

# A Practical Guide to DP Deep Learning

...Using the Pseudo Posterior Mechanism

**Rob Chew**<sup>1</sup>, Sandy Preiss<sup>1</sup>, Amanda Konet<sup>1</sup>,  
Matt Williams<sup>1</sup>, Elan Segarra<sup>2</sup>, David Oh<sup>2</sup>, Erin  
Boon<sup>2</sup>, Terrance Savitsky<sup>2</sup>

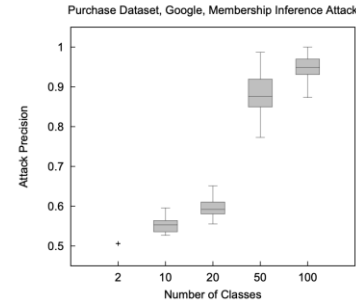
<sup>1</sup> RTI International

<sup>2</sup> U.S. Bureau of Labor Statistics

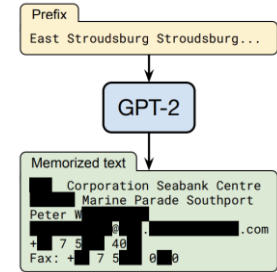


# Motivation

- Statistical agencies are increasingly interested in releasing internally developed autocoding models.
- However, releasing deep learning (DL) models can pose serious privacy risks.
- Differential privacy (DP) can help, but existing DP DL methods are difficult to apply in practice.
- Guidance for newer DP DL paradigms is largely absent.
- This talk presents practical guidance for diagnosing SWAG-PPM, a DP DL method developed with autocoding as a primary use case.



Shokri et al., 2017



Carlini et al., 2021

|   |    |
|---|----|
| Contents  |    |
| 1 Introduction  | 4  |
| 1.1 Prover of the Last Section                                      | 4  |
| 2 Differential Privacy: Definitions, Intuition and Properties       | 6  |
| 2.1 Differential Privacy: Naïve Laplace Mechanism                   | 6  |
| 2.2 Properties of DP  | 6  |
| 2.3 Alternative Stronger Definitions of DP                          | 6  |
| 2.4 Basic DP Mechanisms   | 9  |
| 3 DP-Resistant Models: Settings and Methods                         | 18 |
| 3.1 DP Settings: Exact Models and Release Boundaries                | 18 |
| 3.2 Release Agents  | 18 |
| 3.3 DP at the Data Level  | 18 |
| 3.4 DP at the Protocol Level: Privacy-Preserving Publishing         | 18 |
| 3.5 DP During The Training Phase: Preserving Only Labels (Label-DP) | 17 |
| 4 DP-Resistant: Protecting ML Training Data                         | 18 |
| 4.1 Review of DP-Resistant Methods                                  | 18 |
| 4.1.1 General Model: Naïve Laplace Mechanism                        | 18 |
| 4.1.2 Objective-Cost Modification Methods                           | 19 |
| 4.1.3 Gradient Noise Injection Techniques                           | 19 |
| 4.1.4 Information-Theoretic Methods for DP Training                 | 21 |
| 4.1.5 Generative Models for DP-Resistant                            | 21 |
| 4.2 DP-Resistant: Privacy Guarantees: Theory *                      | 25 |
| 4.3 Privacy Amplification via Sampling *                            | 25 |
| 4.4 Modifications for Low-Level DP Training                         | 25 |
| 4.5 Challenges with DP Training                                     | 26 |
| 5 Practicalities of DP-Resistant                                    | 26 |
| 5.1 Choosing the Right DP-Resistant                                 | 26 |
| 5.2 What is a Cost $c$ for an ML Model                              | 26 |
| 5.2.1 Cost Decomposition for Values for ML models                   | 26 |
| 5.2.2 Discussion and Assumptions                                    | 26 |
| 5.3 Calculating the Required Privacy Guarantees                     | 26 |
| 5.3.1 Calculating Privacy Guarantees for DP-Resistant               | 26 |
| 5.3.2 Calculating Privacy Guarantees for DP-Resistant               | 26 |
| 5.3.3 Reporting Privacy Guarantees for ML Models                    | 26 |
| 5.4 Hyperparameter Tuning   | 27 |
| 5.4.1 How to Tune the Hyperparameters for DP-Resistant              | 27 |
| 5.4.2 How Hyperparameter Tuning Can Improve $\epsilon$              | 27 |
| 5.4.3 Model Performance Considerations: Data Privacy                | 27 |
| 5.4.4 Model Performance Considerations: Model Quality               | 27 |
| 5.4.5 Design Choices Affecting Model Quality                        | 27 |
| 5.5 Model Quality   | 27 |
| 5.6 Model Quality   | 27 |
| 5.7 Frameworks and Libraries for DP                                 | 27 |
| 6 Conclusion  | 33 |

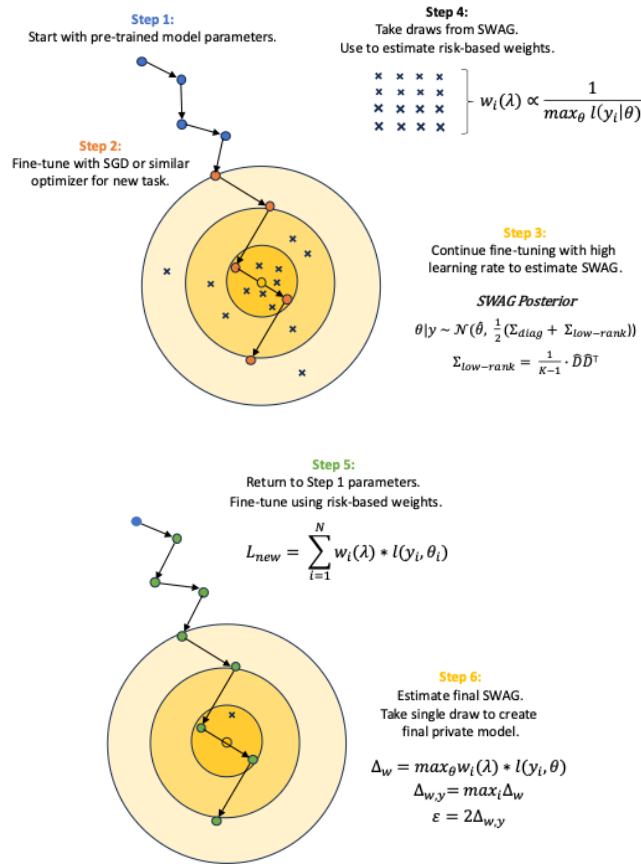
|  |    |
|--|----|
| A. DP-Training for non-differential models                       | 44 |
| A.1. Invariant algorithms  | 44 |
| A.2. Clustering algorithms                                       | 44 |
| B. Derivation of DP-Resistant user profiles                      | 71 |
| C. Example comparison of Approximate DP using accounting methods | 72 |
| C.1 DP-Resistant   | 72 |
| C.2 DP-Resistant   | 72 |
| C.3 DP-Resistant   | 72 |
| C.4 Derivation of DP-Resistant user profiles                     | 72 |
| D.1.1 DP-Resistant   | 72 |
| D.1.2 DP-Resistant   | 72 |
| D.1.3 DP-Resistant   | 72 |
| D.1.4 DP-Resistant   | 72 |
| D.1.5 DP-Resistant   | 72 |
| D.1.6 DP-Resistant   | 72 |
| D.1.7 DP-Resistant   | 72 |
| D.1.8 DP-Resistant   | 72 |
| D.1.9 DP-Resistant   | 72 |
| D.1.10 DP-Resistant  | 72 |
| D.1.11 DP-Resistant  | 72 |
| D.1.12 DP-Resistant  | 72 |
| D.1.13 DP-Resistant  | 72 |
| D.1.14 DP-Resistant  | 72 |
| D.1.15 DP-Resistant  | 72 |
| D.1.16 DP-Resistant  | 72 |
| D.1.17 DP-Resistant  | 72 |
| D.1.18 DP-Resistant  | 72 |
| D.1.19 DP-Resistant  | 72 |
| D.1.20 DP-Resistant  | 72 |
| D.1.21 DP-Resistant  | 72 |
| D.1.22 DP-Resistant  | 72 |
| D.1.23 DP-Resistant  | 72 |
| D.1.24 DP-Resistant  | 72 |
| D.1.25 DP-Resistant  | 72 |
| D.1.26 DP-Resistant  | 72 |
| D.1.27 DP-Resistant  | 72 |
| D.1.28 DP-Resistant  | 72 |
| D.1.29 DP-Resistant  | 72 |
| D.1.30 DP-Resistant  | 72 |
| D.1.31 DP-Resistant  | 72 |
| D.1.32 DP-Resistant  | 72 |
| D.1.33 DP-Resistant  | 72 |
| D.1.34 DP-Resistant  | 72 |
| D.1.35 DP-Resistant  | 72 |
| D.1.36 DP-Resistant  | 72 |
| D.1.37 DP-Resistant  | 72 |
| D.1.38 DP-Resistant  | 72 |
| D.1.39 DP-Resistant  | 72 |
| D.1.40 DP-Resistant  | 72 |
| D.1.41 DP-Resistant  | 72 |
| D.1.42 DP-Resistant  | 72 |
| D.1.43 DP-Resistant  | 72 |
| D.1.44 DP-Resistant  | 72 |
| D.1.45 DP-Resistant  | 72 |
| D.1.46 DP-Resistant  | 72 |
| D.1.47 DP-Resistant  | 72 |
| D.1.48 DP-Resistant  | 72 |
| D.1.49 DP-Resistant  | 72 |
| D.1.50 DP-Resistant  | 72 |
| D.1.51 DP-Resistant  | 72 |
| D.1.52 DP-Resistant  | 72 |
| D.1.53 DP-Resistant  | 72 |
| D.1.54 DP-Resistant  | 72 |
| D.1.55 DP-Resistant  | 72 |
| D.1.56 DP-Resistant  | 72 |
| D.1.57 DP-Resistant  | 72 |
| D.1.58 DP-Resistant  | 72 |
| D.1.59 DP-Resistant  | 72 |
| D.1.60 DP-Resistant  | 72 |
| D.1.61 DP-Resistant  | 72 |
| D.1.62 DP-Resistant  | 72 |
| D.1.63 DP-Resistant  | 72 |
| D.1.64 DP-Resistant  | 72 |
| D.1.65 DP-Resistant  | 72 |
| D.1.66 DP-Resistant  | 72 |
| D.1.67 DP-Resistant  | 72 |
| D.1.68 DP-Resistant  | 72 |
| D.1.69 DP-Resistant  | 72 |
| D.1.70 DP-Resistant  | 72 |
| D.1.71 DP-Resistant  | 72 |
| D.1.72 DP-Resistant  | 72 |
| D.1.73 DP-Resistant  | 72 |
| D.1.74 DP-Resistant  | 72 |
| D.1.75 DP-Resistant  | 72 |
| D.1.76 DP-Resistant  | 72 |
| D.1.77 DP-Resistant  | 72 |
| D.1.78 DP-Resistant  | 72 |
| D.1.79 DP-Resistant  | 72 |
| D.1.80 DP-Resistant  | 72 |
| D.1.81 DP-Resistant  | 72 |
| D.1.82 DP-Resistant  | 72 |
| D.1.83 DP-Resistant  | 72 |
| D.1.84 DP-Resistant  | 72 |
| D.1.85 DP-Resistant  | 72 |
| D.1.86 DP-Resistant  | 72 |
| D.1.87 DP-Resistant  | 72 |
| D.1.88 DP-Resistant  | 72 |
| D.1.89 DP-Resistant  | 72 |
| D.1.90 DP-Resistant  | 72 |
| D.1.91 DP-Resistant  | 72 |
| D.1.92 DP-Resistant  | 72 |
| D.1.93 DP-Resistant  | 72 |
| D.1.94 DP-Resistant  | 72 |
| D.1.95 DP-Resistant  | 72 |
| D.1.96 DP-Resistant  | 72 |
| D.1.97 DP-Resistant  | 72 |
| D.1.98 DP-Resistant  | 72 |
| D.1.99 DP-Resistant  | 72 |
| D.1.100 DP-Resistant   | 72 |

Ponomareva et al., 2023

# SWAG-PPM (Chew et al., 2025)

- **PPM** (Pseudo Posterior Mechanism)
  - Provides DP by randomly drawing from the posterior of model parameters
  - Downweights riskier records during training, limiting how much any single record can influence the model
- **SWAG** (Stochastic Weight Averaging – Gaussian)
  - Fits a multivariate Gaussian from additional training updates around a local minima
  - This Gaussian acts as an approximate (pseudo) posterior that can then be sampled from for a model release

Article  
Link



# Key Hyperparameters

- **Fine-tuning Epochs** (“*FT Epochs*”)
  - # of full passes over the training data
    - **Too few:** leave parameters off-mode, biasing SWAG
    - **Too many:** overfitting, model memorize training data
- **Disclosure Risk Weight Parameter** (“*c*”)
  - How quickly downweighing increases with increased disclosure risk
    - **High c:** high risk records downweighed more relative to low risk
    - **Low c:** flattens decline, lowering max contribution of all records
- **Hyperparameter Interactions**
  - *FT epochs* & *c* are not independent
  - Changing *c* alters both curvature and location of posterior, requiring re-running fine tuning

# Training Flowchart

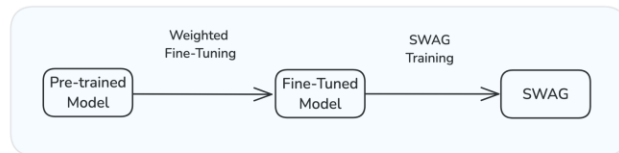
## Three steps:

1. Initial privacy & utility assessment
2. Model convergence ( “*FT epochs*” )
3. Privacy tuning ( “*c*” )

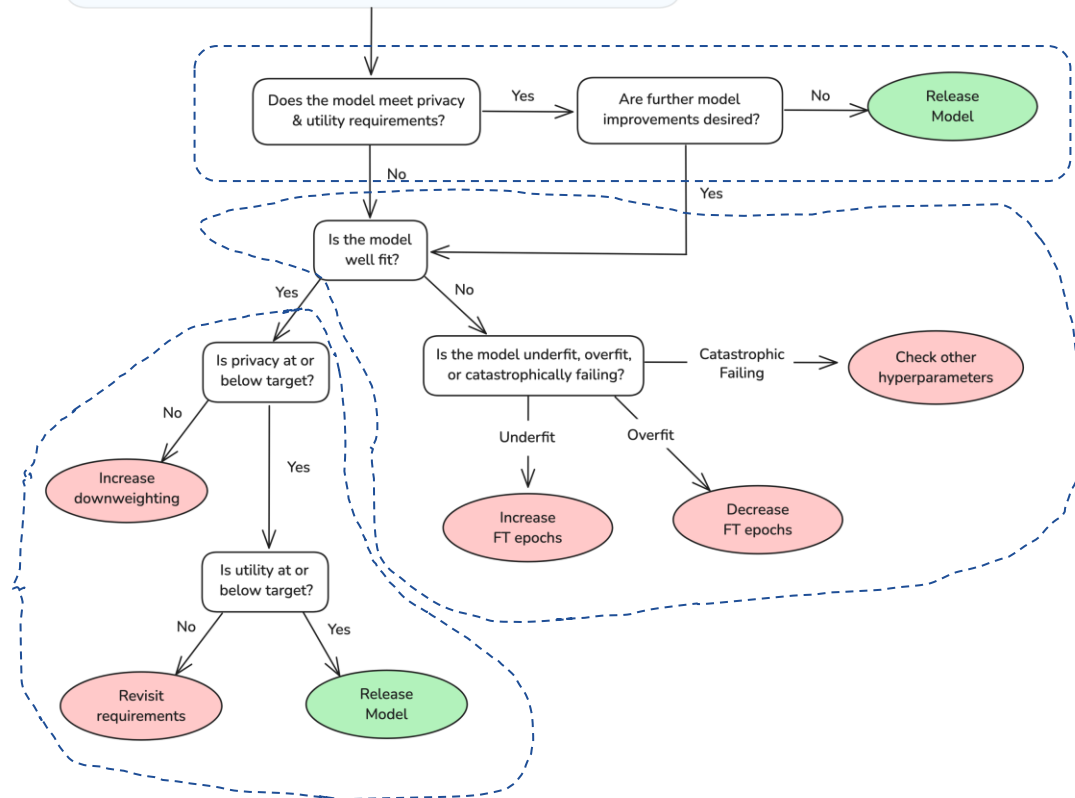
To illustrate these steps, we present a case study that trains an autocoder using the OSHA severe injury reports data.

- <https://www.osha.gov/severe-injury-reports>

Training



Evaluation



# Step 1: Initial Assessment

- **Targets**

- **Privacy:**  $\epsilon \leq 10$
- **Utility:**  $\leq 25\%$  degradation in F1 from non-private baseline

- **Privacy**

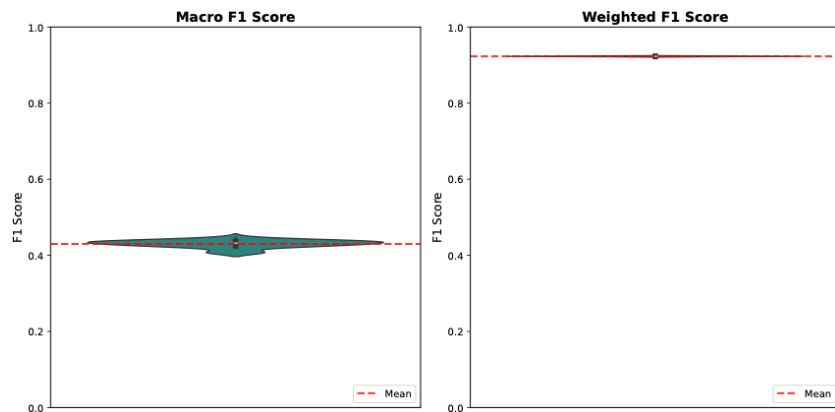
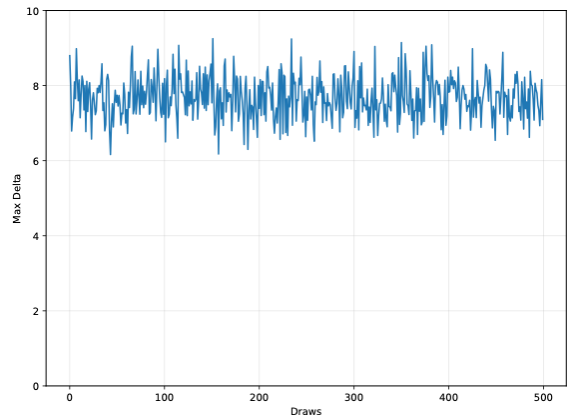
- “Max Delta” plot
  - X-axis: Posterior draws
  - Y-axis: Max loss training data
- $\epsilon = (2 * Max Loss) = 18.5$

- **Utility**

- “Utility” plot
  - F1 score (macro and weighted) across 30 posterior draws

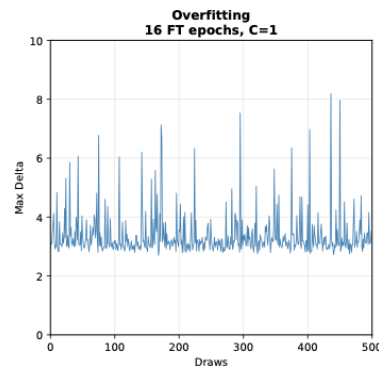
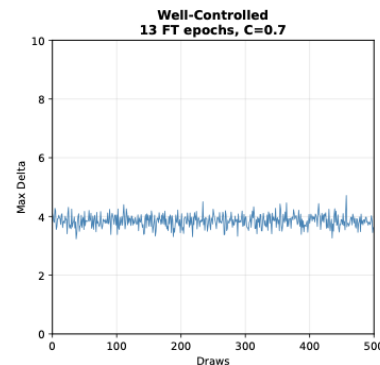
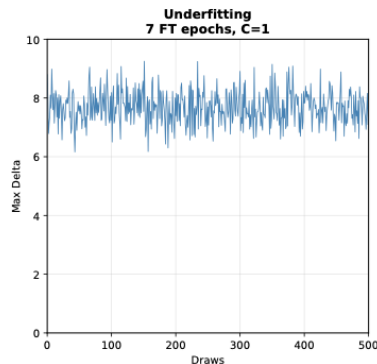
- **Decision**

- Does not meet privacy requirements



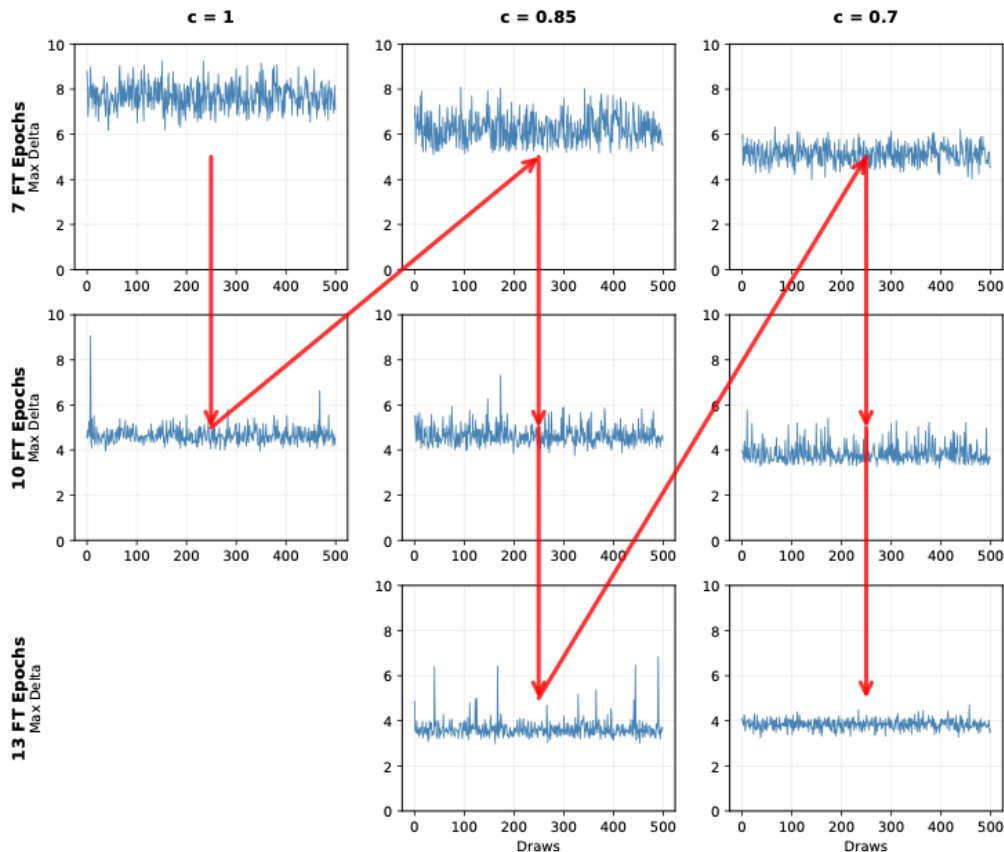
# Step 2: Convergence

- Model convergence is important because both privacy and utility depend on the stability of the model's parameter estimates.
- Unlike mature Bayesian procedures, SWAG does not offer diagnostics to assess convergence
- We repurpose the max delta plots to help:
  - **[Top]** “Underfit” models produce high, persistent variability across draws
  - **[Bottom]** “Overfit” models produce low variability punctuated by sudden spikes
  - **[Middle]** “Well-fit” models exhibit low, stable max deltas across draws
- **Decision**
  - Vary the *FT epochs* until model is well-controlled



# Step 3: Privacy Tuning

- If “converged” model does not meet privacy threshold, we can lower  $c$  to improve privacy at the expense of utility
- Start with high values of  $c$  and train until convergence, checking if privacy reqs are met along the way.
- If we see signs of overfitting, we reduce  $c$  and retrain.
- Since our privacy guarantee only relies on the final model release, intermediate tuning does not penalize privacy accounting
- **Decision**
  - Vary *FT epochs* and  $c$  until privacy and utility requirements are met.
  - If not possible, revisit requirements with stakeholders.



# Discussion

- We presented a procedural framework for SWAG-PPM to help practitioners make decisions about training and privacy-utility trade-offs.
- **Future Work**
  - Adaptive hyperparameter optimization that jointly updates  $c$  and FT epochs during training.
  - Systematic evaluations of SWAG-PPM across different data distributions (class imbalance, text lengths, feature modalities, etc.) would help test robustness and potentially reveal domain-specific tuning patterns
  - Extend the theory of PPM to support finite-sample guarantees to help strengthen confidence in practical deployments

# Further Reading

Check upcoming GASP special issue in *Journal of Data Science*



JOURNAL OF DATA SCIENCE 0 (0), 1–19  
July 2025

DOI: 0000

## A Practical Guide to Differentially Private Deep Learning using the Pseudo Posterior Mechanism

ALEXANDER J. PREISS<sup>1,\*</sup>, AMANDA KONET<sup>1</sup>, ROBERT CHEW<sup>1</sup>, MATTHEW R. WILLIAMS<sup>1</sup>,  
ELAN A. SEGARRA<sup>2</sup>, DAVID H. OH<sup>2</sup>, ERIN BOON<sup>2</sup>, AND TERRANCE D. SAVITSKY<sup>2</sup>

<sup>1</sup>Research Triangle Park, NC 27709, RTI International, USA

<sup>2</sup>Hillcrest Heights, MD 20746, U.S. Bureau of Labor Statistics, USA

### Abstract

Privacy-preserving machine learning methods seek to train useful models that do not disclose information about the data on which they were trained. Such methods are vital when organizations train neural networks on sensitive individual-level data and seek to release the models publicly. Their goal poses a trade-off between predictive performance (utility) and privacy protection. That trade-off makes privacy-preserving machine learning methods difficult to apply in practice, usually requiring extensive iteration and hyperparameter tuning. Yet, practitioners often have little guidance for navigating competing statistical, computational, and privacy demands. We present an implementation algorithm for the Stochastic Weight Averaging–Gaussian Pseudo Posterior Mechanism (SWAG-PPM), a Bayesian differentially private deep learning method. The implementation algorithm focuses on the joint tuning of two key hyperparameters whose interaction governs model convergence and the privacy–utility trade-off. We introduce novel diagnostic tools to evaluate convergence and guide hyperparameter adjustments. Using a transformer model for occupational injury classification, we demonstrate that diagnostic-guided tuning with SWAG-PPM can achieve strong privacy protection and utility. While our case study uses a specific dataset and model architecture, all methodological steps can apply to other settings where privacy risk is heterogeneously distributed.

**Keywords** *Bayesian deep learning; Differential Privacy; Imbalanced learning; Official statistics; Pseudo posterior distribution*



# Thank you

Rob Chew | [rchew@rti.org](mailto:rchew@rti.org)