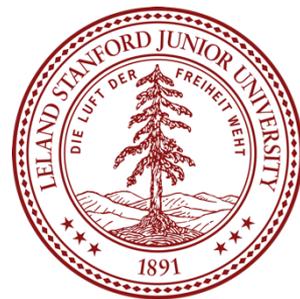# Slalom:
# Fast, Verifiable and Private Execution of Neural Networks in Trusted Hardware
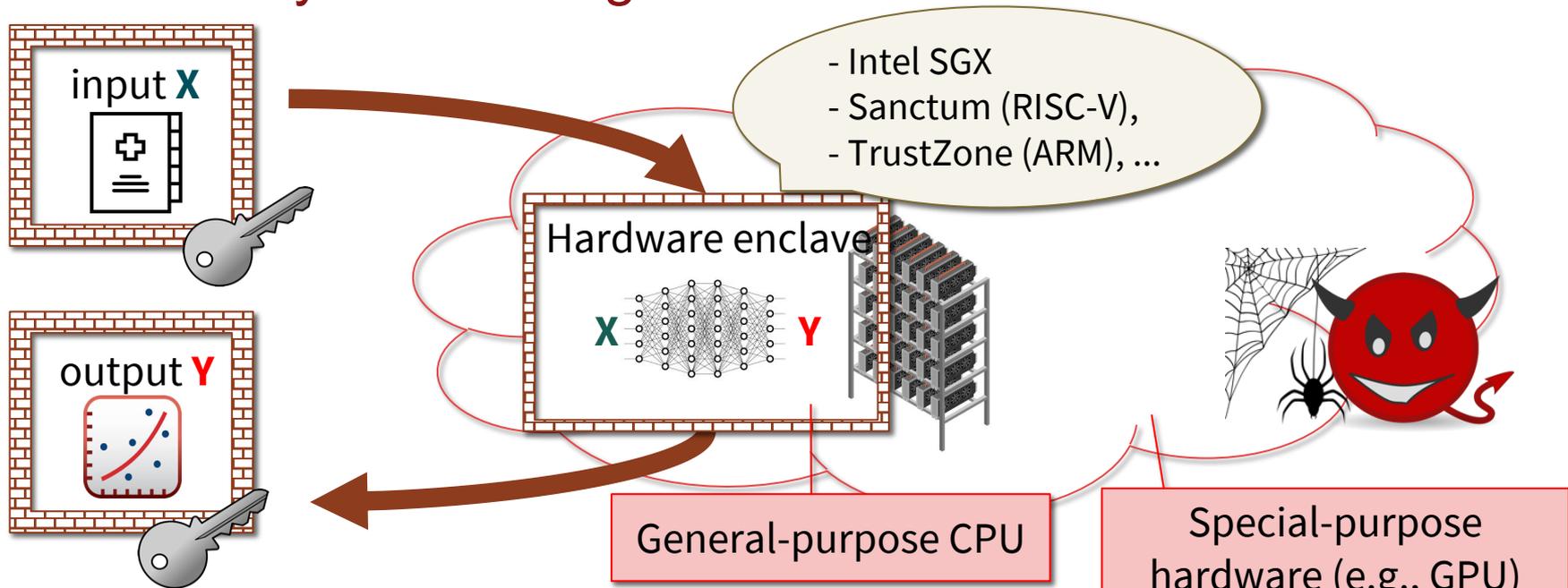
FLORIAN TRAMÈR & DAN BONEH

ICLR, New Orleans

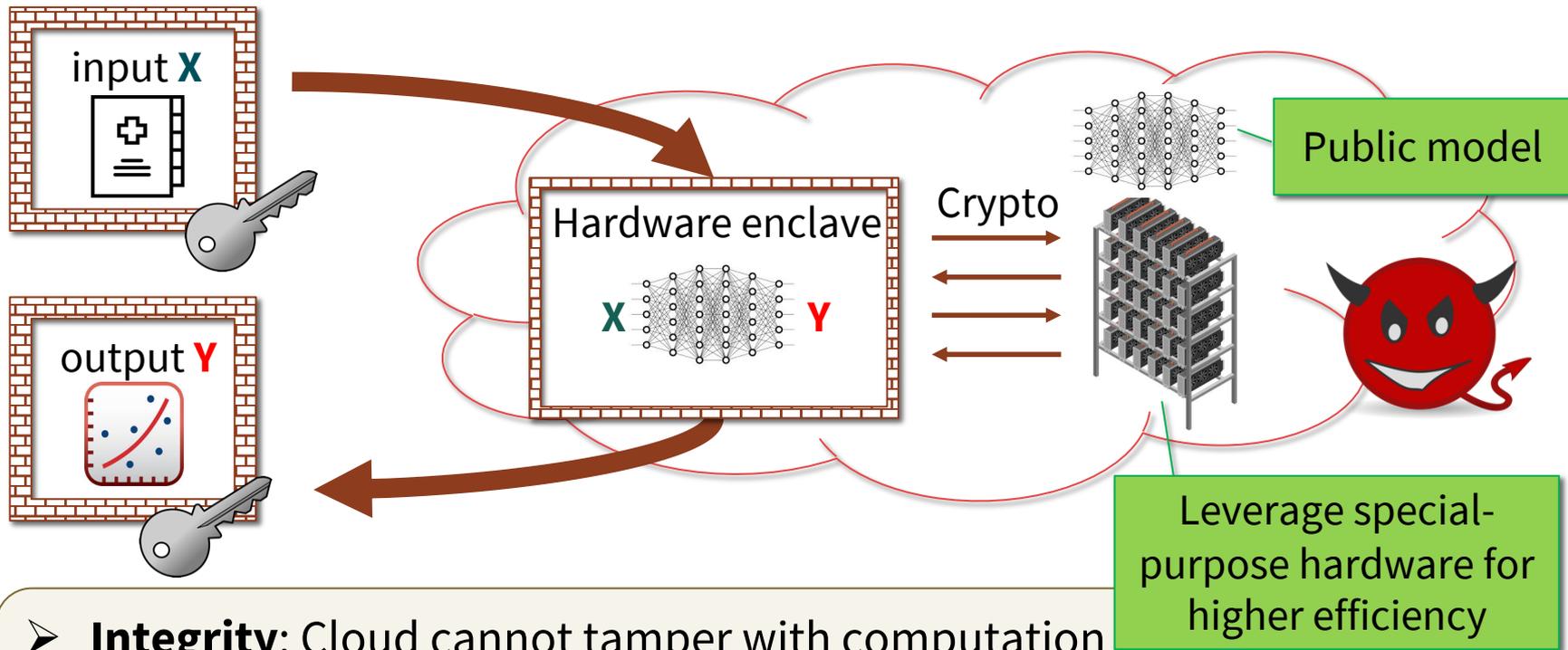May 7th 2019

# Securely outsourcing ML inference with hardware isolation

input **X**

- Intel SGX
- Sanctum (RISC-V),
- TrustZone (ARM), ...

Hardware enclave

**X**    **Y**

output **Y**

General-purpose CPU

Special-purpose hardware (e.g., GPU) provides no security

➤ **Integrity**: Cloud cannot tamper with computation
➤ **Privacy**: Integrity + Cloud does not learn inputs
➤ **Model Privacy**: Cloud does not learn model

**Stanford University**

# Slalom: Outsource ML from CPU enclave to special-purpose hardware



input **X**

output **Y**

Hardware enclave

**X**    **Y**

Crypto

Public model
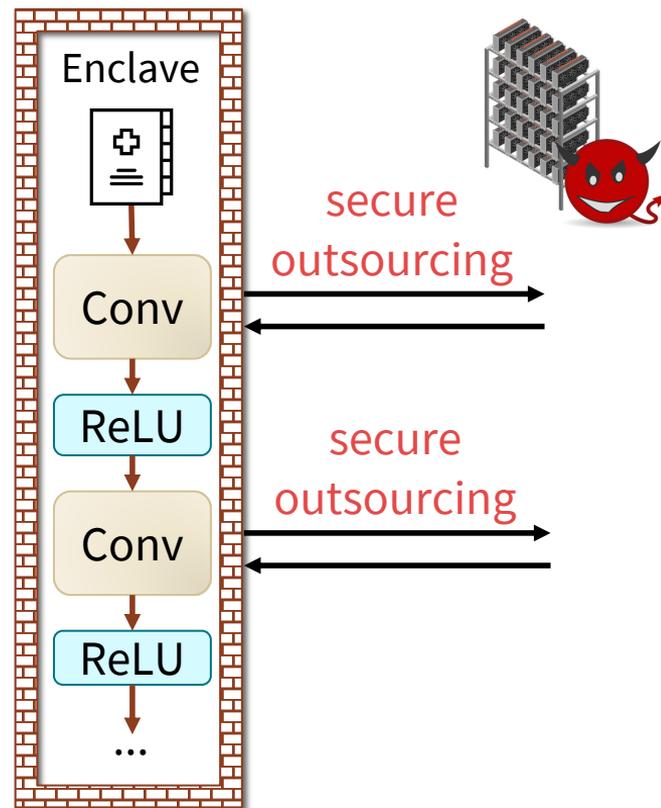
Leverage special-purpose hardware for higher efficiency

➢ **Integrity**: Cloud cannot tamper with computation
➢ **Privacy**: Integrity + Cloud does not learn inputs
➢ ~~**Model Privacy**: Cloud does not learn model~~

**Stanford University**

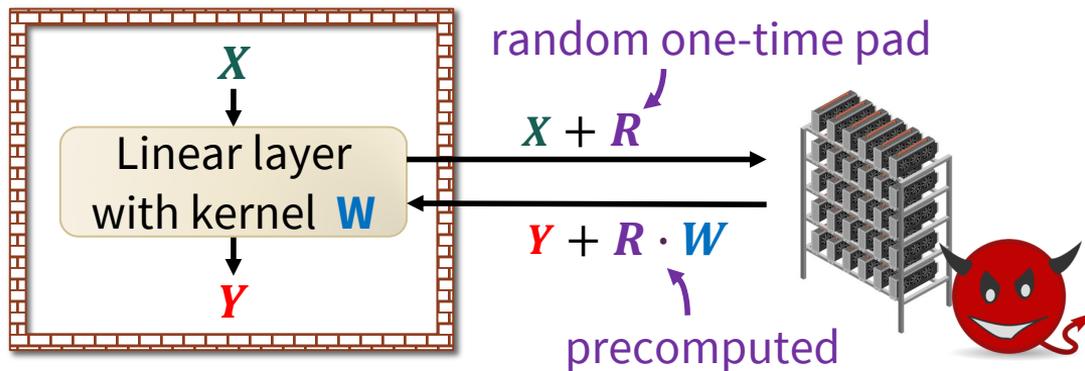# Outsourcing ML inference using cryptography

*Slalom uses cryptographic protocols to securely outsource all <u>linear layers</u> from the enclave to a GPU.*

⚠️ ▪ Crypto protocols have high communication costs

💡 › Enclave processor and GPU are co-located

› For VGG16, Slalom sends **50MB** of data from the enclave to the GPU per inference

⚠️ ▪ Crypto protocols are very efficient for securely outsourcing linear functions

💡 › Most of the computation in a DNN is linear (convolutions, dense, etc.)

› E.g., **~99%** for VGG16 and MobileNet

**Stanford University**

# How to securely outsource a matrix product



random one-time pad

$X + R$

$Y + R \cdot W$

precomputed

Linear layer with kernel $W$

$X$

$Y$

- **Integrity**:
  - › Verify that $Y = X \cdot W$
  - › Check $Y \cdot \vec{r} \stackrel{?}{=} X \cdot (W \cdot \vec{r})$   [Freivalds 1977]

Verify a matrix product with a few inner products
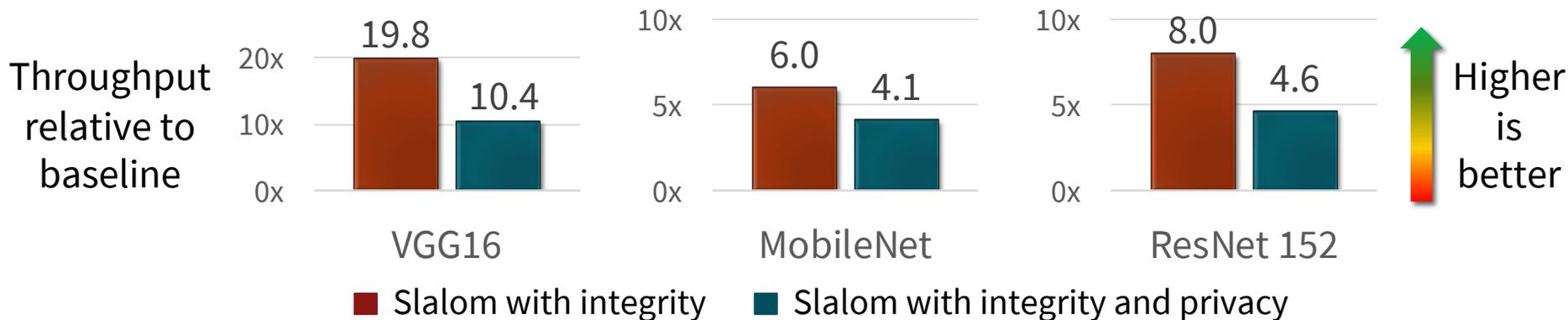(generalizes to arbitrary linear layer)

- **Privacy**:
  - › Evaluate model on random data $R$ in offline pre-processing phase
  - › Store $(R, R \cdot W)$ in the enclave and use these to encrypt & decrypt the communication with the GPU

Stanford University

# Evaluation

- Intel SGX + Nvidia Titan XP
- Throughput for ImageNet inference
- **Goal: Slalom (TEE↔GPU) ≫ TEE$_{baseline}$**

Evaluate DNN in TEE



Throughput relative to baseline

VGG16 — 19.8 / 10.4
MobileNet — 6.0 / 4.1
ResNet 152 — 8.0 / 4.6

Higher is better

■ Slalom with integrity     ■ Slalom with integrity and privacy

Slalom is **10-20x slower** than evaluating on GPU **(with no security guarantees)**
⇒ But, Slalom only utilizes the GPU ~10% of the time
⇒ Multiple CPU enclaves can outsource to the same GPU

*Stanford University*

# Conclusions & Open Problems

- **Slalom allows efficient and secure outsourcing of sensitive DNN computations to the cloud**
  - › Hardware isolation protects privacy & integrity, but is slow
  - › Slalom uses cryptography to leverage fast special-purpose hardware without any isolation guarantees

- **What about training?**
  - › Integrity: Freivalds' still works ☺
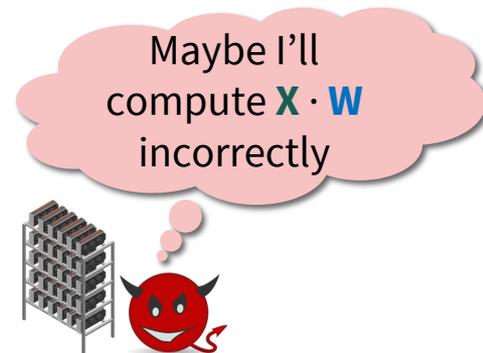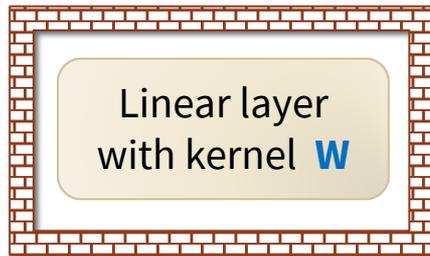  - › Privacy: Model itself should remain secret ☹

https://arxiv.org/abs/1806.03287

https://github.com/ftramer/slalom

https://floriantramer.com

Poster @4:30 - Great Hall BC #44

Stanford University

# How to securely outsource a linear layer

- **Quantization**: Evaluate a DNN over $\mathbb{Z}_p$ for a large prime p

- **Integrity**: Freivalds' 1977

$$Y \stackrel{?}{=} X \cdot W$$

$$\text{check } Y \cdot r \stackrel{?}{=} X \cdot (W \cdot r)$$

random vector

Linear layer with kernel **W**

**X**

**Y**

Maybe I'll compute **X** · **W** incorrectly

Verify any linear layer with a few inner products $\approx O(n^2)$ instead of $O(n^3)$

- **Privacy**: precomputed "one-time pads"
  - › See paper for details

Evaluate model on **random** data in offline preprocessing phase

# Privacy with precomputed one-time pads



random one-time pad

$X + R$

$Y = W \cdot (X + R)$

Linear layer with kernel $W$

$W \cdot X = Y - (W \cdot R)$

precompute this