



# **Blinded Memory**

N. Asokan, Secure Systems Group □ <u>https://asokan.org/asokan/</u>
✓ X @nasokan

(Joint work with Hossam ElAtali, Lachlan J. Gunn, John Z. Jekel, Hans Liljestrand)

### This talk in a nutshell

### 1. Outsourced computing is everywhere...

• Machine learning models kept behind remote APIs

### 2. ...but this introduces security risks...

- Providers don't expose models/code to clients
- Clients expose sensitive data to providers
- Existing solutions like FHE/TEEs have drawbacks

#### 3. ...so we propose a new solution, Blinded Memory

- Attestation + standard encryption + hardware-assisted taint tracking
- Sensitive data not exposed to output devices or covert channels

### **Scenario: outsourced computation**

#### Goal: run the server's confidential code over client's confidential data

• Initial target: Outsourced ML inference and/or training



How can the client avoid revealing data to the service provider?

- Fully-Homomorphic Encryption: slow due to computational overhead
- Multi-Party Computation: slow due to network overhead
- Hardware-based isolation + remote attestation: fast

### Hardware-assisted TEEs are pervasive





### TEEs as an idea date back to the 1980s



"Security bit for designating the security status of information stored in a nonvolatile memory"

# Deployment of mobile TEEs date back to the 2000s



#### First deployment: Nokia 6630 ("Charlie")

• first 3G phone with TI OMAP 1710 processor (June 2004)

#### ARM TrustZone currently widely deployed

• <u>TrustZone-M for Cortex-M class microcontrollers (2016)</u>

#### Ca. 2008, TEE unheard of in academic circles

• first papers in FC 2008, ASIACCS 2009 [AE08] <u>A Platform for OnBoard Credentials</u>, Financial Cryptography and Data Security (2008) [KEAR09] <u>On-board credentials with open provisioning</u>, ACM ASIACCS (2009)

#### Intel SGX

• SkyLake (2015); wide availability of SDK "democratized" TEE research



### More on the history of TEEs

CCS 2019 keynote <a href="https://youtu.be/hHYoGn5PSI4">https://youtu.be/hHYoGn5PSI4</a>



2022 book https://ssg.aalto.fi/publications/hardware-platform-security-for-mobile-devices/

Foundations and Trends® in Privacy and Security 3:3:4

Hardware Platform Security for Mobile Devices

Lachlan J. Gunn, N. Asokan, Jan-Erik Ekberg, Hans Liljestrand, Vijayanand Nayani and Thomas Nyman

> new the essence of knowledge

#### 2019 blog article: https://blog.ssg.aalto.fi/2019/06/historical-insight-into-development-of.html

[A20] <u>"Hardware-assisted Trusted Execution Environments: Look Back, Look Ahead</u>", ACM CCS Keynote (2019) [GAE+22] "<u>Hardware Platform Security for Mobile Devices</u>", Foundations and Trends® in Privacy and Security 3(3-4):214-394, NOW publishers (2022) [MNA19] "<u>Historical insight into the development of Mobile TEEs</u>", SSG blog (2019)

### **Protection provided by TEEs comes with caveats**

TEEs provide an isolated environment for execution of software

TEEs are unsuitable when server code is confidential or unverfiable

• TEEs intended for clients to run code they trust and can verify

**Confidentiality of client data in TEEs is hampered by:** 

- Large TEE code base → vulnerable to software flaws
- Sharing resources → vulnerable to side channels

### Is Confidentiality vs. Performance a tradeoff?



### What can be done?

**1. Prevent** application software from leaking sensitive data

- Use hardware-assisted taint-tracking
- Need not verify trustworthiness of application s/w

### 2. Minimize resource sharing

- Move critical operations to a fixed-function, isolated processor (HSM)
- All HSM code analyzed in advance, guaranteed not to be malicious

### **Prevent leakage of sensitive data via CPU extensions**

### "Safe" streams of instructions don't expose sensitive data

#### Allowed:

• Computation on sensitive data by arbitrary, unattested, untrusted software

### **Prohibited:**

• Leaking sensitive data into any observable state, e.g.: peripherals outside security boundary, microarchitectural state

#### Use taint-tracking-based security policy to limit sensitive data to safe places

### **Combine with attestable HSM to assure clients**

**Remote attestation assures use of client data is subject to security policy** 



### **Taint tracking policy**

Registers/memory have an associated "sensitive" bit ("Blinded") Ideal rule:

Blinded(output<sub>m</sub>)  $\leftarrow \exists n,m$ : Blinded(input<sub>n</sub>)  $\land$  Depends(output<sub>m</sub> on input<sub>n</sub>)

**Goal: changes in sensitive state never affect non-sensitive state** 



# Thinking beyond registers and memory

### Taint-propagation rule must consider many different observable outputs

- Registers
- Memory values
- Control flow

### **Taint tracking policy**

Registers/memory have an associated "sensitive" bit ("Blinded") Ideal rule:

Blinded(output<sub>m</sub>)  $\leftarrow \exists n,m$ : Blinded(input<sub>n</sub>)  $\land$  Depends(output<sub>m</sub> on input<sub>n</sub>) **Goal: changes in sensitive state never affect non-sensitive state** 



# Thinking beyond registers and memory

### Taint-propagation rule must consider many different observable outputs

- Registers
- Memory values
- Control flow
- Exceptions
- Memory access patterns

Not all of these outputs can be marked as sensitive

Data flows from sensitive values to "un-markable" outputs must yield a fault

### How to deal with exceptions

### **Examples of data-dependent exceptions:**

- Division by zero
- Floating-point exceptions
- ...

### Instructions must not raise an exception based on data-dependent conditions

#### Solutions:

- Unconditional faults (i.e., division by sensitive values always fails)
- Set a sensitive error flag and continue computation

### **BliMe Architecture**

- 1. Handshake (incl. remote attestation)
- 2. Shared secret key
- 3. Atomic data import (inputs)
  - Decrypt & blind (Blinded  $\leftarrow$  true)
- 4. Safe ("blinded") computation
  - Enforced by BliMe HW extensions
- 5. Atomic data export (result)
  - Encrypt & unblind (Blinded ← false)



### **BliMe-BOOM Implementation**



### On speculative OoO RISC-V core BOOM

Tagged memory: each byte can be marked as blinded Instructions to mark physical memory as

• Blinded or non-Blinded

### Implements taint-tracking for all instructions

Ideal rule: Blinded(output<sub>m</sub>) ← ∃n,m: Blinded(input<sub>n</sub>) ∧ Depends(output<sub>m</sub> on input<sub>n</sub>)
 Approx. to: Blinded(outputs) ← Blinded(input<sub>1</sub>) ∨ Blinded(input<sub>2</sub>) ∨ ...

instr

### Approximation can be overridden for specific instructions



# **Encryption Engine**

### Encryption engine uses the RoCC accelerator interface in BOOM

• RoCC exposes custom logic as instructions



### Handling multiple clients simultaneously

#### **Problem: So far, one Blinded bit for many clients**

Server can send sensitive data to the wrong client

#### We need a separate sensitivity domain for each client

- Prevent clients accessing each other's sensitive data
- Keys need to be swapped in and out for each client



### Handling multiple clients simultaneously

### Solution 1: BliMe-BOOM-1 + Isolation by honest-but-curious server OS

- OS keeps track of sensitivity domains
- Requires only single Blinded bit from HW: low memory overhead
- Rely on remote attestation of the entire OS to convince client

### Solution 2: BliMe-BOOM-N -- Hardware support for multiple clients

- Hardware keeps track of sensitivity domains: multibit Blindedness tag
- Secure despite malicious OS
- Needs extra memory/logic to keep track of domain identifier for each granule

# **BliMe-BOOM-N Implementation**

**BOOM RTL** 

Data tagged with client-specific tag

1 tag per granule

Tag size = 8 bits, granule size = 8 bytes

Future work:

• parameterize tag size and granule size





#### **Compatibility:** Tested with side-channel-resistant crypto library (TweetNaCI)

• Side-channel-resistant crypto runs without modifications

#### **Overheads**

TypeBOOM-1BOOM-8LUTs & Registers+4.0%+9.0%Power+0.9%+1.4%Max clock frequencyNo reductionNo reduction

**FPGA** 

#### Performance: SPEC2017

BOOM-1	BOOM-8
+23%	+23%

### **BliMe-gem5 optimization**

BliMe-BOOM uses same memory request size for data and tags

Using correct request size (1/8<sup>th</sup>) needs extensive changes to baseline

Solution: Use gem5 simulator to perform evaluation with correct size

• BliMe-gem5-optimized

Could  $\Delta$  in performance just be caused by moving to gem5?

- Implement BliMe-gem5 with BliMe-BOOM configs
- BliMe-gem5 matches BliMe-BOOM in average performance (SPEC 2017)



#### **Compatibility:** Tested with side-channel-resistant crypto library (TweetNaCI)

• Side-channel-resistant crypto runs without modifications

#### **Overheads**



**FPGA** 

#### **Performance: SPEC2017**



**Performance:** 8% average overhead on gem5 after optimization

### **Security: Formal verification in F\***

#### Goal: changes in blinded state never affect non-blinded state

let equivalent\_inputs\_yield\_equivalent\_states (exec:execution\_unit) (pre1 pre2 : systemState) =
 equiv\_system pre1 pre2 ⇒ equiv\_system (step exec pre1) (step exec pre2)

# **Generating compliant code with LLVM**

#### **Problem: software might not run as-is**

• BliMe hardware extensions will abort non-compliant code

### Creating compliant code by hand is error prone

- High-level verification often insufficient
- Challenge exacerbated due to obtuse compiler behavior
- Usability/deployability challenge, not security

### Challenge: solutions like Constantine<sup>[B+21]</sup> are not applicable as-is

• Uses dynamic profiling; under-approximates taint (best-effort approach)

[B+21] "Constantine: Automatic Side-Channel Resistance Using Efficient Control and Data Flow Linearization", ACM CCS (2021)

TensorFlow Lite handported to run on BliMe

 Dolma: BliMe for ML accelerators

 HW accelerators common in outsourced ML workloads

 Genmini is a prominent RISC-V ML accelerator framework

 • Main component: systolic array

 Doma adapts BliMe taint-fracking policy to Gemmini

 • Optimized tag propagation for systolic array

 Formally verified in F\*

 ELOURAY Data Collecter ML Accelerators ung Nuclease Beauly Economics; HORT 2004 (place table)2015(1505)

ompiler support: Improving usability/deploy

aress: summar

ardware improvements: Implementing tag cache

valuation: Experimenting with more TensorFlow models on BliMe

### Generating compliant code with LLVM: our solution

#### Solution: Use static analysis to propagate taint

• Trade-off: over-approximation

### Use SVF<sup>[S+16]</sup> as a starting point

#### SVF provides static value-flow graph

• Shows value dependencies within program

### **Identify and transform potential violations**

• Apply data- and control-flow linearization

[S+16] "SVF: interprocedural static value-flow analysis in LLVM", ACM International Conference on Compiler Construction (2016)

### Adapting TensorFlow Lite to BliMe

#### Summary Bible provides FHE-style security, but efficiently Server can safely run untrusted code on sensitive data Incorporated into speculative OoO RISC-V core BOON In progress: compiler support, tag cache, TensorFlow Paper, source code etc. at <u>https://ssa.research.aithub.loblime</u>/

### **Compiled image classification example**

### Some manual fixes required in TensorFlow Lite library source code

• e.g., array access expansion for softmax lookup table

### In progress:

- For TensorFlow Lite: try more example models
- For the compiler
  - Ensure soundness
  - Produce warnings for untransformed libraries

### **Dolma: BliMe for ML accelerators**

#### HW accelerators common in outsourced ML workloads

#### **Gemmini** is a prominent RISC-V ML accelerator framework

• Main component: systolic array

#### **Dolma adapts BliMe taint-tracking policy to Gemmini**

- Ensures accelerator cannot leak blinded data
- Optimized tag propagation for systolic array



#### Formally verified in F\*

### In progress: summary

**Compiler support:** Improving usability/deployability

Hardware improvements: Implementing tag cache

**Evaluation:** Experimenting with more TensorFlow models on BliMe



BliMe provides FHE-style security, but efficiently

Server can safely run untrusted code on sensitive data

Incorporated into speculative OoO RISC-V core BOOM

In progress: compiler support, tag cache, TensorFlow

Paper, source code etc. at <u>https://ssg-research.github.io/blime/</u>





BliMe provides FHE-style security, but efficiently

Server can safely run untrusted code on sensitive data

Incorporated into speculative OoO RISC-V core BOOM

In progress: compiler support, tag cache, TensorFlow

Paper, source code etc. at <u>https://ssg-research.github.io/blime/</u>

If this type of work interests you, come work with us! https://asokan.org/asokan/research/SecureSystems-open-positions-Jan2024.php

