# Plinius: Secure and Persistent Machine Learning Model Training

**Peterson Yuhala**[1]  **Pascal Felber**[1]  **Valerio Schiavoni**[1]  **Alain Tchana**[2]

[1]University of Neuchâtel, Switzerland

[2]ENS Lyon, France
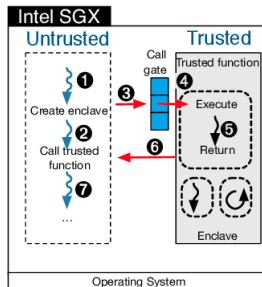
**Context**

Intel SGX

PM

Architecture

Evaluation

Conclusion

- Increasing popularity of cloud-based ML services (*e.g.,* Amazon ML, MS Azure AI).

- Security and privacy issues, *i.e.,* sensitive training data and models.

- DRAM scalability issues and high-access times of secondary storage = bottlenecks for ML.

- We need practical solutions to both problems.

- We solve security issues with TEEs (*e.g.,* Intel SGX).

- Secure *enclaves*: no system functionality, *i.e.,* system calls

- Legacy applications must be re-written/partitioned.

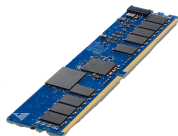# Persistent Memory

- We use persistent memory (PM) to solve DRAM/storage related issues. PM is:

  - Byte-addressable (like DRAM), and accessed via Load/Store.

  - Fast (low-latency, faster than SSD)

  - Persistent (like SSD)
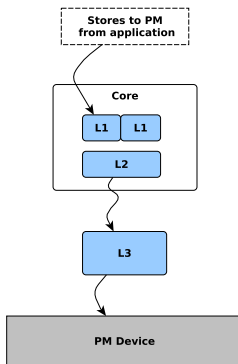
  - Higher capacity than DRAM

# How to use PM

- Like secondary storage: no program changes but smaller performance improvements.

- Leverage byte-addressability: requires program changes but better performance.

# Plinius in a nutshell
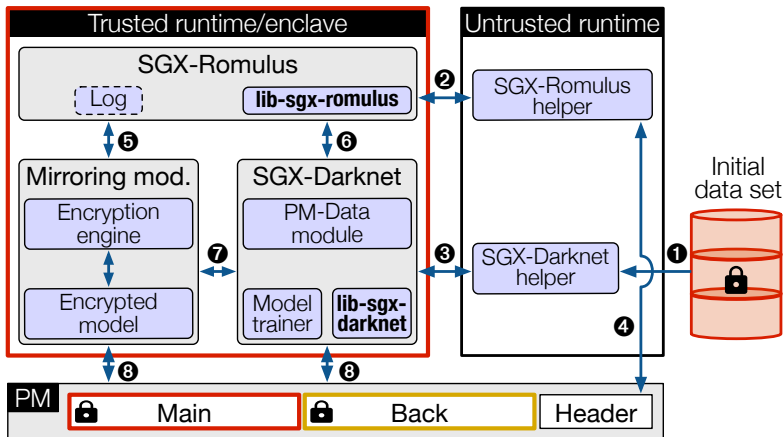
- Plinius ports a PM and ML library into SGX.

- It leverages the byte-addressability of PM for fast access to data in PM.

- Models trained in the enclave are mirrored to/from PM.

- How does Plinius improve save/restore performance ?

- How scalable is Plinius with varying model sizes ?

# Evaluation

sgx−emlPM – Mirroring Step: Save / sgx−emlPM – Mirroring Step: Restore

- Emulated PM + real SGX server: saves $3.5\times$ and restores $2.5\times$ faster vs SSD.

# Evaluation

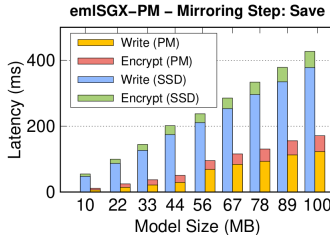- Performance drops at the EPC limit due to page swapping operations.

emISGX–PM – Mirroring Step: Save

Legend: Write (PM), Encrypt (PM), Write (SSD), Encrypt (SSD). Latency (ms) vs Model Size (MB): 10, 22, 33, 44, 56, 67, 78, 89, 100

emISGX–PM – Mirroring Step: Restore

Legend: Decrypt (PM), Read (PM), Decrypt (SSD), Read (SSD). Latency (ms) vs Model Size (MB): 10, 22, 33, 44, 56, 67, 78, 89, 100

- Real PM + sim SGX server: saves $3.2\times$ and restores $3.7\times$ faster vs SSD.

- Plinius is the first framework to leverage SGX for security and PM for fault tolerance.

- We leverage a mirroring mechanism for fault tolerance.

- Model and training data in memory $\rightarrow$ near instantaneous recovery after crashes.

- Test Plinius on github: