

Securing Cloud-assisted Services

N. Asokan

http://asokan.org/asokan/@nasokan

Services are moving to "the cloud"



http://dilbert.com/strip/2012-05-25

Services are moving to "the cloud"

Example: cloud-based malware scanning service Example: cloud storage

. . .

Cloud-based malware scanning service

Needs to learn about apps installed on client devices Can therefore infer personal characteristics of users

Predicting User Traits From a Snapshot of Apps Installed on a Smartphone

> Suranga Seneviratne^{a,b} suranga.seneviratne@nicta.com.au Prasant Mohapatra^c prasant@cs.ucdavis.edu

Aruna Seneviratne^{a,b} aruna.seneviratne@nicta.com.au Anirban Mahanti^b anirban.mahanti@nicta.com.au

^aSchool of EET, University of New South Wales, Australia

^bNICTA, Australia ^cDepartment of Computer Science, University of California, Davis

http://dx.doi.org/10.1145/2636242.2636244

Proceedings of the Tenth International AAAI Conference on Web and Social Media (ICWSM 2016)

You Are What Apps You Use: Demographic Prediction Based on User's Apps

Eric Malmi Verto Analytics and Aalto University Espoo, Finland eric.malmi@aalto.fi **Ingmar Weber** Qatar Computing Research Institute Doha, Qatar iweber@qf.org.qa

http://www.aaai.org/ocs/index.php/ICWSM/ICWSM16/paper/view/13047

Securing cloud storage

Client-side encryption of user data is desirable

But naïve client-side encryption conflicts with

- Storage provider's business requirement: deduplication ([LPA15] ACM CCS '15)
- End user's usability requirement: multi-device access ([P+18] IEEE IC '18, CeBIT '16)



New privacy and security concerns arise

Example: cloud-based malware scanning service Example: cloud storage

Naïve solutions conflict with other requirements

• privacy, usability, deployability

CloSer project: the big picture

Cloud Security Services

• 2014-2016, funded by Academy of Finland

Security

- 2016-2018, funded by Tekes
- Academics collaborating with Industry





https://wiki.aalto.fi/display/CloSeProject/CloSer+Project+Public+Homepage



Deployability/Cost





The Circle Game: Scalable Private Membership Test Using Trusted Hardware

Sandeep Tamrakar¹ Jan-Erik Ekberg² Jian Liu ¹ Benny Pinkas ³ Andrew Paverd¹

<u>N. Asokan¹</u>

1. Aalto University, Finland

2. Huawei (work done while at Trustonic)

3. Bar-Ilan University, Israel

Malware checking



On-device checking

- High communication and computation costs
- Database changes frequently
- Database is revealed to everyone

Cloud-based checking

- Minimal communication and computation costs
- Database can change frequently
- Database is not revealed to everyone
- User privacy at risk!

Private Membership Test (PMT)

The problem: How to preserve end user privacy when querying cloud-hosted databases?



Server must not learn contents of client query (q).

Current solutions (e.g. private set intersection, private information retrieval):

- Single server: expensive in both computation and/or communication
- Multiple independent servers: unrealistic in commercial setting

Can hardware-assisted trusted execution environments provide a practical solution? 10

Trusted Execution Environments are pervasive



Hardware support for

- Isolated execution: Trusted Execution Environment
- Protected storage: Sealing
- Ability to report status to a remote verifier: Remote Attestation



[EKA14] "<u>Untapped potential of trusted execution environments</u>", IEEE S&P Magazine, 12:04 (2014)

Background: Kinibi on ARM TrustZone



Kinibi

• Trusted OS from Trustonic

Remote attestation

• Establish a trusted channel

Private memory

- Confidentiality
- Integrity
- Obliviousness

Background: Intel SGX



Trusted Untrusted

CPU enforced TEE (enclave)

Remote attestation

Secure memory

- Confidentiality
- Integrity

Obliviousness only within 4 KB page granularity

System model





Stefanov et al. ACM CCS 2013, https://dl.acm.org/citation.cfm?id=2516660

Android app landscape



On average a user installs 95 apps (Yahoo Aviate) Yahoo Aviate study Source:

https://yahooaviate.tumblr.com/image/95795838933

Unique new Android malware samples

Source: G Data

2015: https://secure.gd/dl-en-mmwr201504 2018: https://www.gdatasoftware.com/blog/2018/02/30491-some-343new-android-malware-samples-every-hour-in-2017

Current dictionary size $< 2^{24}$ entries

Even comparatively "high" FPR (e.g., ~2⁻¹⁰) may have negligible impact on privacy

Cloud-scale PMT

Verify Apps: cloud-based service to check for harmful Android apps prior to installation

"... over 1 billion devices protected by Google's security services, and over 400 million device security scans were conducted per day" Android Security 2015 Year in Review

"2 billion+ Android devices checked per day" https://www.android.com/security-center/

(c.f. < 17 million malware samples)



Requirements

Query Privacy: Adversary cannot learn/infer query or response content

• User can always choose to reveal query content

Accuracy: No false negatives

• However, some false positives are tolerable (i.e. non-zero false positive rate)

Response Latency: Respond quickly to each query

Server Scalability: Maximize overall throughput (queries per second)

Requirements revisited

Query Privacy: Adversary cannot learn/infer query or response content

• User can always choose to reveal queries

Accuracy: No false negatives

• However, some false positives are tolerable (i.e. non-zero false positive rate)

Response Latency: Respond quickly to each query

Server Scalability: Maximize overall throughput (queries per second)

Dictionary size^{*} = 2^{26} entries (~ 67 million entries)

* parameters suggested by a major anti-malware vendor

 $FPR^* = 2^{-10}$

Latency* ~ 1s

Carousel design pattern



Carousel caveats

- 1. Adversary can measure dictionary processing time
 - Spend equal time processing each dictionary entry
- 2. Adversary can measure query-response time
 - Only respond after one full carousel cycle

Both impact response latency (recall Requirements)

Therefore, aim to minimize carousel cycle time

How to minimize carousel cycle time?

Represent dictionary using efficient data structure

Various existing data structures support membership test:

- Bloom Filter
- Cuckoo hash

Experimental evaluation required for carousel approach

Carousel design pattern



Experimental evaluation

Kinibi on ARM TrustZone

- Samsung Exynos 5250 (Arndale)
- 1.7 GHz dual-core ARM Cortex-A17
- Android 4.2.1
- ARM GCC compiler and Kinibi libraries
- Maximum TA private memory: 1 MB
- Maximum shared memory: 1 MB

Intel SGX

- HP EliteDesk 800 G2 desktop
- 3.2 GHz Intel Core i5 6500 CPU
- 8 GB RAM
- Windows 7 (64 bit), 4 KB page size
- Microsoft C/C++ compiler
- Intel SGX SDK for Windows

Note: Different CPU speeds and architectures

Performance: batch queries



Performance: steady state



Kinibi on ARM TrustZone

Intel SGX

Beyond *breakdown point* query response latency increases over time

Other applications of PMT

Discovery of leaked passwords

Private contact discovery in messaging apps

[KLSAP17] PETS 2017

. . .

Signal private contact discovery, Sep 2017

This is much faster. The above code still iterates across the entire set of registered users, but it only does so once for the entire collection of submitted client contacts. By keeping one big linear scan over the registered user data set, access to unencrypted RAM remains "oblivious," since the OS will simply see the enclave touch every item once for each contact discovery request.

The full linear scan is fairly high latency, but by batching many pending client requests together, it can be high throughput.

Carousel approach: pros and cons

Better scalability (for malware checking)

Possible other application scenarios

Membership test may not be enough

Hardware security guarantees may fail





Oblivious Neural Network Predictions via MiniONN Transformations

N. Asokan

http://asokan.org/asokan/
@nasokan

(Joint work with Jian Liu, Mika Juuti, Yao Lu)



By Source, Fair use, https://en.wikipedia.org/w/index.php?curid=54119040

Cloud-assisted malware lookup: recap



violation of clients' privacy

Machine learning as a service (MLaaS)



violation of clients' privacy

Running predictions on client-side



model theft evasion model inversion

Oblivious Neural Networks (ONN)

Given a neural network, is it possible to make it oblivious?

• server learns nothing about clients' input;

• clients learn nothing about the model.

Example: CryptoNets



- High throughput for batch queries from same client
- High overhead for single queries: 297.5s and 372MB (MNIST dataset)
- Cannot support: high-degree polynomials, comparisons, …

[GDLLNW16] CryptoNets, ICML 2016

FHE: Fully homomorphic encryption (<u>https://en.wikipedia.org/wiki/Homomorphic_encryption</u>)

MiniONN: Overview By Source, Fair use, https://en.wikipedia.org/w/index. Blinded input oblivious protocols **Blinded predictions**

- Low overhead: ~1s
- Support all common neural networks

Example $z = W' \bullet f(W \bullet x + b) + b'$





Skip to performance Core idea: use secret sharing for oblivious computation Z y'° **y**'^s (y'' + y'' = y') $W' \bullet [] + b'$ $\mathbf{X}^{\mathbf{c}}$ X^{†S} $(\mathbf{x'}^{c} + \mathbf{x'}^{s} = \mathbf{x'})$

client & server have shares y^{c} and y^{s} s.t. $y^{s}+y^{c}=y$ Client & server have shares x^{c} and x^{s} s.t. $x^{s}+x^{c}=x$ X^{c} X^{c} X^{s} $(x^{c}+x^{s}=x)$ X^{s} $(x^{c}+x^{s}=x)$

Use efficient cryptographic primitives (2PC, additively homomorphic encryption) ³⁷

Secret sharing initial input **x**







$$x_1^s \coloneqq x_1 - x_1^c, \quad x_2^s \coloneqq x_2 - x_2^c$$



Note that **x**^c is independent of **x**. Can be **pre-chosen**

Oblivious linear transformation $W \bullet x + b$



Oblivious linear transformation: dot-product



Oblivious linear transformation $W \bullet x + b$

$$= \begin{bmatrix} w_{1,1} & w_{1,2} \\ w_{2,1} & w_{2,2} \end{bmatrix} \bullet \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} = \begin{bmatrix} w_{1,1} & w_{1,2} \\ w_{2,1} & w_{2,2} \end{bmatrix} \bullet \begin{bmatrix} x_1^s + x_1^c \\ x_2^s + x_2^c \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$$
$$= \begin{bmatrix} w_{1,1}(x_1^s + x_1^c) + w_{1,2}(x_2^s + x_2^c) + b_1 \\ w_{2,1}(x_1^s + x_1^c) + w_{2,2}(x_2^s + x_2^c) + b_2 \end{bmatrix} = \begin{bmatrix} w_{1,1}x_1^s + w_{1,2}x_2^s + b_1 \\ w_{2,1}x_1^s + w_{2,2}x_2^s + b_2 \end{bmatrix} + \begin{bmatrix} w_{1,1}x_1^s + w_{2,2}x_2^s + b_1 \\ w_{2,1}x_1^s + w_{2,2}x_2^s + b_2 \end{bmatrix}$$

Oblivious linear transformation $W \bullet x + b$

$$= \begin{bmatrix} w_{1,1} & w_{1,2} \\ w_{2,1} & w_{2,2} \end{bmatrix} \bullet \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} = \begin{bmatrix} w_{1,1} & w_{1,2} \\ w_{2,1} & w_{2,2} \end{bmatrix} \bullet \begin{bmatrix} x_1^s + x_1^c \\ x_2^s + x_2^c \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$$
$$= \begin{bmatrix} w_{1,1}(x_1^s + x_1^c) + w_{1,2}(x_2^s + x_2^c) + b_1 \\ w_{2,1}(x_1^s + x_1^c) + w_{2,2}(x_2^s + x_2^c) + b_2 \end{bmatrix} = \begin{bmatrix} w_{1,1}x_1^s + w_{1,2}x_2^s + b_1 + w_{1,1}x_1^c + w_{1,2}x_2^c \\ w_{2,1}x_1^s + w_{2,2}x_2^s + b_2 + w_{2,1}x_1^c + w_{2,2}x_2^c \end{bmatrix}$$
$$= \begin{bmatrix} w_{1,1}x_1^s + w_{1,2}x_2^s + b_1 + u_1 \\ w_{2,1}x_1^s + w_{2,2}x_2^s + b_2 + u_2 \end{bmatrix} + \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} y_1^s \\ y_2^s \end{bmatrix} + \begin{bmatrix} y_1^c \\ y_2^c \end{bmatrix}$$

Recall: use secret sharing for oblivious computation





Oblivious activation/pooling functions f(y)

Piecewise linear functions e.g.,

- ReLU: $x := \max(y, 0)$
- Oblivious ReLU: $x^s + x^c := \max(y^s + y^c, 0)$
 - easily computed obliviously by a garbled circuit

Oblivious activation/pooling functions f(y)

Smooth functions e.g.,

- Sigmoid: $x := 1/(1 + e^{-y})$
- Oblivious sigmoid: $x^s + x^c := 1/(1 + e^{-(y^s + y^c)})$
 - approximate by a piecewise linear function
 - then compute obliviously by a garbled circuit
 - empirically: ~14 segments sufficient



Combining the final result



They can jointly calculate $max(y_1,y_2)$ (for minimizing information leakage)

 $y_1 \coloneqq y_1^s + y_1^c$ $y_2 \coloneqq y_2^s + y_2^c$

 y_{1}^{s}, y_{2}^{s}

Recall: use secret sharing for oblivious computation



Performance (for single queries)

| Model | Latency (s) | Msg sizes (MB) | Loss of accuracy |
|---------------|---------------|----------------|----------------------------------------|
| MNIST/Square | 0.4 (+ 0.88) | 44 (+ 3.6) | none |
| CIFAR-10/ReLU | 472 (+ 72) | 6226 (+ 3046) | none |
| PTB/Sigmoid | 4.39 (+ 13.9) | 474 (+ 86.7) | Less than 0.5% (cross-entropy loss) |

Pre-computation phase timings in parentheses

PTB = Penn Treebank

Skip to End

MiniONN pros and cons

300-700x faster than CryptoNets

Can transform any given neural network to its oblivious variant

Still ~1000x slower than without privacy

Server can no longer filter requests or do sophisticated metering

Assumes online connectivity to server

Reveals structure (but not params) of NN

Using a client-side TEE to vet input



Using a client-side TEE to run the model





MiniONN + policy filtering + advanced metering

- disconnected operation + performance + better model secrecy

- harder to reason about client input privacy

MiniONN: Efficiently transform any given neural network into oblivious form with no/negligible accuracy loss Try at: https://github.com/SSGAalto/minionn

Trusted Computing can help realize improved security and privacy for ML

ML is very fragile in adversarial settings



Conclusions

Cloud-assisted services raise new security/privacy concerns

• But naïve solutions may conflict with privacy, usability, deployability, ...

Cloud-assisted malware scanning

Carousel approach is promising

Generalization to privacy-preserving ML predictions

[TLPEPA17] <u>Circle Game</u>, ASIACCS 2017 [LJLA17] <u>MiniONN</u>, ACM CCS 2017





https://ssg.aalto.fi/