Alyafeai, Z., Talat, Z., Raja, A., Heinzerling, B., Si, C., Salesky, E., Mielke, S. J., Lee, W. Y., Sharma, A., Santilli, A., Chaffin, A., Stiegler, A., Datta, D., Szczechla, E., Chhablani, G., Wang, H., Pandey, H., Strobelt, H., Fries, J. A., Rozen, J., Gao, L., Sutawika, L., Bari, M. S., Al-Shaibani, M. S., Manica, M., Nayak, N. V., Teehan, R., Albanie, S., Shen, S., Ben-David, S., Bach, S. H., Kim, T., Bers, T., Févry, T., Neeraj, T., Thakker, U., Raunak, V., Tang, X., Yong, Z.-X., Sun, Z., Brody, S., Uri, Y., Tojarieh, H., Roberts, A., Chung, H. W., Tae, J., Phang, J., Press, O., Li, C., Narayanan, D., Bourfoune, H., Casper, J., Rasley, J., Ryabinin, M., Mishra, M., Zhang, M., Shoeybi, M., Peyrounette, M., Patry, N., Tazi, N., Sanseviero, O., von Platen, P., Cornette, P., Lavall'ee, P. F., Lacroix, R., Rajbhandari, S., Gandhi, S., Smith, S., Requena, S., Patil, S., Dettmers, T., Baruwa, A., Singh, A., Cheveleva, A., Ligozat, A.-L., Subramonian, A., N'ev'eol, A., Lovering, C., Garrette, D., Tunuguntla, D. R., Reiter, E., Taktasheva, E., Voloshina, E., Bogdanov, E., Winata, G. I., Schoelkopf, H., Kalo, J.-C., Novikova, J., Forde, J. Z., Tang, X., Kasai, J., Kawamura, K., Hazan, L., Carpuat, M., Clinciu, M., Kim, N., Cheng, N., Serikov, O., Antverg, O., van der Wal, O., Zhang, R., Zhang, R., Gehrmann, S., Mirkin, S., Pais, S. O., Shavrina, T., Scialom, T., Yun, T., Limisiewicz, T., Rieser, V., Protasov, V., Mikhailov, V., Pruksachatkun, Y., Belinkov, Y., Bamberger, Z., Kasner, Z., Kasner, Z., Pestana, A., Feizpour, A., Khan, A., Faranak, A., Santos, A. S. R., Hevia, A., Unldreaj, A., Aghagol, A., Abdollahi, A., Tammour, A., HajiHosseini, A., Behroozi, B., Ajibade, B. A., Saxena, B. K., Ferrandis, C. M., Contractor, D.,

6

Lansky, D. M., David, D., Kiela, D., Nguyen, D. A., Tan, E., Baylor, E., Ozoani, E., Mirza, F. T., Ononiwu, F., Rezanejad, H., Jones, H., Bhattacharya, I., Solaiman, I., Sedenko, I., Nejadgholi, I., Passmore, J., Seltzer, J., Sanz, J. B., Fort, K., Dutra, L., Samagaio, M., Elbadri, M., Mieskes, M., Gerchick, M., Akinlolu, M., McKenna, M., Qiu, M., Ghauri, M., Burynok, M., Abrar, N., Rajani, N., Elkott, N., Fahmy, N., Samuel, O., An, R., Kromann, R. P., Hao, R., Alizadeh, S., Shubber, S., Wang, S. L., Roy, S., Viguier, S., Le, T.-C., Oyebade, T., Le, T. N. H., Yang, Y., Nguyen, Z., Kashyap, A. R., Palasciano, A., Callahan, A., Shukla, A., Miranda-Escalada, A., Singh, A. K., Beilharz, B., Wang, B., de Brito, C. M. F., Zhou, C., Jain, C., Xu, C., Fourrier, C., Perin'an, D. L., Molano, D., Yu, D., Manjavacas, E., Barth, F., Fuhrimann, F., Altay, G., Bayrak, G., Burns, G., Vrabec, H. U., Bello, I. I., Dash, I., Kang, J. S., Giorgi, J., Golde, J., Posada, J. D., Sivaraman, K., Bulchandani, L., Liu, L., Shinzato, L., de Bykhovetz, M. H., Takeuchi, M., Pàmies, M., Castillo, M. A., Nezhurina, M., Sanger, M., Samwald, M., Cullan, M., Weinberg, M., Wolf, M., Mihaljcic, M., Liu, M., Freidank, M., Kang, M., Seelam, N., Dahlberg, N., Broad, N. M., Muellner, N., Fung, P., Haller, P., Haller, P., Eisenberg, R., Martin, R., Canalli, R., Su, R., Su, R., Cahyawijaya, S., Garda, S., Deshmukh, S. S., Mishra, S., Kiblawi, S., Ott, S., Sang-aroonsiri, S., Kumar, S., Schweter, S., Bharati, S. P., Laud, T., Gigant, T., Kainuma, T., Kusa, W., Labrak, Y., Bajaj, Y., Venkatraman, Y., Xu, Y., Xu, Y., Xu, Y., Tan, Z. X., Xie, Z., Ye, Z., Bras, M., Belkada, Y., and Wolf, T. Bloom: A 176b-parameter open-access multilingual language model. ArXiv, abs/2211.05100, 2022. URL https: //api.semanticscholar.org/CorpusID:253420279.

- Shin, T., Razeghi, Y., Logan IV, R. L., Wallace, E., and Singh, S. AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts. In Webber, B., Cohn, T., He, Y., and Liu, Y. (eds.), Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 4222–4235, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.346. URL https://aclanthology.org/2020.emnlp-main.346/.
- Soldaini, L., Kinney, R., Bhagia, A., Schwenk, D., Atkinson, D., Authur, R., Bogin, B., Chandu, K., Dumas, J., Elazar, Y., Hofmann, V., Jha, A., Kumar, S., Lucy, L., Lyu, X., Lambert, N., Magnusson, I., Morrison, J., Muennighoff, N., Naik, A., Nam, C., Peters, M., Ravichander, A., Richardson, K., Shen, Z., Strubell, E., Subramani, N., Tafjord, O., Walsh, E., Zettlemoyer, L., Smith, N., Hajishirzi, H., Beltagy, I., Groeneveld, D., Dodge, J., and Lo, K. Dolma: an open corpus of three trillion tokens for language model pretraining research. In Ku, L.-W., Martins, A., and Srikumar, V. (eds.), *Proceedings of the 62nd Annual Meeting of*

the Association for Computational Linguistics (Volume 1: Long Papers), pp. 15725–15788, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.840. URL https://aclanthology.org/2024.acl-long.840/.

- Subramani, N. and Suresh, N. Discovering useful sentence representations from large pretrained language models. *ArXiv*, abs/2008.09049, 2020.
- Subramani, N., Bowman, S. R., and Cho, K. Can unconditional language models recover arbitrary sentences? In *NeurIPS*, 2019.
- Subramani, N., Suresh, N., and Peters, M. Extracting latent steering vectors from pretrained language models. In Muresan, S., Nakov, P., and Villavicencio, A. (eds.), *Findings of the Association for Computational Linguistics: ACL 2022*, pp. 566–581, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-acl.48. URL https: //aclanthology.org/2022.findings-acl.48/.
- Subramani, N., Luccioni, S., Dodge, J., and Mitchell, M. Detecting personal information in training corpora: an analysis. In Ovalle, A., Chang, K.-W., Mehrabi, N., Pruksachatkun, Y., Galystan, A., Dhamala, J., Verma, A., Cao, T., Kumar, A., and Gupta, R. (eds.), *Proceedings* of the 3rd Workshop on Trustworthy Natural Language Processing (TrustNLP 2023), pp. 208–220, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.trustnlp-1.18. URL https://aclanthology.org/2023.trustnlp-1.18/.

A. Limitations

Annotating personal information is time-consuming. Since the data is private, out-sourcing the annotation process should not be done because it could expose PI. As a result, the sets we can annotate are small. Prior work annotated just a few hundred examples (Subramani et al., 2023; Elazar et al., 2023), whereas we annotated 1750 total detections. We hope that larger studies can more comprehensively annotate without exposing privacy risks. Most modern language models do not have open pretraining data, so figuring out what data a model has seen can be challenging. As a result, we focused on using the Pythia model suite because it was one of the only models that had a variety of model sizes, checkpoints throughout pretraining, and open pretraining data. OLMo also has different model sizes, checkpoints, and open pretraining data (Groeneveld et al., 2024), but Dolma (Soldaini et al., 2024), its pretraining corpus, contains a PI filtering and anonymization step using the WIMBD detectors.

During the annotation process, we found that both detectors identify strings of 10 numbers that could be phone numbers, but they are not phone numbers. For example MAXINT=2147483647. 214 is also a Dallas area code, so this could be flagged as a phone number. Additional rules to automatically eliminate detections such as these could help us build better detectors. A further extension of the post processing rules that we did not apply is to filter out subsets of data based on likelihood of false positives. With further study, adding rules about which subset of the Pile certain detectors operate on could decrease the false positive rate greatly.

B. Impact Statement

We hope that our R&R detectors help the community better anonymize and curate pretraining datasets such that the LMs that we deploy in the real world do not expose personal information. In addition, we hope that our analysis showing how significant personal information parroting is by models of all sizes convinces more large language modeling teams to think more carefully about sanitizing pretraining data.

C. R&R Specifics

Here we present the specifics for each detector.

C.1. Email Addresses

The regex used is

(?:[a-z0-9]+(?:\.[a-z0-9!#\$%&'*+/=?^_'{|} ~-]+)*|"(?:[\x01-\x08\x0b\x0c\x0e-\x1f\x21 \x23-\x5b\x5d-\x7f]|\\[\x01-\x09\x0b\x0c \x0e-\x7f])*")@(?:(?:[a-z0-9](?:[a-z0-9-] *[a-z0-9])?\.)+[a-z0-9](?:[a-z0-9-]*[a-z 0-9])?|\[(?:(?:2(?:5[0-5]][0-4][0-9])|1 [0-9][0-9]|[1-9]?[0-9])\.){3}(?:2(?:5 [0-5]|[0-4][0-9])|1[0-9][0-9]|[1-9]?[0-9])|[a-z0-9-]*[a-z0-9]:(?:[\x01-\x08\x0b \x0c\x0e-\x1f\x21-\x5a\x53-\x7f]|\\[\x01 -\x09\x0b\x0c\x0e-\x7f])+\])

As mentioned before, we check for matches with the regular expression and then begin our set of rules to further filter detections. We check whether an "@" character exists to split the addressee and domain. We make sure these are nonempty strings. Next we check whether there exists a starting or trailing period in the domain. If the detected instance has this, we flag this as a false detection. Lastly we make sure that there exists a period (".") in the domain.

C.2. IP Addresses

For IP addresses we have two regexes, an IPv4 and an IPv6 pattern:

ipv4 = (?:(?:25[0-5]|2[0-4][0-9]|[01]?[0-9] [0-9]?)\.){3}(?:25[0-5]|2[0-4][0-9]|[01]? [0-9][0-9]?)

```
ipv6 = (?:^|(?<=\s))(?:(?:[0-9a-fA-F]{1,4}:)
\{7,7\}[0-9a-fA-F]\{1,4\}|(?:[0-9a-fA-F]\{1,4\}:)
\{1,7\}: | (?: [0-9a-fA-F] \{1,4\}:) \{1,6\}: [0-9a-f
A-F]{1,4}|(?:[0-9a-fA-F]{1,4}:){1,5}(?::
[0-9a-fA-F]{1,4}{1,2}|(?:[0-9a-fA-F]{1,4}:)
\{1,4\}(?::[0-9a-fA-F]\{1,4\})\{1,3\}|(?:[0-9a-f
A-F]{1,4}:){1,3}(?::[0-9a-fA-F]{1,4}){1,4}|
(?: [0-9a-fA-F] \{1,4\}:) \{1,2\} (?:: [0-9a-fA-F]
\{1,4\} \{1,5\} [0-9a-fA-F] \{1,4\} : (?:(?::[0-9a-f
A-F]{1,4}{1,6}|:(?:(?::[0-9a-fA-F]{1,4})
{1,7}:)|fe80:(?:(?::[0-9a-fA-F]{0,4}){0,4}
%[0-9a-zA-Z]{1,})::(?:ffff(?::0{1,4}){0,1}
:){0,1}(?:(?:(?:25[0-5]|(?:2[0-4]|1{0,1}
[0-9] \{0,1\} [0-9] \} \{3,3\} (?:25[0-5] | (?:2
[0-4]|1{0,1}[0-9]){0,1}[0-9]))|(?:(?:
[0-9a-fA-F]{1,4}:){1,4}(?:(?:25[0-5])(?:2)
[0-4]|1{0,1}[0-9]){0,1}[0-9]).}{3,3}(?:25
[0-5]|(?:2[0-4]|1{0,1}[0-9]){0,1}[0-9])))
(?=\s|$)
```

After running the regular expression based detectors, we filter the detected IP addresses using the following set of rules. First, we check if any of the following common words occur in the micro context window of 20 characters preceding the detected PI span:

```
'isbn', 'doi', '#', 'grant', 'award', 'nsf',
'patent', 'usf', 'edition', 'congress',
'appeal', 'claim', 'exhibit', 'serial',
```

```
'pin', 'receipt', 'case', 'tracking',
'ticket', 'route', ' wo ', 'volume'
```

These words often indicate a type of number that could have syntax similar to an IPv4 address. This is a much larger problem for phone numbers, so we also do this for phone numbers. Next, we check whether the prefix has alphabet characters. We take at most the 50 characters preceding the detected pi span and see if at least 10% of them are alpha numeric. This is to filter out arbitrary sets of numbers.

C.3. Phone Numbers

For phone numbers we have two regular expression based detectors *phone numbers*:

\s+\(?(\d{3})\)?[-\.]*(\d{3})[-.]?(\d{4})
(?!\d)

and *phone numbers global* (US/Canadian phone numbers in a global context with the country code):

```
\s+(?:\+1|1)[-\. ]*\(?(\d{3})\)?[-\. ]*
(\d{3})[-\. ]?(\d{4})(?!\d)
```

After running the regular expression based detectors, we filter using a set of rules. First, we check if any of the common words (the words used when filtering IP addresses above) occur in the micro context window of 20 characters preceding the detected PI span. These words often indicate a type of number that can often have 9, 10, or 11 digits. Next, we check whether the prefix has siffucient alphabet characters, identical to how IP addresses were processed. For phone numbers, this helps filter out things like html polygons or random sets of numbers without context like dumps of arbitrary numbers. Next, we standardize the detected number and validate the area code: excluding numbers with an area code starting with a 0 or 1 and verify that the area code is a valid one. After doing this, we validate the central office code. Lastly, we exclude a set of placeholder numbers including 1234567890, 2345678910, MAXINT, 7373737373, and 3141592653 (the digits of π).

D. Ablation Analysis

Here, we look at the impact of pretraining steps and prefix length on memorization. In Figure 3, we find that models parrot even when only halfway through training. The Pythia models are trained for 143,000 steps and even from 70,000 steps as mentioned before, PARROTSCOREremains constant. Additionally, prefix length is highly correlated with PAR-ROTSCORE. However, even with as little as 10 tokens in the prefix, PI memorization is rampant, indicating severe risk.



Figure 3: The effect of the number of pretraining steps (top) and prefix length (bottom) on PARROTSCORE across PI types for the Pythia-6.9b model.

E. Memorization of Constituent Parts:

Here, we measure how each constituent part of a type of PI is verbatim parroted by each model. To do this, we first parse IP addresses (IPv4) into their four constituent groups separated by a period (e.g. 8.8.8.8 turns into [8,8,8,8]). Each of these 4 groups are measured separately. A candidate generation that outputs "12.8.8 abcd" will turn into [12, 8, 8 abcd, ""] and comparing that to 8.8.8.8 will lead to verbatim parroting of only group 2. We parse email addresses into two groups: the username and domain separated by the '@' symbol because email addresses are normally separated into these two groups. We parse phone numbers into two groups: area code and the rest of the digits following that. Since we are only looking at US/Canada phone numbers that always start with a 1, we did not emphasize splitting out the country code. We measure verbatim parroting for all model sizes at the 143,000 iteration checkpoint for a prefix length of 80. This is an extension of Table 2, where we report percent of instances that are verbatim parroted.

Personal Information Parroting in Language Models

model sizes	username	domain
160m	11.22%	12.15%
410m	15.42%	20.09%
1b	20.09%	24.77%
1.4b	20.56%	28.50%
2.8b	25.70%	31.31%
6.9b	26.17%	30.84%

model sizes area code rest of the number 160m 8.20% 1.64% 4.92% 410m 13.11% 1b 11.48% 3.28% 1.4b 14.75% 3.28% 2.8b 9.84% 1.64% 6.9b 19.67% 3.28%

Table 3: Percent of total instances that are verbatim parroted for the constituent parts of email addresses (the username and domain). Note that this corresponds to the percent of instances with that component achieving a PARROTSCORE of 1. We find that verbatim parroting increases with model size.

Table 5: Percent of total instances that are verbatim parroted for the constituent parts of phone numbers (split by area code and rest of the number). Note that this corresponds to the percent of instances with that component achieving a PARROTSCORE of 1. We find that verbatim parroting increases with model size for area codes generally, but not for the rest of the number.

In Table 3, we find that for email addresses, both the usernames and domains are verbatim parroted often, as much as 31% of instances have a parroted domain. Usernames are slightly less parroted, but for even those up to 26.17% of instances have a verbatim parroted username, underscoring significant risk. For IP addresses in Table 4, each successive constituent was slightly less likely to be parroted than the previous group. We hypothesize that this is due to the nature of left-to-right autoregressive decoding and not due to any other confounding factors. For the 6.9b model, in about 32% of instances either group 1 or group 2 were verbatim parroted, whereas only 14% of IP addresses overall were verbatim parroted.

model sizes	grp1	grp2	grp3	grp4
160m	18.11%	13.39%	13.39%	9.45%
410m	18.90%	19.69%	14.96%	11.02%
1b	21.26%	21.26%	20.47%	11.02%
1.4b	25.20%	24.41%	21.26%	14.17%
2.8b	26.77%	25.20%	22.83%	18.90%
6.9b	32.28%	31.50%	24.41%	19.69%

Table 4: Percent of total instances that are verbatim parroted for the constituent parts of IP addresses (split on each '.' character grp1-grp4). Note that this corresponds to the percent of instances with that component achieving a PAR-ROTSCORE of 1. We find that verbatim parroting increases with model size. For phone numbers in Table 5, which had a relatively low verbatim parrot rate ($\sim 3\%$), we find that area codes are much more likely to be parroted, increasing privacy risk as this can be tied closely to location. We find that for the 6.9b model, even when only 3.28% of instances are parroted by the model, in 19.67% of cases the area code is parroted exactly. Taken together, these results further underscore our point that PI memorization and parroting is a risk that needs to be mitigated.