

| OpenDP Privacy Attacks & Auditing Working Group   |   |  |   |  |              |                  |                      |                      |   |  |   |
|---|---|--|---|--|--------------|------------------|----------------------|----------------------|---|--|---|
| Privacy Attacks Repository  |   |  |   |  |              |                  |                      |                      |   |  |   |
| URL   | BitTex (Please add a bittex entry for this paper to facilitate writing our summary document)  | Authors  | Title   | Short Description  | Type of Data | Type of Release  | Threat Model         | Research Type        | Links to Artifacts  | Comments   | Submitter (your name, affiliation)                                  |
| <a href="https://dl.acm.org/doi/10.1145/77153.771771">https://dl.acm.org/doi/10.1145/77153.771771</a>   | @inproceedings{10.1145/77153.771771, author = {Dinar, Iti and Kobbi, Nisim}, title = {Revealing Information While Preserving Privacy}, year = {2003}, }   | Iti Dinar and Kobbi Nisim  | Revealing Information While Preserving Privacy  | Seminal paper on the theory of data reconstruction attacks.  | Tabular      | Linear-Queries   | Reconstruction       | Theoretical          |   |  | Jon Ullman, Northeastern University                                 |
| <a href="https://arxiv.org/abs/1810.05620">https://arxiv.org/abs/1810.05620</a>   | @article{john2018linear, title={Linear program reconstruction in practice}, author={Cohen, Aloni and Nisim, Kobbi}, year={2018}}  | Aloni Cohen and Kobbi Nisim  | Linear program reconstruction in practice.  | Implemented linear reconstruction attacks against a production private query system called Diffrax   | Tabular      | Linear-Queries   | Reconstruction       | Applications         |   |  | Jon Ullman, Northeastern University                                 |
| <a href="https://arxiv.org/abs/2405.10999">https://arxiv.org/abs/2405.10999</a>   | @article{janamal2024you, title="What do you want from theory alone?" Experimenting with Tight Auditing of Differentially Private Synthetic Data Generation}, author={Janamal, Meenatchi Sundaram Muthu Selva and Ganey, Georgi and De Cristofaro, Emiliano}, journal={USENIX Security}, year={2024}}  | Meenatchi Sundaram Muthu Selva Ananias, Georgi Ganey, Emiliano De Cristofaro                               | "What do you want from theory alone?" Experimenting with Tight Auditing of Differentially Private Synthetic Data Generation | Audits six implementations of DP synthetic data generative models using different datasets and threat models and finds that commonly used black-box MIA are severely limited in power, yielding remarkably loose empirical privacy estimates. Considers MIA in stronger threat models, i.e., passive and active white-box, using both existing and newly proposed attacks. | Tabular      | Generative-Model | Membership-Inference | Empirical            | <a href="https://github.com/koala/sharvict-h-audit">https://github.com/koala/sharvict-h-audit</a> |  | Georgi Ganey, UCL   |
| <a href="https://proceedings.mlr.press/v160/jagielski2022a.html">https://proceedings.mlr.press/v160/jagielski2022a.html</a>   | @inproceedings{10.1145/3546158.3546167, author={Jagielski, Matthew and Ullman, Jonathan and Oprea, Alina}, journal={Advances in Neural Information Processing Systems}, year={2022}, }  | Matthew Jagielski, Jonathan Ullman, Alina Oprea  | Auditing Differentially Private Machine Learning: How Private is Private SGD?   | Connects the success rate of a membership inference adversary to a lower bound on the privacy loss of the underlying DP mechanism.   | Image        | Predictive-Model | Membership-Inference | Empirical            |   |  | Tudor Cabero, Inria   |
| <a href="https://www.pnas.org/doi/10.1073/pnas.2118605120">https://www.pnas.org/doi/10.1073/pnas.2118605120</a>   | @article{dik2022confidence, title={Confidence-ranked reconstruction of census microdata from published statistics}, author={Dik, Travis and Dwork, Cynthia and Kearns, Michael and Liu, Terrance and Roth, Aaron and Vietri, Giuseppe and Wu, Zhilue Steven}, journal={Proceedings of the National Academy of Sciences}, volume={119}, number={35}, pages={62218605120}, year={2022}, publisher={National Acad Sciences}} | Travis Dik, Cynthia Dwork, Michael Kearns, Terrance Liu, Aaron Roth, Giuseppe Vietri, and Zhilue Steven Wu | Confidence-ranked reconstruction of census microdata from published statistics  | Ranking rows in reconstructed microdata by how confident the adversary is that are in the true dataset.  | Tabular      | Linear-Queries   | Membership-Inference | Empirical            |   | I didn't want to give it its own row but linking a rebuttal paper which is interesting in that it is indicative of arguments people make against reconstruction attacks: <a href="https://dl.acm.org/doi/10.1145/3546158.3546171">https://dl.acm.org/doi/10.1145/3546158.3546171</a> | Audra McMillan, Apple   |
| <a href="https://dl.acm.org/doi/10.1145/348698.3486981">https://dl.acm.org/doi/10.1145/348698.3486981</a>   | @inproceedings{mhu2022querysnout, title={QuerySnout: Automating the discovery of attribute inference attacks against query-based systems}, author={Cretu, Ana-Maria and Houssiau, Florimond and Cully, Antoine and de Montjoye, Yves-Alexandre}, booktitle={Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security}, pages={623–637}, year={2022}}   | Ana-Maria Cretu, Florimond Houssiau, Antoine Cully, Yves-Alexandre de Montjoye                             | QuerySnout: Automating the Discovery of Attribute Inference Attacks against Query-Based Systems                             | A method to automatically discover attacks against interactive query systems   | Tabular      | Linear-Queries   | Attribute-Inference  | Empirical            |   |  | Ana-Maria Cretu, EPFL   |
| <a href="https://www.ndss-symposium.org/ndss-conference/attachments/18_058-5_Pyrgelis.pdf">https://www.ndss-symposium.org/ndss-conference/attachments/18_058-5_Pyrgelis.pdf</a> | @article{pyrgelis2017knock, title={Knock knock, who's there? Membership inference on aggregate location data}, author={Pyrgelis, Carmela and Troncoso, Carmela and De Cristofaro, Emiliano}, journal={arXiv preprint arXiv:1708.06145}, year={2017}}  | Apostolos Pyrgelis, Carmela Troncoso, Emiliano De Cristofaro   | Knock Knock, Who's There? Membership Inference on Aggregate Location Data   | Membership inference attacks against location aggregates   | Tabular      | Linear-Queries   | Membership-Inference | Empirical            |   |  | Ana-Maria Cretu, EPFL, Yves-Alexandre de Montjoye, Imperial College |
| <a href="https://arxiv.org/abs/2003.10053">https://arxiv.org/abs/2003.10053</a>   | @inproceedings{janamal2023near, title={A linear reconstruction approach for attribute inference attacks against synthetic data}, author={Janamal, Meenatchi Sundaram Muthu Selva and Ganey, Georgi and Reher, Luc}, booktitle={Usenix Security}, year={2024}}   | Meenatchi Sundaram Muthu Selva Ananias, Andrea Galotti, Luc Reher  | A linear reconstruction approach for attribute inference attacks against synthetic data                                     | Introduces a new attribute inference attack against synthetic data based on linear reconstruction methods for aggregate statistics, which target all records in the dataset, not only outliers.  | Tabular      | Generative-Model | Attribute Inference  | Empirical            | <a href="https://github.com/yvrtic/synthetic-synth">https://github.com/yvrtic/synthetic-synth</a> |  | Georgi Ganey, UCL   |
| <a href="https://arxiv.org/abs/2206.10465">https://arxiv.org/abs/2206.10465</a>   | @article{carlini2022privacy, title={The privacy onion effect: Memorization is relative}, author={Carlini, Nicholas and Jagielski, Matthew and Zhang, Chiyuan and Papernot, Nicolas and Terzis, Andreas and Tramèr, Florian}, journal={Advances in Neural Information Processing Systems}, volume={35}, pages={13263–13276}, year={2022}}  | Carlini et al.   | The Privacy Onion Effect: Memorization is Relative  | Removing vulnerable records makes other records vulnerable   | Tabular      | Predictive-Model | Membership-Inference | Empirical            |   |  | Yves-Alexandre de Montjoye, Imperial College                        |
| <a href="https://arxiv.org/abs/2112.03970">https://arxiv.org/abs/2112.03970</a>   | @inproceedings{carlini2022membership, title={Membership inference attacks from first principles}, author={Carlini, Nicholas and Chen, Steve and Nasr, Maud and Song, Shuang and Terzis, Andreas and Tramèr, Florian}, booktitle={2022 IEEE Symposium on Security and Privacy (SP)}, pages={1897–1914}, year={2022}, organization={IEEE}}  | Carlini et al.   | Membership Inference Attacks From First Principles  | Takes a step back - new attack (LIRA)  | Tabular      | Predictive-Model | Membership-Inference | Theoretical          |   |  | Yves-Alexandre de Montjoye, Imperial College                        |
| <a href="https://arxiv.org/abs/2020.14051">https://arxiv.org/abs/2020.14051</a>   | @article{geiping2020inverting, title={Inverting gradients-how easy is it to break privacy in federated learning?}, author={Geiping, Jonas and Saemundsson, Hartmut and DiTolone, Hannah and Meibler, Michael}, journal={Advances in neural information processing systems}, volume={33}, pages={16937–16947}, year={2020}}  | Geiping et al.   | Inverting Gradients – How easy is it to break privacy in federated learning?  | Gradient inversion attack (reconstruct data point from gradient) - application to federated learning   | Image        | Predictive-Model | Reconstruction       | Empirical            | <a href="https://github.com/koala/sharvict-h-audit">https://github.com/koala/sharvict-h-audit</a> |  | Aurélien Bellet, Inria  |
| <a href="https://dl.acm.org/doi/10.1145/3546158.3546167">https://dl.acm.org/doi/10.1145/3546158.3546167</a>   | @inproceedings{ormrod2013empirical, title={Empirical privacy and empirical utility of anonymized data}, author={Ormrod, Graham and Procopiac, Cecilia M and Shen, Erling and Shmatikov, Divyesh and Yu, T}, booktitle={2013 IEEE 20th International Conference on Data Engineering Workshops (ICDEW)}, pages={77–82}, year={2013}, organization={IEEE}}   | Ormrod, G., Procopiac, C.M., Shen, E., Srivastava, D. and Yu, T  | Empirical privacy and empirical utility of anonymized data.   | Naive Bayes attacks against simple workloads of statistics   | Tabular      | Linear-Queries   | Reconstruction       | Empirical            |   |  | James Honaker, Anonym   |
| <a href="https://proceedings.mlr.press/v48/dwork17.html">https://proceedings.mlr.press/v48/dwork17.html</a>   | @article{dwork2017exposed, title={Exposed: a survey of attacks on private data}, author={Dwork, Cynthia and Smith, Adam and Steinke, journal={Annual Review of Statistics and Its Application}, volume={4}, number={1}, pages={81–84}, year={2017}, publisher={Annual Reviews}}   | Dwork, C., Smith, A., Steinke, T. and Ullman, data.  | Exposed: a survey of attacks on private data.   | Survey of privacy attacks  | Tabular      |                  |                      | Theoretical          |   |  | James Honaker, Anonym   |
| <a href="https://proceedings.mlr.press/v160/staierke2022a.html">https://proceedings.mlr.press/v160/staierke2022a.html</a>   | @article{staierke2022a, title={Privacy auditing with one (1) training run}, author={Staierke, Thomas and Nasr, Maud and Jagielski, journal={Advances in Neural Information Processing Systems}, volume={35}, year={2022}}   | Staierke, T., Nasr, M. and Jagielski, M.   | Privacy Auditing with One (1) Training Run.   |  |              |                  |                      | Membership-Inference |   |  | James Honaker, Anonym   |
| <a href="https://arxiv.org/abs/2017.05820">https://arxiv.org/abs/2017.05820</a>   | @inproceedings{shokri2017membership, title={Membership inference attacks against machine learning models}, author={Shokri, Reza and Stronati, Marco and Song, C}, booktitle={2017 IEEE Symposium on Security and Privacy (SP)}, pages={1–18}, year={2017}, organization={IEEE}}   | Shokri, R., Stronati, M., Song, C. and Shmatik   | Membership inference attacks against machine learning models.   | The original membership inference based on shadow models.  |              |                  |                      | Membership-Inference |   |  | James Honaker, Anonym   |
| <a href="https://arxiv.org/abs/1811.00927">https://arxiv.org/abs/1811.00927</a>   | @misc{wagner2018technical, title={Technical privacy metrics: a systematic survey}, author={Wagner, I. and Eckhoff, D.}}   | Wagner, I. and Eckhoff, D.   | Technical privacy metrics: a systematic survey.   | Survey of privacy-loss metrics   | Tabular      | Linear-Queries   | Information Leakage  | Theoretical          |   |  | James Honaker, Anonym   |
| <a href="https://arxiv.org/abs/2022.04599">https://arxiv.org/abs/2022.04599</a>   | @misc{gliomi2022unified, title={A Unified Framework for Quantifying Privacy Risk in Synthetic Data}, author={Matteo Gliomi and Franziska Bonisch and Christoph Wöhner and Borbála Tászadi}, eprint={2211.10459}, archivePrefix={arXiv}, primaryClass={cs.CR}, url={https://arxiv.org/abs/2211.10459}}   | Gliomi, M., Boenisch, F., Wöhner, C., and Tászadi, B.  | A Unified Framework for Quantifying Privacy Risk in Synthetic Data  | Adversarial evaluation of singling out, linkability, and inference risk in tabular synthetic data  | Tabular      |                  |                      | Empirical            | <a href="https://github.com/StarForce04/unified">https://github.com/StarForce04/unified</a>       |  | Matteo Gliomi, Anoncs   |
| <a href="https://arxiv.org/abs/2211.10459">https://arxiv.org/abs/2211.10459</a>   | @inproceedings{gliomi2022unified, title={A unified framework for quantifying privacy risk in synthetic data}, author={Gliomi, Matteo and Boenisch, Franziska and Wöhner, Christoph and Tászadi, Borbála}, booktitle={PETS}, year={2022}}  | Matteo Gliomi, Franziska Boenisch, Christoph Wöhner, Borbála Tászadi                                       | A Unified Framework for Quantifying Privacy Risk in Synthetic Data  | Present Anonymizer, a statistical framework to jointly quantify different types of privacy risks in synthetic tabular datasets. Equips this framework with attack-based evaluations for the singling out, linkability, and inference risks, the three key indicators of factual anonymization according to the GDPR.   | Tabular      | Generative-Model | Membership-Inference | Empirical            |   |  | Georgi Ganey, UCL   |

OpenDP Privacy Attacks & Auditing Working Group

Privacy Attacks Repository

| URL   | BitTex (Please add a bittex entry for this paper to facilitate writing our summary document)   | Authors   | Title   | Short Description   | Type of Data | Type of Release  | Threat Model         | Research Type    | Links to Artifacts  | Comments | Submitter (your name, affiliation)      |
|---|--|---|---|---|--------------|------------------|----------------------|------------------|---|----------|---|
| <a href="https://arxiv.org/abs/2110.16789">https://arxiv.org/abs/2110.16789</a>                 | <pre>@article{sh2023detecting,   title={Detecting pretraining data from large language models},   author={Shi, Weijia and Ajith, Anirudh and Xia, Mengzhou and Huang, Yixuan and Liu, Daogang and Blewett, Tara and Chen, Dang and Zettlemoyer, Luke},   journal={arXiv preprint arXiv:2310.16789},   year={2023} }</pre>  | Weijia Shi et al.   | Detecting Pretraining Data from Large Language Models                               | introduces Min-K prob attack: Membership Inference attack against LLMs using average of lowest k probable tokens in the target sequence.  | Text         | Generative-Model | Membership-Inference | Applications     | <a href="https://arxiv.org/abs/2110.16789">https://arxiv.org/abs/2110.16789</a>   |          | Hamid Mozaffari, Oracle Labs            |
| <a href="https://arxiv.org/abs/2402.07841">https://arxiv.org/abs/2402.07841</a>                 | <pre>@misc{duan2024membership,   title={Do Membership Inference Attacks Work on Large Language Models?},   author={Michael Duan and Anshuman Suri and Nilofar Mirshghalali and Sewon Min and Weijia Shi and Luke Zettlemoyer and Yulia Tsvetkov and Yiqin Choi and David Evans and Hananeh Hajishirzi},   year={2024},   eprint={2402.07841},   archivePrefix={arXiv},   primaryClass={cs.CL},   full_name="Computation and language" is_active=True alt_name="lp" in_archive="cs" is_general=False description="Covers natural language processing. Roughly includes material in ACM Subject Class I.2.7. Note that work on artificial languages [programming languages, logics, formal systems] that does not explicitly address natural-language issues broadly construed (natural-language processing, computational linguistics, speech, text retrieval, etc.) is not appropriate for this area." }</pre> | Michael Duan, Anshuman Suri, Nilofar Mirshghalali, Sewon Min, Weijia Shi, Luke Zettlemoyer, Yulia Tsvetkov, Yiqin Choi, David Evans, Hananeh Hajishirzi                                   | Do Membership Inference Attacks Work on Large Language Models?                      | Examining MIA attacks on LLM, proposing a gptub rep set with set of standard attacks  | Text         | Generative-Model | Membership-Inference | Empirical        | <a href="https://github.com/lanamgoc42/gptub">https://github.com/lanamgoc42/gptub</a>   |          | Danil Filenko, University of Washington |
| <a href="https://arxiv.org/abs/2101.13188">https://arxiv.org/abs/2101.13188</a>                 | <pre>@article{carlini2023extracting,   title={Extracting training data from diffusion models},   author={Carlini, Nicholas and Hayes, Jamie and Nasr, Mads and Jagielski, Matthew and Selwag, Florian and Tramèr, Vít and Florian and Balta, Borja and Ippolito, Daphne and Wallace, Eric},   journal={IEEE Security &amp; Privacy},   year={2023} }</pre>   | Nicholas Carlini, Jamie Hayes, Mads Nasr, Matthew Jagielski, Vikash Selwag, Florian Tramèr, Borja Balta, Daphne Ippolito, Eric Wallace  | Extracting training data from diffusion models                                      | Show that diffusion models memorize individual images from their training data and emit them at generation time. With a generate-and-filter pipeline, extracts over a thousand training examples from state-of-the-art models, ranging from photographs of individual people to trademarked company logos.                                    | Image        | Generative-Model | Data-Extraction      | Empirical        |   |          | Georgi Ganev, UCL                       |
| <a href="https://arxiv.org/abs/2012.07826">https://arxiv.org/abs/2012.07826</a>                 | <pre>@inproceedings{carlini2021extracting,   title={Extracting training data from large language models},   author={Carlini, Nicholas and Tramèr, Florian and Wallace, Eric and Jagielski, Matthew and Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Tom and Song, Dawn and Erlingsson, Úlfar and others},   year={2021} }</pre>   | Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Úlfar Erlingsson, Alina Oprea, and Colin Raffel | Extracting training data from large language models                                 | Showed how to prompt models like GPT to reveal specific training examples   | Text         | Generative-Model | Data-Extraction      | Applications     |   |          | Jon Ullman, Northeastern University     |
| <a href="https://arxiv.org/abs/2211.06564">https://arxiv.org/abs/2211.06564</a>                 | <pre>@article{houssiau2023tapas,   title={Tapas: a toolbox for adversarial privacy auditing of synthetic data},   author={Houssiau, Florimond and Jordan, James and Cohen, Samuel N and Daniel, Owen and Elliott, Andrew and Geddes, James and Mole, Callum and Rangeli-Smith, Camilla and Szpruch, Lukasz},   journal={arXiv preprint arXiv:2211.06564},   year={2022} }</pre>  | Florimond Houssiau, James Jordan, Samuel N. Cohen, Owen Daniel, Andrew Elliott, James Geddes, Callum Mole, Camilla Rangeli-Smith, Lukasz Szpruch  | Tapas: Toolbox for Adversarial Privacy Auditing of Synthetic Data                   | A system of classification of MIA attacks presented as part of a toolbox of attacks to evaluate synthetic data privacy under a wide range of scenarios  | Tabular      |                  |                      |                  |   |          | Danil Filenko, University of Washington |
| <a href="https://arxiv.org/abs/2109.03931">https://arxiv.org/abs/2109.03931</a>                 | <pre>@inproceedings{chen2020gan,   title={GAN-lacks: a taxonomy of membership inference attacks against generative models},   author={Chen, Dingfan and Yu, Ning and Zhang, Yang and Wang, Yizhou and Chen, Yizhou},   booktitle={ACM CCS},   year={2020} }</pre>  | Dingfan Chen, Ning Yu, Yang Zhang, Maria Fritz  | Gan-lacks: a taxonomy of membership inference attacks against generative models     | Presents a taxonomy of membership inference attacks, encompassing not only existing attacks but also our novel ones. Moreover, provides a theoretically grounded attack calibration technique, which consistently boosts the attack performance in all cases, across different attack settings, data modalities, and training configurations. | Image        | Generative-Model | Membership-Inference | Empirical        |   |          | Georgi Ganev, UCL                       |
| <a href="https://arxiv.org/abs/2109.07661">https://arxiv.org/abs/2109.07661</a>                 | <pre>@inproceedings{hayes2019logan,   title={LOGAN: membership inference attacks against generative models},   author={Hayes, Jamie and Melis, Luca and Danzic, George and De Cristofaro, Emiliano},   booktitle={PPE'19},   year={2019} }</pre>   | Jamie Hayes, Luca Melis, George Danzic, Emiliano De Cristofaro  | LOGAN: Membership Inference Attacks Against Generative Models                       | Presents the first membership inference attacks against generative models (GANs) given a data point, the adversary determines whether or not it was used to train the model.  | Image        | Generative-Model | Membership-Inference | Empirical        |   |          | Georgi Ganev, UCL                       |
| <a href="https://arxiv.org/abs/2112.05301.pdf">https://arxiv.org/abs/2112.05301.pdf</a>         | <pre>@inproceedings{jin2022we,   title={Are we there yet? Timing and floating-point attacks on differential privacy systems},   author={Jin, Jiarui and McMurtry, Eleanor and Rubinfeld, Benjamin and Shrivastava, Abhinav and Shrivastava, Abhinav and Shrivastava, Abhinav},   pages={473--488},   year={2022},   organization={IEEE} }</pre>  | Jin, McMurtry, Rubinfeld, Shrivastava   | Are We There Yet? Timing and Floating-Point Attacks on Differential Privacy Systems | Attacks on DP implementations that use floating-point arithmetic: time side channels on discrete samplers   |              |                  |                      | Reconstruction   |   |          | Zachary Ratliff, Harvard + OpenDP       |
| <a href="https://arxiv.org/abs/2109.18462">https://arxiv.org/abs/2109.18462</a>                 | <pre>@article{mattam2023membership,   title={Membership inference attacks against language models via neighbourhood comparison},   author={Mattam, Justus and Mirshghalali, Fahimeh and Jin, Zhijing and Shi, Yiqin and Taylor, Thomas and Sachan, Mihir and Berg-Kirkpatrick, Taylor},   journal={arXiv preprint arXiv:2305.18462},   year={2023} }</pre>   | Justus Mattam et al.  | Membership Inference Attacks against Language Models via Neighbourhood Comparison   | Membership inference attack against LLMs by generating neighbours which are generated using a mask model like BERT  | Text         | Generative-Model | Membership-Inference | Applications     | <a href="https://github.com/justusmattam/neighborhood-mia">https://github.com/justusmattam/neighborhood-mia</a>   |          | Hamid Mozaffari, Oracle Labs            |
| <a href="https://arxiv.org/abs/2112.03922">https://arxiv.org/abs/2112.03922</a>                 | <pre>@inproceedings{zarfaden2024low,   title={Low-Cost High-Power Membership Inference Attacks},   author={Zarfaden, Sajjad and Liu, Philippe and Shari, Reza},   booktitle={Forty-first International Conference on Machine Learning},   year={2024} }</pre>  | Sajjad Zarfaden, Philippe Liu, Reza Shari   | Low-Cost High-Power Membership Inference Attacks                                    | Efficient MIA with shadow models and auxiliary reference data. Outperforms UPIA in cases where one has plenty of reference data and almost no reference models  | Tabular      |                  |                      | Predictive-Model |   |          | Luca Melis, Meta                        |
| <a href="https://arxiv.org/abs/2107.07569">https://arxiv.org/abs/2107.07569</a>                 | <pre>@article{van2023membership,   title={Membership inference attacks against synthetic data through overfitting detection},   author={van Bruegel, Boris and Sun, Hao and Qian, Zhaoyan and van der Schaar, Minko},   journal={JSTATS},   year={2023} }</pre>  | Boris van Bruegel, Hao Sun, Zhaoyan Qian, Minko van der Schaar  | Membership Inference attacks against synthetic data through overfitting detection   | Proposes DOMAS, a density-based MIA model that aims to infer membership by targeting local overfitting of the generative model assuming the attack has some knowledge of the underlying data distribution.  | Tabular      | Generative-Model | Membership-Inference | Empirical        |   |          | Georgi Ganev, UCL                       |
| <a href="https://arxiv.org/abs/2111.08629">https://arxiv.org/abs/2111.08629</a>                 | <pre>@inproceedings{ye2022enhanced,   title={Enhanced membership inference attacks against machine learning models},   author={Ye, Jiyuan and Maddi, Sasi Kumar Marakonda, Vincent Bindhuwader, Raza Shokri, Reza},   booktitle={Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security},   pages={3093--3106},   year={2022} }</pre>   | Jiyuan Ye, Aadyaa Maddi, Sasi Kumar Marakonda, Vincent Bindhuwader, Raza Shokri   | Enhanced Membership Inference Attacks against Machine Learning Models               | Framework for MIA based on approximate ERT. Theory behind general class of attacks + a few instantiations.  | Tabular      | Predictive-Model | Membership-Inference | Empirical        |   |          |   |
| <a href="https://arxiv.org/abs/2208.14932">https://arxiv.org/abs/2208.14932</a>                 | <pre>@inproceedings{liu2022membership,   title={Membership inference attacks by exploiting loss},   author={Liu, Yiyong and Zhao, Zhengyu and Bades, M},   booktitle={Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security},   pages={2088--2098},   year={2022} }</pre>  | Yiyong Liu, Zhengyu Zhao, Michael Backes, Y   | Membership inference attacks by exploiting loss trajectories                        | Membership inference attack based on knowledge distillation to form signals to attack via loss trajectories over multiple epochs.   | Image        | Predictive-Model | Membership-Inference | Empirical        | <a href="https://github.com/Danish142022/Membership-Inference-Attacks-by-Exploiting-Loss-Trajectory">https://github.com/Danish142022/Membership-Inference-Attacks-by-Exploiting-Loss-Trajectory</a> |          | Johan Östman, AI Sweden                 |
| <a href="https://arxiv.org/abs/2107.07636">https://arxiv.org/abs/2107.07636</a>                 | <pre>@article{bertram2023scalable,   title={Scalable membership inference attacks via quantile regression},   author={Bertram, Martin and Tang, Shuai and Roth, Aaron and Kearns, Michael and Morgenstern, James H and Wu, Steven W},   journal={Advances in Neural Information Processing Systems},   volume={36},   year={2024} }</pre>  | Martin Bertram, Shuai Tang, Michael Kearns, James Morgenstern, Aaron Roth, Zhewei Steven Wu   | Scalable Membership Inference Attacks via Quantile Regression                       | Membership inference attack that does not require shadow models but only to train a single regression model to predict quantiles of the logits.   | Image        | Predictive-Model | Membership-Inference | Empirical        |   |          | Johan Östman, AI Sweden                 |
| <a href="https://arxiv.org/abs/2007.14321">https://arxiv.org/abs/2007.14321</a>                 | <pre>@inproceedings{choquetta2021label,   title={Label-only membership inference attacks},   author={Choquetta, Choo, Christopher A and Tramèr, Florian and Carlini, Nicholas and Papernot, Nicolas},   booktitle={International conference on machine learning},   pages={1964--1974},   year={2021},   organization={PMLR} }</pre>   | Christopher A. Choquetta-Choo, Florian Tramèr, Nicholas Carlini, Nicolas Papernot   | Label-only membership inference attacks   | Black-box membership inference attack with label-only signals. Signals are created by probing the model with permutations around a given datapoint.   | Image        | Predictive-Model | Membership-Inference | Empirical        |   |          | Johan Östman, AI Sweden                 |
| <a href="https://openreview.net/forum?id=VwvYwv95">https://openreview.net/forum?id=VwvYwv95</a> | <pre>@inproceedings{wu2024you,   title={You Only Query Once: An Efficient Label-Only Membership Inference Attack},   author={Wu, Yutong and Qiu, Han and Guo, Shangwei and Li, Jiewei and Zhang, Tianwei},   booktitle={The Twelfth International Conference on Learning Representations},   year={2024} }</pre>   | Yutong Wu, Han Qiu, Shangwei Guo, Jiewei Li, Tianwei Zhang  | You only query once an efficient label-only membership inference attack             | Strategies to craft query examples to reduce the required number of queries   | Image        | Predictive-Model | Membership-Inference | Empirical        | <a href="https://github.com/WuYutong/YouOnlyQueryOnce">https://github.com/WuYutong/YouOnlyQueryOnce</a>   |          | Johan Östman, AI Sweden                 |

| OpenDP Privacy Attacks & Auditing Working Group                                     |   |   |  |  |                  |                      |                      |               |   |          |   |
|---|---|---|--|--|------------------|----------------------|----------------------|---------------|---|----------|---|
| Privacy Attacks Repository  |   |   |  |  |                  |                      |                      |               |   |          |   |
| URL   | BitTex (Please add a bittex entry for this paper to facilitate writing our summary document)  | Authors   | Title  | Short Description  | Type of Data     | Type of Release      | Threat Model         | Research Type | Links to Artifacts  | Comments | Submitter (your name, affiliation)        |
| <a href="https://arxiv.org/abs/2021.10.0587">https://arxiv.org/abs/2021.10.0587</a> | @article{fowl2021robbing, title={Robbing the Fed: Directly Obtaining Private Data in Federated Learning with Modified Models}, author={Fowl, Weiping, Wu, Cezayirli, M. Goldblum, and Goldstein}, journal={arXiv preprint arXiv:2110.10587}, year={2021}}   | L Fowl, L Geiping, W Czaia, M Goldblum, Goldstein   | Robbing the Fed: Directly Obtaining Private Data in Federated Learning with Modified Models                            | Malicious adversary to perform exact reconstruction of the training data in FL   | Image            | Predictive-Model     | Reconstruction       | Empirical     | <a href="https://github.com/ImperialCollegeLondon/robbing-the-fed">https://github.com/ImperialCollegeLondon/robbing-the-fed</a>             |          | Dmitri Ushyn, TUM/Imperial College London |
| <a href="https://arxiv.org/abs/2201.12675">https://arxiv.org/abs/2201.12675</a>     | @article{fowl2022deceptions, title={Deceptions: Corrupted Transformers Breach Privacy in Federated Learning for Language Models}, author={Fowl, Liam and Geiping, Jonas and Reich, Steven and Wu, Yuxin and Czaia, Wojtek and Goldblum, Michal and Goldstein, Tom}, journal={arXiv preprint arXiv:2201.12675}, year={2022}}   | Liam Fowl, Jonas Geiping, Steven Reich, Yuxin Wu, Wojtek Czaia, Michal Goldblum, Tom Goldstein  | Deceptions: Corrupted Transformers Breach Privacy in Federated Learning for Language Models                            | Model inversion attacks using a malicious attacker in transformer-based FL settings  | Text             | Predictive-Model     | Reconstruction       | Empirical     | <a href="https://github.com/ImperialCollegeLondon/deceptions">https://github.com/ImperialCollegeLondon/deceptions</a>                       |          | Dmitri Ushyn, TUM/Imperial College London |
| <a href="https://arxiv.org/abs/2202.00580">https://arxiv.org/abs/2202.00580</a>     | @article{win2022shing, title={Fishing for user data in large-batch federated learning via gradient magnification}, author={Win, Yuxin and Geiping, Jonas and Fowl, Liam and Goldblum, Michal and Goldstein, Tom}, journal={arXiv preprint arXiv:2202.00580}, year={2022}}   | Yuxin Win, Jonas Geiping, Liam Fowl, Michal Goldblum, Tom Goldstein   | Fishing for User Data in Large-Batch Federated Learning via Gradient Magnification                                     | Malicious adversary to perform exact reconstruction of the training data in FL, this time with (almost) arbitrarily large batch sizes  | Image            | Predictive-Model     | Reconstruction       | Empirical     | <a href="https://github.com/ImperialCollegeLondon/shing">https://github.com/ImperialCollegeLondon/shing</a>                                 |          | Dmitri Ushyn, TUM/Imperial College London |
| <a href="https://arxiv.org/abs/2111.02918">https://arxiv.org/abs/2111.02918</a>     | @inproceedings{boenisch2023curious, title={When the curious abandon honesty: Federated learning is not private}, author={Boenisch, Franziska and Dziedzic, Adam and Schuster, Roni and Shamsabadi, Ali Shahin and Shumailov, Ilya and Papernot, Nicolas}, booktitle={2023 IEEE European Symposium on Security and Privacy (EuroS&P)}, pages={175–189}, year={2023}, organization={IEEE}}  | Franziska Boenisch, Adam Dziedzic, Roni Schuster, Ali Shahin Shamsabadi, Ilya Shumailov, Nicolas Papernot   | When the Curious Abandon Honesty: Federated Learning is Not Private  | Malicious FL input reconstruction using trap weights (model modification)  | Image            | Predictive-Model     | Reconstruction       | Empirical     | <a href="https://github.com/ImperialCollegeLondon/curious">https://github.com/ImperialCollegeLondon/curious</a>                             |          | Dmitri Ushyn, TUM/Imperial College London |
| <a href="https://arxiv.org/abs/2202.00580">https://arxiv.org/abs/2202.00580</a>     | @inproceedings{lyn2022see, title={See through gradients: Image batch recovery via gradient inversion}, author={Yin, Hongyu and Malliyil, Arun and Vahdat, Arash and Alvarez, Jose M and Kautz, Jan and Madhava, Pavlo}, booktitle={Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition}, pages={16337–16346}, year={2022}}  | Hongyu Yin, Arun Malliyil, Arash Vahdat, Jose M Alvarez, Jan Kautz, Pavlo Molchanov   | See Through Gradients: Image Batch Recovery via Gradient Inversion   | First demonstration of large batch-size model inversion in FL  | Image            | Predictive-Model     | Reconstruction       | Empirical     | <a href="https://github.com/ImperialCollegeLondon/see-through-gradients">https://github.com/ImperialCollegeLondon/see-through-gradients</a> |          | Dmitri Ushyn, TUM/Imperial College London |
| <a href="https://arxiv.org/abs/2202.00580">https://arxiv.org/abs/2202.00580</a>     | @inproceedings{karlyappa2022cocktail, title={Cocktail party attack: Breaking aggregation-based privacy in federated learning using independent component analysis}, author={Karlyappa, Sarinay and Guo, Chuan and Maeng, Kwang and Wang, Wenjie and Gu, Edward and Qureshi, Mounuiddin K and Lee, Hsiun-Hsin S}, booktitle={International Conference on Machine Learning}, pages={15884–15899}, year={2023}, organization={PMLR}}   | Sarinay Karlyappa, Chuan Guo, Kwang Maeng, Wenjie Wang, G. Edward Suh, Mounuiddin K Qureshi, Hsiun-Hsin S Lee   | Cocktail Party Attack: Breaking Aggregation-Based Privacy in Federated Learning Using Independent Component Analysis   | Theoretical attempts to prevent (and attack) secure aggregation in FL  | Image            | Predictive-Model     | Reconstruction       | Theoretical   |   |          | Dmitri Ushyn, TUM/Imperial College London |
| <a href="https://arxiv.org/abs/2202.00580">https://arxiv.org/abs/2202.00580</a>     | @article{ushyn2022beyond, title={Beyond gradients: Exploiting adversarial priors in model inversion attacks}, author={Ushyn, Dmitri and Ruecker, Daniel and Kassis, Georgios}, journal={ACM Transactions on Privacy and Security}, volume={26}, number={3}, pages={1–30}, year={2023}, publisher={ACM New York, NY}}  | Dmitri Ushyn, Daniel Ruecker, Georgios Kassis   | Beyond Gradients: Exploiting Adversarial Priors in Model Inversion Attacks   | More successful model inversion using HBC adversary with the knowledge of context and style of the training data   | Image            | Predictive-Model     | Reconstruction       | Empirical     |   |          | Dmitri Ushyn, TUM/Imperial College London |
| <a href="https://arxiv.org/abs/2404.02936">https://arxiv.org/abs/2404.02936</a>     | @article{zhang2024min, title={Min-KL++: Improved Baseline for Detecting Pre-Training Data from Large Language Models}, author={Zhang, Jingyao and Sun, Jingwei and Yeats, Eric and Duayang, Yang and Kuo, Martin and Zhang, Jimmy and Hsu, Wei-Chih}, journal={arXiv preprint arXiv:2404.02936}, year={2024}}   | Jingyao Zhang et al.  | Min-KL++: Improved Baseline for Detecting Pre-Training Data from Large Language Models                                 | Improved the Min-KL attack by normalizing the log prob of tokens   | Text             | Generative-Model     | Membership-Inference | Applications  | <a href="https://github.com/ImperialCollegeLondon/min-kl">https://github.com/ImperialCollegeLondon/min-kl</a>                               |          | Hamid Mofazzafi, Oracle Labs              |
| <a href="https://arxiv.org/abs/2019.06670">https://arxiv.org/abs/2019.06670</a>     | @inproceedings{hale2019monte, title={Monte carlo and reconstruction membership inference attacks against generative models}, author={Hale, Benjamin and Sun, Jingwei and Martin, Daniel}, booktitle={PACIS}, year={2019}}   | Benjamin Hale, Benjamin Sun, Jingwei and Daniel Martin  | Monte Carlo and Reconstruction Membership Inference Attacks against Generative Models                                  | Present two information leakage attacks that outperform previous work on membership inference against generative models.   | Image            | Generative-Model     | Membership-Inference | Empirical     |   |          | Georgi Ganev, UCL                         |
| <a href="https://arxiv.org/abs/2111.05114">https://arxiv.org/abs/2111.05114</a>     | @article{ganev2021inadequacy, title={On the Inadequacy of Similarity-based Privacy Metrics: Reconstruction Attacks against “Truly Anonymous Synthetic Data”}, author={Ganev, Georg and De Cristofaro, Emiliano}, journal={arXiv:2111.05114}, year={2021}}   | Georgi Ganev, Emiliano De Cristofaro  | On the Inadequacy of Similarity-based Privacy Metrics: Reconstruction Attacks against “Truly Anonymous Synthetic Data” | Revises the privacy metrics offered by leading companies in this space and sheds light on a few critical flaws in reasoning about privacy entirely via empirical evaluations. Presents a reconstruction attack, ReconSyn, which successfully recovers at least 75% of the low-density train records (or outliers) with only black-box access to a single fixed generative model and the privacy metrics. | Tabular          | Generative-Model     | Reconstruction       |               |   |          | Georgi Ganev, UCL                         |
| <a href="https://arxiv.org/abs/2107.03114">https://arxiv.org/abs/2107.03114</a>     | @inproceedings{opriescu2021on, title={On utility and privacy in synthetic genomic data}, author={Opriescu, Brindana and Ganev, Georgi and De Cristofaro, Emiliano}, booktitle={NDS3}, year={2021}}  | Brindana Opriescu, Georgi Ganev, Emiliano De Cristofaro   | On utility and privacy in synthetic genomic data   | Provides the first evaluation of both utility and privacy protection of a state-of-the-art models for generating synthetic genomic data. The experiments show that no single approach to generate synthetic genomic data yields both high utility and strong privacy across the board.   | Tabular          | Generative-Model     | Membership-Inference | Empirical     |   |          | Georgi Ganev, UCL                         |
| <a href="https://arxiv.org/abs/2202.03098">https://arxiv.org/abs/2202.03098</a>     | @article{andrew2022one, title={One-shot empirical privacy estimation for federated learning}, author={Andrew, Galen and Karouz, Peter and Oh, Seungwon and Ojha, Abhinav and McHale, Brendan and Suriyakumar, Vinith}, journal={arXiv preprint arXiv:2202.03098}, year={2022}}  | Galen Andrew, Peter Karouz, Seungwon Oh, Abhinav Ojha, H. Brendan McHale, Vinith M. Suriyakumar   | One-shot Empirical Privacy Estimation for Federated Learning   | Presents an approach for performing a strong white-box attack for measuring the DP epsilon of ML training algos in 1 training run.   | Generative-Model | Membership-Inference |                      |               |   |          | Peter Karouz, Google                      |
| <a href="https://arxiv.org/abs/2018.08027">https://arxiv.org/abs/2018.08027</a>     | @INPROCEEDINGS{8429311, author={Yeom, Samuel and Giacomelli, Irene and Fredrikson, Matt and Jha, Somesh}, booktitle={2018 IEEE 31st Computer Security Foundations Symposium (CSF)}, title={Privacy Risk in Machine Learning: Analyzing the Connection to Overfitting}, year={2018}, volume={1}, number={1}, pages={268–282}, keywords={Privacy;Machine learning algorithms; Data models;Training data;Machine learning;Training; privacy;machine-learning;inference-attacks}, doi={10.1109/CSF.2018.00027}} | Samuel Yeom; Irene Giacomelli; Matt Fredrikson; Somesh Jha  | Privacy Risk in Machine Learning: Analyzing the Connection to Overfitting  | MIA attack based on final loss for LLMs  | Text             | Generative-Model     | Membership-Inference | Empirical     |   |          | Daniel Flato, University of Washington    |
| <a href="https://arxiv.org/abs/2203.03929">https://arxiv.org/abs/2203.03929</a>     | @article{mishra2022quantifying, title={Quantifying privacy risks of masked language models using membership inference attacks}, author={Mishra, Fatemeh and Goyal, Kartik and Ushyn, Dmitri and Berg-Kirkpatrick, Taylor and Shari, Reza}, journal={arXiv preprint arXiv:2203.03929}, year={2022}}  | Fatemehdad Mreshghallah, Kartik Goyal, Arhit Ushyn, Taylor Berg-Kirkpatrick, Reza Shari   | Quantifying Privacy Risks of Masked Language Models Using Membership Inference Attacks                                 | Membership inference attack based on likelihood ratio hypothesis testing that involves an additional reference MLM to more accurately quantify the privacy risks of memorization in MLMs.  | Text             | Generative-Model     | Membership-Inference | Applications  |   |          | Daniel Flato, University of Washington    |
| <a href="https://arxiv.org/abs/2003.10920">https://arxiv.org/abs/2003.10920</a>     | @article{naor2020scalable, title={Scalable extraction of training data from production language models}, author={Naor, Mital and Carlini, Nicholas and Hayes, Jonathan and Jagielski, Matthew and Cooper, Daphne and Ippolito, Christopher A. and Wallace, Eric and Tramèr, Florian and Lee, Katherine}, journal={arXiv:2003.10920}, year={2020}}   | Mital Naor, Nicholas Carlini, Jonathan Hayes, Matthew Jagielski, A. Feder Cooper, Daphne Ippolito, Christopher A. Choquette-Choo, Eric Wallace, Florian Tramèr, Katherine Lee | Scalable extraction of training data from production language models   | Studies extractable memorization: training data that an adversary can efficiently extract by querying a machine learning model without prior knowledge of the training dataset. Shows an adversary can extract gigabytes of training data from open-source language models like Pythia or GPT-Neo, semi-open models like LLaMA or Falcon, and closed models like ChatGPT.                                | Text             | Generative-Model     | Data-Extraction      | Empirical     |   |          | Georgi Ganev, UCL                         |
| <a href="https://arxiv.org/abs/2101.04543">https://arxiv.org/abs/2101.04543</a>     | @inproceedings{hsu2021adversary, title={Adversary instantiation: lower bounds for differentially private machine learning}, author={Hsu, Mital and Song, Shuang and Thakurta, Abhradeep and Papernot, Nicolas and Carlini, Nicholas}, booktitle={IEEE S&P}, year={2021}}  | Mital Nasir, Shuang Song, Abhradeep Thakurta, Nicolas Papernot, Nicholas Carlini  | Adversary instantiation: lower bounds for differentially private machine learning                                      | Instantiates a hypothetical adversary in order to establish lower bounds on the probability that the distinguishing game can be won. Uses this adversary to evaluate the importance of the adversary capabilities allowed in the privacy analysis of DP training algorithms.   | Image            | Predictive-Model     | Membership-Inference | Empirical     |   |          | Georgi Ganev, UCL                         |

| OpenDP Privacy Attacks & Auditing Working Group   |   |   |   |   |              |                     |                      |               |   |          |  |
|---|---|---|---|---|--------------|---------------------|----------------------|---------------|---|----------|--|
| Privacy Attacks Repository  |   |   |   |   |              |                     |                      |               |   |          |  |
| URL   | BitTex (Please add a bitTex entry for this paper to facilitate writing our summary document)  | Authors   | Title   | Short Description   | Type of Data | Type of Release     | Threat Model         | Research Type | Links to Artifacts  | Comments | Submitter (your name, affiliation)           |
| <a href="https://arxiv.org/abs/2101.03634">https://arxiv.org/abs/2101.03634</a>   | @article{carlini2021stealing, title={Stealing part of a production language model}, author={Carlini, Nicholas and Paleka, Daniel and Dvijotham, Krishnamurthy D and Steinko, Thomas and Hayes, Jonathan and Cooper, A Feder Cooper, Katherine Lee, Matthew Jagielski, Milad Nasr, Arthur Conroy, Eric Wallace, David Rolnick, Florian Tramèr} year={2021} } | Nicholas Carlini, Daniel Paleka, Krishnamurthy (D) Dvijotham, Thomas Steinko, Jonathan Hayes, A. Feder Cooper, Katherine Lee, Matthew Jagielski, Milad Nasr, Arthur Conroy, Eric Wallace, David Rolnick, Florian Tramèr | Stealing Part of a Production Language Model  | Recovers the embedding projection layer (ignoring the rest of the transformer model) given typical black-box (API) access   | Text         | Generative-Model    | Data-Extraction      | Applications  |   |          | Daniil Filenko, University of Washington     |
| <a href="https://arxiv.org/pdf/2011.07018v4.pdf">https://arxiv.org/pdf/2011.07018v4.pdf</a>                                       | @inproceedings{stadler2022synthetic, title={Synthetic data - anonymisation groundhog day}, author={Stadler, Theresa and Opiranis, Bristen and Troncoso, Carmela}, booktitle={USENIX Security Symposium (USENIX Security 22)}, pages={1451–1468}, year={2022} }  | Stadler et al.  | Synthetic Data - Anonymisation Groundhog Day  | Black-box attack on synthetic data and quantitative evaluation of privacy of synthetic data   | Tabular      | Generative-Model    | Membership-Inference | Empirical     |   |          | Yves-Alexandre de Montjoye, Imperial College |
| <a href="https://arxiv.org/abs/2102.07956">https://arxiv.org/abs/2102.07956</a>   | @inproceedings{nasr2021tight, title={Tight Auditing of Differentially Private Machine Learning}, author={Milad Nasr, Jamie Hayes and Thomas Steinko and Borja Balle and Florian Tramèr (U) and Matthew Jagielski and Nicholas Carlini and Andreas Terzis}, booktitle={USENIX Security}, year = {2021} }   | Milad Nasr, Jamie Hayes, Thomas Steinko, Borja Balle, Florian Tramèr, Matthew Jagielski, Nicholas Carlini, Andreas Terzis   | Tight Auditing of Differentially Private Machine Learning   | Designs an improved auditing scheme that yields tight privacy estimates for natural (not adversarially crafted) datasets -- if the adversary can see all model updates during training. Moreover, the auditing scheme requires only two training runs (instead of thousands) to produce tight privacy estimates, by adapting recent advances in tight composition theorems for differential privacy.                              | Image        | Predictive-Model    | Membership-Inference | Empirical     |   |          | Georgi Ganev, UCL                            |
| <a href="https://arxiv.org/abs/2011.11455v2">https://arxiv.org/abs/2011.11455v2</a>   | @article{garfinkel2019understanding, title={Understanding database reconstruction attacks on public data}, author={Garfinkel, Simon and Abowd, John M and Martindale, Christian}, journal={ACM Queue}, year={2019} }  | Simon Garfinkel, John M Abowd, Christian Martindale   | Understanding database reconstruction attacks on public data  | Database reconstruction attacks can be performed by using published statistical tables to create a set of mathematical constraints and then solving the resulting set of simultaneous equations. Shows how such an attack can be addressed by adding noise to published tabulations, so that the reconstruction no longer results in the original data. This has implications for the 2020 Census.                                | Tabular      | Linear-Queries      | Reconstruction       | Empirical     |   |          | Georgi Ganev, UCL                            |
| <a href="https://arxiv.org/abs/2107.11701">https://arxiv.org/abs/2107.11701</a>   | @article{gulpin2022synthetic, title={Synthetic is all you need: removing the auxiliary data assumption for membership inference attacks against synthetic data}, author={Florent Gulpin, Matthieu Meun, Ana-Maria Cretu, Yves-Alexandre de Montjoye, Jean-Luc Arlot}, year={2022} }   | Florent Gulpin, Matthieu Meun, Ana-Maria Cretu, Yves-Alexandre de Montjoye  | Synthetic is all you need: removing the auxiliary data assumption for membership inference attacks against synthetic data | Develops new MIs performed using only the synthetic data in three different scenarios: (S1) Black-box access to the generator; (S2) only access to the released synthetic dataset and (S3) a theoretical setup as upper bound for the attack performance.   | Tabular      | Generative-Model    | Membership-Inference | Empirical     |   |          | Georgi Ganev, UCL                            |
| <a href="https://arxiv.org/abs/2106.07758">https://arxiv.org/abs/2106.07758</a>   | @inproceedings{ham2022reconstructing, title={Reconstructing training data from trained neural networks}, author={Ham, Niv and Vardi, Gal and Yehudai, Gilad and Eran, Michael and Shamir, Ohad}, booktitle={NeurIPS}, year = {2022} }   | Niv Ham, Gal Vardi, Gilad Yehudai, Ohad Shamir, Michael Eran  | Reconstructing Training Data from Trained Neural Networks   | Shows that in some cases a significant fraction of the training data can in fact be reconstructed from the parameters of a trained neural network classifier. Proposes a novel reconstruction scheme that stems from recent theoretical results about the implicit bias in training neural networks with gradient-based methods.  | Image        | Predictive-Model    | Reconstruction       | Empirical     | <a href="https://github.com/ohadshamir/reconstruction">https://github.com/ohadshamir/reconstruction</a>   |          | Georgi Ganev, UCL                            |
| <a href="https://arxiv.org/abs/2101.04845">https://arxiv.org/abs/2101.04845</a>   | @inproceedings{balle2022reconstructing, title={Reconstructing training data with informed adversaries}, author={Balle, Borja and Chenab, Giovanni and Hayes, Jamie}, booktitle={IEEE S&P}, year={2022} }  | Borja Balle, Giovanni Chenab, Jamie Hayes   | Reconstructing training data with informed adversaries  | Studies how given access to a machine learning model an adversary can reconstruct the training data from the lens of a powerful, informed adversary who knows all the training data points except one.  | Image        | Predictive-Model    | Reconstruction       | Empirical     |   |          | Georgi Ganev, UCL                            |
| <a href="https://www.usenix.org/conference/ussm21/paper/2101.04845">https://www.usenix.org/conference/ussm21/paper/2101.04845</a> | @inproceedings{islem2020updates, title={Updates-leak: Data set inference and reconstruction attacks in online learning}, author={Islem, Ahmed Mohamed Gamal and Bhatnagar, Arvind and Backes, Michael and Fritz, Mario and Zhang, Yang}, booktitle={USENIX Security}, year={2020} }   | Ahmed Salem, Apratim Bhatnagar, Michael Backes, Mario Fritz, Yang Zhang   | Updates-leak: Data set inference and reconstruction attacks in online learning  | Investigate whether the change in the output of a black-box ML model before and after being updated can leak information of the dataset used to perform the update, namely the updating set. Proposes four attacks following an encoder-decoder formulation, which allows inferring diverse information of the updating set.  | Image        | Predictive-Model    | Reconstruction       | Empirical     |   |          | Georgi Ganev, UCL                            |
| <a href="https://arxiv.org/abs/2105.16160">https://arxiv.org/abs/2105.16160</a>   | @inproceedings{jokna2023group, title={Group and Attack: Auditing Differential Privacy}, author={Jokna, Johan and Paradi, Anouk and Dimitrov, Dimitar and Vechev, Martin}, year={2023} }   | Johan Jokna, Anouk Paradi, Dimitar I Dimitrov, Martin Vechev  | Group and Attack: Auditing Differential Privacy   | Present a novel method to efficiently discover $k$ differentially private violations based on the key insight that many $k$ -tuples can be grouped as they occur in the same algorithm. Crucially, the method is orthogonal to existing approaches and, when combined, results in a faster and more precise violation search.   | Image        | Information Leakage | Information Leakage  | Empirical     | <a href="https://github.com/paradi/auditdiffp">https://github.com/paradi/auditdiffp</a>                   |          | Georgi Ganev, UCL                            |
| <a href="https://arxiv.org/abs/2107.08923">https://arxiv.org/abs/2107.08923</a>   | @inproceedings{carlini2019secret, title={The secret sharer: Evaluating and testing unintended memorization in neural networks}, author={Carlini, Nicholas and Jagielski, Matthew and Erlingsson, (Ulfarar and Kos, Jernge and Song, Dawn)}, booktitle={28th USENIX security symposium USENIX security 18}, pages={267–284}, year={2019} }                   | Nicholas Carlini, Chang Liu, Ulfar Erlingsson, Jernge Kos, Dawn Song  | The Secret Sharer: Evaluating and Testing Unintended Memorization in Neural Networks                                      | Describes a testing methodology for quantitatively assessing the risk that one or more parameters in a model are unintentionally memorized by generative sequence models.   | Text         | Generative-Model    | Data-Extraction      | Empirical     |   |          | Georgi Ganev, UCL                            |
| <a href="https://arxiv.org/abs/2109.12711">https://arxiv.org/abs/2109.12711</a>   | @article{tramer2022debugging, title={Debugging Differential Privacy: A Case Study for Privacy Auditing}, author={Tramer, Florian and Terzis, Andreas and Steinko, Thomas and Song, Shuang and Jagielski, Matthew and Carlini, Nicholas}, journal={arXiv:2202.12119}, year={2022} }  | Florian Tramèr, Andreas Terzis, Thomas Steinko, Shuang Song, Matthew Jagielski, Nicholas Carlini  | Debugging Differential Privacy: A Case Study for Privacy Auditing   | Inspired by recent advances in auditing which have been used for estimating lower bounds on differentially private algorithms, shows that auditing can also be used to find in (purportedly) differentially private schemes.  | Image        | Predictive-Model    | Information Leakage  | Empirical     |   |          | Georgi Ganev, UCL                            |
| <a href="https://arxiv.org/abs/2106.05193">https://arxiv.org/abs/2106.05193</a>   | @inproceedings{zanella2023bayesian, title={Bayesian estimation of differential privacy}, author={Zanella, B\u00e9n\u00e9dikt, Lukas Wutschitz, Lukas and Topke, Shrut and Salem, Ahmed and R\u00f9hle, Victor and Pavoni, Andrew and Nazer, Mohammad and K\u00f6pf, Boris and Jones, Daniel}, booktitle={ICML}, year={2023} }                               | Santiago Zanella-B\u00e9n\u00e9dikt, Lukas Wutschitz, Shrut Topke, Ahmed Salem, Victor R\u00f9hle, Andrew Pavoni, Mohammad Nazer, Boris K\u00f6pf, Daniel Jones   | Bayesian Estimation of Differential Privacy   | Proposes a novel Bayesian method that greatly reduces sample size, and adapts and validates a heuristic to draw more than one sample per trained model. The Bayesian method exploits the hypothesis testing interpretation of differential privacy to obtain a posterior for $\epsilon$ (not just a confidence interval) from the joint posterior of the false positive and false negative rates of membership-inference attacks. | Image        | Predictive-Model    | Information Leakage  | Empirical     |   |          | Georgi Ganev, UCL                            |
| <a href="https://arxiv.org/abs/2103.05111">https://arxiv.org/abs/2103.05111</a>   | @inproceedings{bichsel2021dp, title={DP-Sniper: Black-Box Discovery of Differential Privacy Violations using Classifier}, author={Bichsel, Benjamin and Samet, Samuel and Bogunovic, Ija and Vechev, Martin}, booktitle={IEEE S&P}, year={2021} }   | Benjamin Bichsel, Samuel Steffen, Ija Bogunovic, Martin Vechev  | DP-Sniper: Black-Box Discovery of Differential Privacy Violations using Classifier  | Present DP-Sniper, a practical black-box method that automatically finds violations of differential privacy.  | Tabular      | Information Leakage | Information Leakage  | Empirical     |   |          | Georgi Ganev, UCL                            |
| <a href="https://arxiv.org/pdf/2110.09366v4.pdf">https://arxiv.org/pdf/2110.09366v4.pdf</a>                                       | @article{kandpal2022user, title={User inference attacks on large language models}, author={Kandpal, Nikhil and Pillutla, Krishna and Oprea, Alina and Kairouz, Peter and Choquette-Choo, Christopher A and Xu, Zheng}, journal={arXiv preprint arXiv:2310.09366}, year={2023} }   | Nikhil Kandpal, Krishna Pillutla, Alina Oprea, Peter Kairouz, Christopher A. Choquette-Choo, Zheng Xu   | User Inference Attacks on Large Language Models   | Presents an attack for inferring the preexistence of a user in the fine-tuning set of an LLM. The adversary is not assumed to know all the fine-tuning examples of a user -- only a subset (including some examples that weren't used even if the user participated in the fine-tuning stage).  | Image        | Generative-Model    | Membership-Inference | Empirical     |   |          | Peter Kairouz, Google                        |
| <a href="https://arxiv.org/abs/2406.11544">https://arxiv.org/abs/2406.11544</a>   | @inproceedings{svr2024do, title={Do Parameters Reveal More than Loss for Membership Inference?}, author={Anshuman Suri and Xiao Zhang and David Evans}, booktitle={Workshop on High-dimensional Learning Dynamics (HLD), ICML}, year = {2024}, url = {https://arxiv.org/abs/2406.11544} }   | Anshuman Suri, Xiao Zhang, David Evans  | Do Parameters Reveal More than Loss for Membership Inference?   | The paper shows how prior claims about black-box access sufficing for optimal membership inference do not hold for most useful settings such as SGD   | Tabular      | Predictive-Model    | Membership-Inference | Theoretical   | <a href="https://github.com/anshumansuri/24hld">https://github.com/anshumansuri/24hld</a>                 |          | Anshuman Suri, UVA                           |
| <a href="https://arxiv.org/abs/2402.10001">https://arxiv.org/abs/2402.10001</a>   | @INPROCEEDINGS{EMINx2024a, author = {(El Menni, Abdelhak and ElFerc, Edwige and Bellat, Aur\u00e9lien)}, title = {(Privacy Attacks in (D) centralized Learning)}, booktitle = {(ICML)}, year = {2024} }   | Abdelhak El Menni, Edwige Cyffers, Aur\u00e9lien Bellat   | Privacy Attacks in Decentralized Learning   | The paper designs data reconstruction attacks against Decentralized SGD (where nodes in a communication graph alternate between local gradient steps and averaging steps with their neighbors). They show it is possible for a small subset of honest but curious attacker nodes to reconstruct the data from even distant nodes in the graph.  | Image        | Predictive-Model    | Reconstruction       | Empirical     | <a href="https://github.com/AbdelhakElMenni/eminx2024a">https://github.com/AbdelhakElMenni/eminx2024a</a> |          | Aur\u00e9lien Bellat, Inria                  |

OpenDP Privacy Attacks & Auditing Working Group

| URL   | BitTex (Please add a bittext entry for this paper to facilitate writing our summary document)   | Authors   | Title   | Short Description   | Type of Data   | Type of Release  | Threat Model         | Research Type | Links to Artifacts  | Comments | Submitter (your name, affiliation) |
|---|---|---|---|---|----------------|------------------|----------------------|---------------|---|----------|------------------------------------|
| <a href="https://arxiv.org/abs/2106.03458">https://arxiv.org/abs/2106.03458</a> | @inproceedings{mlak2021labeldp,<br>author = {Malik, Emanel, Mani and Mironov, Ilya and Prasad, Karthik and Shilov, Igor and Tramer, Florian},<br>booktitle = {Advances in Neural Information Processing Systems},<br>editor = {M. Ranzato and A. Beygelzimer and Y. Dauphin and P.S. Liang and J. Wortman Vaughan},<br>pages = {6934–6945},<br>publisher = {Curran Associates, Inc.},<br>title = {Antipodes of Label Differential Privacy: (PATE) and (ALIB)},<br>url = {https://proceedings.neurips.cc/paper_files/paper/2021/file/97cc27608480aa23569a511e638ca74f-Paper.pdf},<br>volume = {34},<br>year = {2021}}} | Mani Malik, Ilya Mironov, Karthik Prasad, Igor Shilov, Florian Tramer   | Antipodes of Label Differential Privacy: PATE and ALIB  | Label-inference attack against two LabelDP mechanisms for image classification  | Image          | Predictive-Model | Attribute inference  | Empirical     | <a href="https://github.com/curran/labeldp">https://github.com/curran/labeldp</a> , <a href="https://arxiv.org/abs/2106.03458">https://arxiv.org/abs/2106.03458</a> , <a href="https://www.robots.ox.ac.uk/~sreeram/papers/labeldp/paper.pdf">https://www.robots.ox.ac.uk/~sreeram/papers/labeldp/paper.pdf</a> |          | Ilya Mironov, Meta                 |
| <a href="https://arxiv.org/abs/2112.11089">https://arxiv.org/abs/2112.11089</a> | @INPROCEEDINGS {islem2023-sok,<br>author = {A. Salem and G. Cherubin and D. Evans and B. Kopf and A. Paverd and A. Suri and S. Tople and S. Zanella Baguein},<br>booktitle = {2023 IEEE Symposium on Security and Privacy (S&P)},<br>title = {[S&K] Let the Privacy Games Begin! A Unified Treatment of Data Inference Privacy in Machine Learning},<br>year = {2023},<br>pages = {327–345},<br>doi = {10.1109/SP4215.2023.10179281},<br>url = {https://doi.ieeecomputersociety.org/10.1109/SP4215.2023.10179281},<br>publisher = {IEEE Computer Society},<br>address = {Los Alamitos, CA, USA},<br>month = {May}}}   | Ahmed Salem, Giovanni Cherubin, David Evans, Boris Kopf, Andrew Paverd, Anshuman Suri, Shruti Tople, Santiago Zanella-Baguein   | S&K: Let the Privacy Games Begin! A Unified Treatment of Data Inference Privacy in Machine Learning | Systematization of Knowledge of data inference attacks using a privacy-game framework.  |                |                  |                      | Theoretical   |   |          | Ilya Mironov, Meta                 |
| <a href="https://arxiv.org/abs/2111.08464">https://arxiv.org/abs/2111.08464</a> | @inproceedings{watson2022-diffculty,<br>author = {Lauren Watson and Chuan Guo and Graham Cormode and Alexandros Sablayrolles},<br>title = {(On the Importance of Difficulty Calibration in Membership Inference Attacks)},<br>booktitle = {The Tenth International Conference on Learning Representations, (ICLR) 2022, Virtual Event, April 25–29, 2022},<br>publisher = {OpenReview.net},<br>year = {2022},<br>url = {https://openreview.net/forum?id=3e1d0TwQj}}}  | Lauren Watson, Chuan Guo, Graham Cormode, Alex Sablayrolles   | On the Importance of Difficulty Calibration in Membership Inference Attacks                         | Improvement in membership inference attacks by taking into account the difficulty of correct classification   | Image, tabular | Predictive-Model | Membership-inference | Empirical     | <a href="https://github.com/curran/labeldp">https://github.com/curran/labeldp</a> , <a href="https://arxiv.org/abs/2111.08464">https://arxiv.org/abs/2111.08464</a> , <a href="https://openreview.net/forum?id=3e1d0TwQj">https://openreview.net/forum?id=3e1d0TwQj</a>   |          | Ilya Mironov, Meta                 |
| <a href="https://arxiv.org/abs/2108.03458">https://arxiv.org/abs/2108.03458</a> | @INPROCEEDINGS {papernot2018-s&k,<br>author = {Papernot, Nicolas and McDaniel, Patrick and Sinha, Arunesh and Weisman, Michael P.},<br>booktitle = {2018 IEEE European Symposium on Security and Privacy (EuroSec)},<br>title = {[S&K: Security and Privacy in Machine Learning]},<br>year = {2018},<br>volume = {1},<br>number = {1},<br>pages = {399–414}}}   | Nicolas Papernot, Patrick McDaniel, Arunesh Sinha, Michael P. Weisman   | S&K: Security and Privacy in Machine Learning   | Systematization of Knowledge of attacks on privacy of ML.   |                |                  |                      | survey        |   |          | Ilya Mironov, Meta                 |
| <a href="https://arxiv.org/abs/2107.09171">https://arxiv.org/abs/2107.09171</a> | @ARTICLE {trux2021,<br>author = {Trux, Stacy and Liu, Ling and Gursoy, Mehmet Emre and Yu, Lei and Wei, Wenqi},<br>journal = {IEEE Transactions on Services Computing},<br>title = {Demystifying Membership Inference Attacks in Machine Learning as a Service},<br>year = {2021},<br>volume = {14},<br>number = {6},<br>pages = {2073–2089},<br>doi = {10.1109/TSC.2019.2897554}}}   | Stacy Trux, Ling Liu, Mehmet Emre Gursoy, Lei Yu, and Wenqi Wei   | Demystifying Membership Inference Attacks in Machine Learning as a Service                          | (1) MIA in the style of Shokri et al.'17 with shadow models with architecture different from that of the target model. (2) MIA with access to training gradients in the federated setting.  | Image, tabular | Predictive-Model | Membership-Inference | Empirical     |   |          | Ilya Mironov, Meta                 |
| <a href="https://arxiv.org/abs/2005.13702">https://arxiv.org/abs/2005.13702</a> | @inproceedings{Rezaei_2021_CVPR,<br>author = {Rezaei, Shahbaz and Liu, Xin},<br>title = {(On the Difficulty of Membership Inference Attacks)},<br>booktitle = {Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)},<br>month = {June},<br>year = {2021},<br>pages = {7892–7900}}}   | Shahbaz Rezaei, Xin Liu   | On the Difficulty of Membership Inference Attacks   | Critique of prior MIAs for their poor false positive rate.  | Image          | Predictive-Model | Membership-Inference | Empirical     | <a href="https://github.com/shahbazrezaei">https://github.com/shahbazrezaei</a> , <a href="https://arxiv.org/abs/2005.13702">https://arxiv.org/abs/2005.13702</a>   |          |                                    |
| <a href="https://arxiv.org/abs/2111.11396">https://arxiv.org/abs/2111.11396</a> | @misc {abowd20232010censusconfidentialityprotections,<br>title = {The 2010 Census Confidentiality Protections Failed, Here's How and Why},<br>author = {John M. Abowd and Tamara Adams and Robert Ashmead and David Danis and Sourya Dey and Simon L. Garfinkel and Nathan Goldschlag and Daniel Kifer and Philip Lederer and Ethan Lew and Scott Moore and Rolando A. Rodriguez and Romy N. Tadro and Lars Vilhuber},<br>year = {2023},<br>eprint = {2111.11396},<br>archivePrefix = {arXiv}}}   | John M. Abowd, Tamara Adams, Robert Ashmead, David Danis, Sourya Dey, Simon L. Garfinkel, Nathan Goldschlag, Daniel Kifer, Philip Lederer, Ethan Lew, Scott Moore, Rolando A. Rodriguez, Romy N. Tadro, Lars Vilhuber | The 2010 Census Confidentiality Protections Failed, Here's How and Why                              | The definitive analysis of the reconstruction-abetted reidentification attack on the 2010 Census data. Microdata were reconstructed from tabular summaries using a series of Integer Programs. An additional Integer Program assesses the solution variability. For 97 million of 308 million persons, the reconstruction is provably perfect. For perfectly reconstructed persons with physical characteristics within their geography, inferences about these characteristics are correct with probability 0.95, far in excess of correct inferences from purely statistical models. Inferences are also correct with similar probability for typical persons within their geography, but such inferences could have been made with a purely statistical model. | Tabular        | Linear-Queries   | Reconstruction       | Empirical     | <a href="https://github.com/abowd/2010censusrebuild">https://github.com/abowd/2010censusrebuild</a> , <a href="https://arxiv.org/abs/2111.11396">https://arxiv.org/abs/2111.11396</a>   |          | John Abowd, Cornell                |
|   |   |   |   |   | Tabular        |                  |                      |               |   |          |                                    |