

# SIGuard: Guarding Secure Inference with Post Data Privacy

Xinqian Wang<sup>\*</sup>, Xiaoning Liu<sup>\*</sup>, Shangqi Lai<sup>+</sup>, Xun Yi<sup>\*</sup>, Xingliang Yuan<sup>‡</sup>
\* RMIT University, <sup>†</sup> CSIRO's Data61, <sup>‡</sup> The University of Melbourne

# Machine Learning as a Service (MLaaS) offerings at the Cloud



Google Vision API Amazon SageMaker

Microsoft Project InnerEye (2020). https://www.microsoft.com/en-us/research/project/medical-image-analysis/
 Google DeepMind Health (2020). https://deepmind.com/blog/announcements/deepmind-health-joins-google-health
 Kuna AI (2017). https://getkuna.com/ blogs / news / 2017 - 05 - 24 - introducing - kuna-ai.

# Machine Learning as a Service (MLaaS) offerings at the Cloud

- Ready-made intelligence for a wide spectrum of applications
  - Medical imaging [1], clinical trial and research [2], home monitoring [3]



Microsoft Project InnerEye (2020). https://www.microsoft.com/en-us/research/project/medical-image-analysis/
 Google DeepMind Health (2020). https://deepmind.com/blog/announcements/deepmind-health-joins-google-health
 Kuna AI (2017). https://getkuna.com/ blogs / news / 2017 - 05 - 24 - introducing - kuna-ai.

# Machine Learning as a Service (MLaaS) offerings at the Cloud

- Ready-made intelligence for a wide spectrum of applications
  - Medical imaging [1], clinical trial and research [2], home monitoring [3]
- Essential offerings for cloud service providers
  - Google Cloud Vision API, Amazon AWS SageMaker



Microsoft Project InnerEye (2020). https://www.microsoft.com/en-us/research/project/medical-image-analysis/
 Google DeepMind Health (2020). https://deepmind.com/blog/announcements/deepmind-health-joins-google-health
 Kuna AI (2017). https://getkuna.com/ blogs / news / 2017 - 05 - 24 - introducing - kuna-ai.

- Cloud capitalizes on neural networks (NNs) to offer prediction services (e.g., medical diagnostics)
- User leverages the service to make a prediction over its data (e.g., brain MRI)





• User feeds data to the model through the API



- User feeds data to the model through the API
- Cloud hosts a pre-trained NN model and runs an inference function over user's data



- User feeds data to the model through the API
- Cloud hosts a pre-trained NN model and runs an inference function over user's data
- User receives confidence vectors to choose a quality prediction





Individual data is sensitive



Individual data is sensitive

NN models are valuable and IP



• *Input Privacy* - sensitive information using proprietary model *W* to classify sensitive individual data *X* 

Individual data is sensitive

NN models are valuable and IP



# Secure inference to guarantee input privacy

- Privacy-preserving machine learning (PPML) inference service, i.e., secure inference f(W, X)
  - A user and a model owner upload encrypted data Enc(X) and model Enc(W) to the cloud
  - Opt for secure multiparty computation (MPC) [1, 2, 3] for efficiency



P. Mishra, R. Lehmkuhl, A. Srinivasan, W. Zheng, and R. A. Popa, "Delphi: A cryptographic inference service for neural networks," in USENIX Security, 2020.
 Mohassel, Payman, and Yupeng Zhang. "Secureml: A system for scalable privacy-preserving machine learning." 2017 IEEE symposium on security and privacy (SP). IEEE, 2017.
 Mohassel, Payman, and Peter Rindal. "ABY3: A mixed protocol framework for machine learning." Proceedings of the 2018 ACM SIGSAC conference on computer and communications security. 2018.

# Secure inference to guarantee input privacy

- Privacy-preserving machine learning (PPML) inference service, i.e., secure inference f(W, X)
  - A user and a model owner upload encrypted data Enc(X) and model Enc(W) to the cloud
  - Opt for secure multiparty computation (MPC) [1, 2, 3] for efficiency

#### User learns the inference result s



P. Mishra, R. Lehmkuhl, A. Srinivasan, W. Zheng, and R. A. Popa, "Delphi: A cryptographic inference service for neural networks," in USENIX Security, 2020.
 Mohassel, Payman, and Yupeng Zhang. "Secureml: A system for scalable privacy-preserving machine learning." 2017 IEEE symposium on security and privacy (SP). IEEE, 2017.
 Mohassel, Payman, and Peter Rindal. "ABY3: A mixed protocol framework for machine learning." Proceedings of the 2018 ACM SIGSAC conference on computer and communications security. 2018.

### Secure inference to guarantee input privacy

- Privacy-preserving machine learning (PPML) inference service, i.e., secure inference f(W, X)
  - A user and a model owner upload encrypted data Enc(X) and model Enc(W) to the cloud
  - Opt for secure multiparty computation (MPC) [1, 2, 3] for efficiency

#### User learns the inference result s





P. Mishra, R. Lehmkuhl, A. Srinivasan, W. Zheng, and R. A. Popa, "Delphi: A cryptographic inference service for neural networks," in USENIX Security, 2020.
 Mohassel, Payman, and Yupeng Zhang. "Secureml: A system for scalable privacy-preserving machine learning." 2017 IEEE symposium on security and privacy (SP). IEEE, 2017.
 Mohassel, Payman, and Peter Rindal. "ABY3: A mixed protocol framework for machine learning." Proceedings of the 2018 ACM SIGSAC conference on computer and communications security. 2018.

Privacy threats on post data - prediction API attacks O' exploit inference result s to learn training dataset information by querying the model [1, 2]

![](_page_16_Figure_2.jpeg)

R. Shokri, M. Stronati, C. Song, and V. Shmatikov, "Membership inference attacks against machine learning models," in S&P, 2017.
 H. Yu, K. Yang, T. Zhang, Y.-Y. Tsai, T.-Y. Ho, and Y. Jin, "Cloudleak: Large-scale deep learning models stealing through adversarial examples," in NDSS, 2020.

Privacy threats on post data - prediction API attacks O' exploit inference result s to learn training dataset information by querying the model [1, 2]

*Output encodes knowledge of training dataset* 

![](_page_17_Figure_3.jpeg)

R. Shokri, M. Stronati, C. Song, and V. Shmatikov, "Membership inference attacks against machine learning models," in S&P, 2017.
 H. Yu, K. Yang, T. Zhang, Y.-Y. Tsai, T.-Y. Ho, and Y. Jin, "Cloudleak: Large-scale deep learning models stealing through adversarial examples," in NDSS, 2020.

Privacy threats on post data - prediction API attacks O' exploit inference result s to learn training dataset information by querying the model [1, 2]

#### *Output encodes knowledge of training dataset*

![](_page_18_Figure_3.jpeg)

R. Shokri, M. Stronati, C. Song, and V. Shmatikov, "Membership inference attacks against machine learning models," in S&P, 2017.
 H. Yu, K. Yang, T. Zhang, Y.-Y. Tsai, T.-Y. Ho, and Y. Jin, "Cloudleak: Large-scale deep learning models stealing through adversarial examples," in NDSS, 2020.

• **Privacy threats on post data** - Membership inference attacks (MIAs) [1] O' can exploit secure inference to infer whether the data X belongs to the encrypted model W's training dataset

![](_page_19_Figure_2.jpeg)

- Privacy threats on post data in PPML MIAs O' can exploit reconstructed inference result s to learn membership information by querying the encrypted model
- **Output Privacy** sensitive information revealed from secure inference output (predictions)

#### *Reconstructed output encodes knowledge of training dataset*

![](_page_20_Figure_4.jpeg)

![](_page_21_Figure_1.jpeg)

[1] Jia, Jinyuan, et al. "Memguard: Defending against black-box membership inference attacks via adversarial examples." CCS 2019.

[2] Chen, Zitao, and Karthik Pattabiraman. "Overconfidence is a dangerous thing: Mitigating membership inference attacks by enforcing less confident prediction." NDSS. 2023.

[3] Abadi, Martin, et al. "Deep learning with differential privacy." CCS 2016.

[4] Tang, Xinyu, et al. "Mitigating membership inference attacks by {Self-Distillation} through a novel ensemble architecture." USENIX Security 2022.

[5] Nasr, Milad, Reza Shokri, and Amir Houmansadr. "Machine learning with membership privacy using adversarial regularization." CCS 2018.

- Training time defenses use specific training approaches
  - However, they inherently reduce prediction accuracy, e.g., differentially private training [3]

![](_page_22_Figure_3.jpeg)

[1] Jia, Jinyuan, et al. "Memguard: Defending against black-box membership inference attacks via adversarial examples." CCS 2019.

[2] Chen, Zitao, and Karthik Pattabiraman. "Overconfidence is a dangerous thing: Mitigating membership inference attacks by enforcing less confident prediction." NDSS. 2023.

[3] Abadi, Martin, et al. "Deep learning with differential privacy." CCS 2016.

[4] Tang, Xinyu, et al. "Mitigating membership inference attacks by {Self-Distillation} through a novel ensemble architecture." USENIX Security 2022.

[5] Nasr, Milad, Reza Shokri, and Amir Houmansadr. "Machine learning with membership privacy using adversarial regularization." CCS 2018.

- Training time defenses use specific training approaches
  - However, they inherently reduce prediction accuracy, e.g., differentially private training [3]
  - Prefer inference time defense without accuracy loss

![](_page_23_Figure_4.jpeg)

[1] Jia, Jinyuan, et al. "Memguard: Defending against black-box membership inference attacks via adversarial examples." CCS 2019.

[2] Chen, Zitao, and Karthik Pattabiraman. "Overconfidence is a dangerous thing: Mitigating membership inference attacks by enforcing less confident prediction." NDSS. 2023.

[3] Abadi, Martin, et al. "Deep learning with differential privacy." CCS 2016.

[4] Tang, Xinyu, et al. "Mitigating membership inference attacks by {Self-Distillation} through a novel ensemble architecture." USENIX Security 2022.

[5] Nasr, Milad, Reza Shokri, and Amir Houmansadr. "Machine learning with membership privacy using adversarial regularization." CCS 2018.

• Inference time defenses inject carefully crafted noises to perturb the inference output

![](_page_24_Figure_2.jpeg)

- Inference time defenses inject carefully crafted noises to perturb the inference output
  - HAMP [2] needs to apply defense in both training and inference, thus needs costly re-training

![](_page_25_Figure_3.jpeg)

- Inference time defenses inject carefully crafted noises to perturb the inference output
  - HAMP [2] needs to apply defense in both training and inference, thus needs costly re-training

![](_page_26_Figure_3.jpeg)

- Inference time defenses inject carefully crafted noises to perturb the inference output
  - HAMP [2] needs to apply defense in both training and inference, thus needs costly re-training
  - MemGuard [1] is chosen as our starting point

![](_page_27_Figure_4.jpeg)

- Inference time defenses inject carefully crafted noises to perturb the inference output
  - HAMP [2] needs to apply defense in both training and inference, thus needs costly re-training
  - MemGuard [1] is chosen as our starting point

![](_page_28_Figure_4.jpeg)

- Inference time defenses inject carefully crafted noises to perturb the inference output
  - HAMP [2] needs to apply defense in both training and inference, thus needs costly re-training
  - MemGuard [1] is chosen as our starting point

![](_page_29_Figure_4.jpeg)

- Inference time defenses inject carefully crafted noises to perturb the inference output
  - HAMP [2] needs to apply defense in both training and inference, thus needs costly re-training
  - MemGuard [1] is chosen as our starting point

![](_page_30_Figure_4.jpeg)

- Inference time defenses inject carefully crafted noises to perturb the inference output
  - HAMP [2] needs to apply defense in both training and inference, thus needs costly re-training
  - MemGuard [1] is chosen as our starting point

![](_page_31_Figure_4.jpeg)

- Advantages
  Preserve prediction accuracy
  No need for costly re-training
  Do not disrupt MLaaS pipeline

• Inference time defenses inject carefully crafted noises to perturb the inference output

![](_page_32_Figure_2.jpeg)

- Inference time defenses inject carefully crafted noises to perturb the inference output
  - However, it cannot directly be adopted in secure inference

![](_page_33_Figure_3.jpeg)

- Preserve prediction accuracy
  No need for costly re-training
  Do not disrupt MLaaS pipeline

- Inference time defenses inject carefully crafted noises to perturb the inference output
  - However, it cannot directly be adopted in secure inference

![](_page_34_Figure_3.jpeg)

- Inference time defenses inject carefully crafted noises to perturb the inference output
  - However, it cannot directly be adopted in secure inference

![](_page_35_Figure_3.jpeg)
### Existing MIA defenses in the plaintext domain

- Inference time defenses inject carefully crafted noises to perturb the inference output
  - However, it cannot directly be adopted in secure inference



Our research

# Building **SIGuard** framework

Secure Inference Guard

#### guarding output privacy

integrated into versatile MPC-based secure inference

Our philosophy

# Harnessing insights of MPC techniques and machine learning

stringent and provable security guarantees preserving model prediction accuracy

- 3-party MPC techniques to achieve privacy guarantees
- SIGuard protocol plugs into secure inference



- 3-party MPC techniques to achieve privacy guarantees
- SIGuard protocol plugs into secure inference



- 3-party MPC techniques to achieve privacy guarantees
- SIGuard protocol plugs into secure inference



- 3-party MPC techniques to achieve privacy guarantees
- SIGuard protocol plugs into secure inference



- 3-party MPC techniques to achieve privacy guarantees
- SIGuard protocol plugs into secure inference



- 3-party MPC techniques to achieve privacy guarantees
- SIGuard protocol plugs into secure inference



- Independent corruption
- Collusion broaden attack surface
  - MIA adversary corrupting user, colludes with secure inference adversary corrupting one cloud server



#### MPC techniques - Replicate secret sharing (RSS)

• Replicate secret sharing encrypts private data as secret shares



### MPC techniques - Replicate secret sharing (RSS)

• Replicate secret sharing encrypts private data as secret shares



### Machine learning - MIAs

• The MIA adversary trains a membership (binary) classifier  $f(\vec{s}) = \begin{cases} 1 & x \in M$ 's training dataset  $x \notin M$ 's training dataset

**Ideal attack** 



Membership dataset D

Ground-truth



Member prediction  $\vec{s}$ 

Non-member prediction  $\vec{s}$ 

100% membership classify accuracy!



• Insight: MIAs classifier is vulnerable to *adversarial examples*, fool MIAs by perturbing confidence to adversarial example

- Insight: MIAs classifier is vulnerable to *adversarial examples*, fool MIAs by perturbing confidence to adversarial example
- Approach: adding a noise vector  $\vec{e}$  to perturb logits  $\vec{z}$ , while keeping final prediction unchanged  $\vec{s}' = \vec{s}$

- Insight: MIAs classifier is vulnerable to *adversarial examples*, fool MIAs by perturbing confidence to adversarial example
- Approach: adding a noise vector  $\vec{e}$  to perturb logits  $\vec{z}$ , while keeping final prediction unchanged  $\vec{s}' = \vec{s}$
- Workflow:
  - Iteratively optimizes the noise vector  $\vec{e}$  by minimizing the loss  $\mathcal{L}$  (c3 controls optimization pace)



- Insight: MIAs classifier is vulnerable to *adversarial examples*, fool MIAs by perturbing confidence to adversarial example
- Approach: adding a noise vector  $\vec{e}$  to perturb logits  $\vec{z}$ , while keeping final prediction unchanged  $\vec{s}' = \vec{s}$
- Workflow:
  - Iteratively optimizes the noise vector  $\vec{e}$  by minimizing the loss  $\mathcal{L}$  (c3 controls optimization pace)
  - Generate perturbed prediction (confidence)  $\vec{s}' = softmax(\vec{z} + \vec{e})$



- Insight: MIAs classifier is vulnerable to *adversarial examples*, fool MIAs by perturbing confidence to adversarial example
- Approach: adding a noise vector  $\vec{e}$  to perturb logits  $\vec{z}$ , while keeping final prediction unchanged  $\vec{s}' = \vec{s}$
- Workflow:
  - Iteratively optimizes the noise vector  $\vec{e}$  by minimizing the loss  $\mathcal{L}$  (*c*3 controls optimization pace)
  - Generate perturbed prediction (confidence)  $\vec{s}' = softmax(\vec{z} + \vec{e})$
  - Terminate when  $h(\vec{s}') = 0$ , i.e., classifier *h* couldn't determine the membership of perturbed prediction and noise  $\vec{e}$  can't be further minimized



- MPC requires to approximate non-linear functions, e.g., softmax (division, exp)
  - E.g., SecureML [1] approximates softmax using ReLU



- MPC requires to approximate non-linear functions, e.g., softmax (division, exp)
  - E.g., SecureML [1] approximates softmax using ReLU
- Approximated softmax in secure inference perturbs confidence scores



- MPC requires to approximate non-linear functions, e.g., softmax (division, exp)
  - E.g., SecureML [1] approximates softmax using ReLU
- Approximated softmax in secure inference perturbs confidence scores



**<u>RQ1</u>**: Whether MIAs can still exploit secure inference with perturbed predictions?

# Observation 2: Secure post inference defense just like MemGuard

- Inspired by MemGuard [1], SIGuard's protocol injects carefully crafted perturbations into the encrypted predictions without harming the inference accuracy
- Takes encrypted logits, and finally returns encrypted perturbed confidence to the user



# Observation 2: Secure post inference defense just like MemGuard

- Inspired by MemGuard [1], SIGuard's protocol injects carefully crafted perturbations into the encrypted predictions without harming the inference accuracy
- Takes encrypted logits, and finally returns encrypted perturbed confidence to the user



# Observation 2: Secure post inference defense just like MemGuard

- Inspired by MemGuard [1], SIGuard's protocol injects carefully crafted perturbations into the encrypted predictions without harming the inference accuracy
- Takes encrypted logits, and finally returns encrypted perturbed confidence to the user

**<u>RQ2</u>**: How to efficiently realize SIGuard with MPC?



#### Observation 3: Broaden MIAs attack surface

• MIAs in plaintext can only access perturbed confidence scores



#### **Observation 3: Broaden MIAs attack surface**

- MIAs in secure inference can obtain knowledge of how the secure defense mechanism operates
  - When the corrupted user colludes with one corrupted cloud server
  - Attempt to leverage auxiliary knowledge to bypass defense



#### **Observation 3: Broaden MIAs attack surface**

- MIAs in secure inference can obtain knowledge of how the secure defense mechanism operates
  - When the corrupted user colludes with one corrupted cloud server
  - Attempt to leverage auxiliary knowledge to bypass defense

**<u>RQ3</u>**: How to mitigate the leakages and achieve the stringent privacy guarantee?





<sup>[1]</sup> Mohassel, Payman, and Yupeng Zhang. "Secureml: A system for scalable privacy-preserving machine learning." S&P, 2017.

<sup>[2]</sup> Knott, Brian, et al. "Crypten: Secure multi-party computation meets machine learning." NeuIPS, 2021.

<sup>[3]</sup> Watson, Jean-Luc, Sameer Wagh, and Raluca Ada Popa. "Piranha: A {GPU} platform for secure computation." USENIX Security. 2022.

<sup>[4]</sup> Aly, Abdelrahaman, and Nigel P. Smart. "Benchmarking privacy preserving scientific operations" ACNS, 2019.

• The risk of MIAs against secure inference remains significant, and in some cases may be even higher.



• The risk of MIAs against secure inference remains significant, and in some cases may be even higher.



• The risk of MIAs against secure inference remains significant, and in some cases may be even higher.





Secure Noise Optimization protocol aims to find an optimal noise vector to perturb the secret-shared confidence vector.



Secure Noise Optimization protocol aims to find an optimal noise vector to perturb the secret-shared confidence vector.



**Secure Noise Validation** protocol aims to validate that whether the secret-shared noise vector is carefully crafted to preserve the inference accuracy.

Secure Noise Optimization protocol aims to find an optimal noise vector to perturb the secret-shared confidence vector.



**Secure Noise Validation** protocol aims to validate that whether the secret-shared noise vector is carefully crafted to preserve the inference accuracy.

**Softmax** needs a careful design for softmax approximation – whether it affects defense performance
- MPC cannot directly compute While loop in the encrypted domain
  - The secret-shared condition needs to be revealed to terminate the loop



- MPC cannot directly compute While loop in the encrypted domain
  - The secret-shared condition needs to be revealed to terminate the loop
- Revealing the condition lets MIA know the number of iterations to find the optimal noise



- MPC cannot directly compute While loop in the encrypted domain
  - The secret-shared condition needs to be revealed to terminate the loop
- Revealing the condition lets MIA know the number of iterations to find the optimal noise
  - Broadened attack surface: MIA colludes with secure inference adversary
  - Condition  $h(softmax(\vec{z} + \vec{e})) = 0$
- Would revealing the number of iterations increase MIA's attack performance?



- MPC cannot directly compute While loop in the encrypted domain
  - The secret-shared condition needs to be revealed to terminate the loop
- Revealing the condition lets MIA know the number of iterations to find the optimal noise
  - Broadened attack surface: MIA colludes with secure inference adversary
  - Condition  $h(softmax(\vec{z} + \vec{e})) = 0$
- Would revealing the number of iterations increase MIA's attack performance?



## Refinement I: Mitigating potential leakages from iterations

- The revealed number of iterations is linked with data sample's membership information
- Even if the condition is protected, MIA still observes functionality to learn the number of iterations due to collusion



Fig. 4: Comparison of iteration frequency between member and non-member samples.

No member data when #iter >= 160 Must be non-member

### Refinement I: Mitigating potential leakages from iterations

- Difference of total iteration numbers is caused by each sample's optimization hardness
- Convert While to For
- Use a fixed number of iterations (inner\_loop, outer\_loop)



### Refinement II: Balanced efficiency & defense effectiveness

- Less iterations: faster v.s. ↓defense effectiveness (↓accuracy)
- More iterations: heavier MPC computation v.s. 1 defense effectiveness
- How to find suitable inner\_Loop, outer\_Loop balance the efficiency and defense?



#### Refinement II: Balanced efficiency & defense effectiveness

• outer\_loop is fixed to 3

inner\_loop = 10 (MemGuard sets 300)
lr = 0.8
defense near 50%

Compared with our basic SIGuard Runtime: basic 64s  $\rightarrow$  1.1s **58x** savings Bandwidth: basic 301MB  $\rightarrow$  5.4MB **55x** savings



Fig. 5: Comparison of SlGuard's defense performance under different learning rates and iterations (*inner\_loop*  $\in$  {10, 19, 38, 75, 150, 300}).

#### Implementations

• We leverage the MP-SPDZ framework as the skeleton of SIGuard.

**Q1**:Can SIGuard effectively mitigate the risk of MIAs in secure inference?

**Q2**:Is SIGuard efficient to be integrated into secure inference pipeline?



## Defense performance of SIGuard

#### SIGuard can effectively mitigate MIAs from balanced accuracy and TPR @ FPR.



Close to random guessing of 50%

## SIGuard's efficiency

SIGuard can effectively defeat MIAs for secure inference

without introducing dominant overhead.



LAN

WAN

# Thanks for listening!

## **Questions?**

Contact me:

xinqian.wang@rmit.edu.au