# Apple's Deployment of Homomorphic Encryption at Scale

Rehan Rishi, Haris Mughees RWC 2025 | Apple | 03/26/2025





#### **Business Services in Mail**





#### **Business Services in Mail**





#### **Business Services in Mail**

# Apple's privacy commitment

Client gueries remain private, even from Apple

# Practical system requirements

Meets client constraints

- Small storage A few MegaBytes
- High quality experience

#### Meets server constraints

- High queries per second Tens of thousands
- Small communication < 1 MegaByte



#### Feature Walkthrough: Enhanced Visual Search Improving on-device search with privately tagging landmarks and places of interest









#### Region-of-interest

#### Detector





#### **Region-of-interest**

#### Detector

Landmark Region Of Interest in an image is converted into embeddings



# ML Embeddings



Embedding: A vector of floats of a fixed dimension that retains semantic meaning about the input

#### **Embedding Space**

# ML Embeddings



#### **Embedding Space**



# **Nearest Neighbor Search with Plaintext Embeddings**

**Query Embedding** 



Max dot-product ID

Similarity metric computed using dot product

#### **Database Index Embeddings**















5 Million



#### Devices are storage constrained

# Apple's privacy commitment

Client queries remain private, even from Apple  $\checkmark$ 

# **Practical system requirements**

Meets client constraints

- Small storage A few MegaBytes
- High quality experience

#### Meets server constraints

- High queries per second Tens of thousands
- Small communication < 1 MegaByte





## **Background** Server assisted Nearest Neighbor Search without Query Privacy





# **Embeddings Leak Original Data**

Session 2A: ML and Information Leakage

CCS '20, November 9-13, 2020, Virtual Event, USA

#### Information Leakage in Embedding Models

Congzheng Song Cornell University & Google Brain cs2296@cornell.edu

Ananth Raghunathan Facebook & Google Brain ananthr@cs.stanford.edu

#### **Information Leakage from Embedding in Large Language Models**

Zhipeng Wan<sup>\*1</sup> Anda Cheng<sup>\*1</sup> Yinggui Wang<sup>1</sup> Lei Wang<sup>1</sup>

#### Deep Private-Feature Extraction

Seyed Ali Osia, Ali Taheri, Ali Shahin Shamsabadi, Kleomenis Katevas, Hamed Haddadi, Hamid R. Rabiee

#### **Analyzing Sensitive Information Leakage in Trajectory Embedding Models**

Jiaxin Ding\* jiaxinding@sjtu.edu.cn Shanghai Jiao Tong University

Pan Liu wslp1999@sjtu.edu.cn Shanghai Jiao Tong University

Shichuan Xi\* Xi Shichuan@sjtu.edu.cn Shanghai Jiao Tong University

Xinbing Wang xwang8@sjtu.edu.cn Shanghai Jiao Tong University

Kailong Wu 1473686097@sjtu.edu.cn Shanghai Jiao Tong University

Chenghu Zhou zhouch@lreis.ac.cn Institute of Geographical Science and Natural Resources Research, Chinese Academy of Sciences

Saeed Mahloujifar Princeton University sfar@princeton.edu

**Text Embeddings Reveal (Almost) As Much As Text** 

John X. Morris, Volodymyr Kuleshov, Vitaly Shmatikov, Alexander M. Rush Department of Computer Science **Cornell University** 



## **Candidate Secure Solution** Using homomorphic encryption only







## **Candidate Secure Solution** Using homomorphic encryption only



For 3.2 million MSMARCO dataset , assuming 10-thousand cores

[1]- Henzinger, Alexandra, et al. "Private web search with Tiptoe." Proceedings of the 29th symposium on operating systems principles. 2023



Queries per	Communication		
second	per querv		
909	~17.4 MB		
Tens of thousands	< 1 MegaBvte		

## **Candidate Secure Solution** Using homomorphic encryption only

# Apple's privacy commitment

Client queries remain private, even from Apple  $\checkmark$ 

# **Practical system requirements**

Meets client constraints

- Small storage A few MegaBytes
- High quality experience

#### Meets server constraints

- High queries per second Tens of thousands
- Small communication < 1 MegaByte







## **Problem:** High computation + communication

server performs HE operation for each entry in the database

#### **Our solution**

- Use Differential Privacy to reduce server computation
- Use efficient HE to reduce communication

# educe server computation communication

## **Problem:** High computation + communication

server performs HE operation for each entry in the database

## **Our solution**

Use Differential Privacy to reduce server computation 

Use efficient HE to reduce communication

**Clustering**: standard ML technique that's a hard requirement for an efficient nearest neighbor search





-0.9 1.5

0.5

![](_page_25_Picture_4.jpeg)

![](_page_25_Picture_5.jpeg)

![](_page_26_Picture_1.jpeg)

![](_page_26_Picture_2.jpeg)

client photo

batch of queries.

![](_page_27_Picture_3.jpeg)

#### **Problem:** Revealing cluster might reveal semantic information about the

#### **Observation:** Many users at any given time. We can hide a user query in

- Provides trade-off between efficiency and privacy
- Formally bounds worst privacy leakage through clusters
- Guarantees ( $\epsilon$ , $\delta$ )-differential privacy at the user level w.r.t user's photo library

#### Intuition of our privacy guarantee

![](_page_29_Picture_2.jpeg)

# We selected $\epsilon=0.8, \delta=10^{-9}$

![](_page_29_Figure_5.jpeg)

# Achieving ( $\epsilon$ , $\delta$ )-DP: First Step

![](_page_30_Figure_1.jpeg)

![](_page_30_Figure_3.jpeg)

#### Works in epochs, with many clients

![](_page_30_Picture_5.jpeg)

# Achieving ( $\epsilon, \delta$ )-DP: Second Step

![](_page_31_Picture_1.jpeg)

#### Request

User IP: 123.123.123.123

![](_page_31_Picture_4.jpeg)

# Anonymization Network

![](_page_31_Picture_6.jpeg)

![](_page_31_Picture_7.jpeg)

#### **Oblivious HTTP**

# Achieving ( $\epsilon, \delta$ )-DP: Third Step

![](_page_32_Picture_1.jpeg)

![](_page_32_Figure_2.jpeg)

![](_page_32_Figure_3.jpeg)

![](_page_32_Picture_4.jpeg)

![](_page_32_Picture_5.jpeg)

![](_page_32_Figure_6.jpeg)

33

#### Fake queries

# Achieving ( $\epsilon, \delta$ )-DP: Fourth Step

![](_page_33_Picture_1.jpeg)

![](_page_33_Figure_3.jpeg)

34

Random query schedule

Proof details: Scalable Private Search with Wally (https://arxiv.org/abs/2406.06761)

#### **Proof Intuition:**

1. We show that server view is a noisy histogram of clusters

2. Prove this noisy histogram is  $(\epsilon, \delta)$ -DP in central model

3. The server gains no extra advantage in distributed model

![](_page_34_Figure_8.jpeg)

![](_page_35_Picture_1.jpeg)

![](_page_35_Picture_2.jpeg)

## **Problem:** High computation + communication

server performs HE operation for each entry in the database

#### **Our solution**

Use Differential Privacy to reduce server computation

Use efficient HE to reduce communication

![](_page_37_Picture_1.jpeg)

BFV HE [1]

![](_page_37_Figure_3.jpeg)

![](_page_37_Picture_5.jpeg)

[1] Jean-Claude Bajard, et al., "A Full RNS Variant of FV-like Somewhat Homomorphic Encryption Schemes," International Conference on Selected Areas in Cryptography, 2016

![](_page_37_Picture_7.jpeg)

![](_page_37_Picture_8.jpeg)

![](_page_38_Picture_1.jpeg)

![](_page_38_Picture_2.jpeg)

![](_page_38_Figure_3.jpeg)

#### HE computation

![](_page_38_Picture_5.jpeg)

![](_page_39_Picture_1.jpeg)

![](_page_39_Picture_2.jpeg)

![](_page_39_Figure_3.jpeg)

![](_page_40_Figure_1.jpeg)

Message is encoded in higher order bits

#### RNS based ciphertext of BFV HE

RNS Limb 1

RNS Limb 2

#### - Modulus switching [1]: Keep single RNS limb

![](_page_41_Figure_2.jpeg)

#### - Dropping LSB [2]: Further drop least significant bits from the remaining limb

![](_page_41_Figure_4.jpeg)

[1] Zvika Brakerski and Vinod Vaikuntanathan, "Efficient Fully Homomorphic Encryption from (Standard) LWE," IEEE Symposium on Foundations of Computer Science, 2011 [2] Zhenyu Huang, et al., "Cheetah: Lean and Fast Secure Two-Party Deep Neural Network Inference," USENIX Security Symposium, 2022

![](_page_41_Figure_6.jpeg)

![](_page_41_Picture_7.jpeg)

# **Our Solution** Other Optimizations

Delayed modular reduction to reduce server compute

![](_page_42_Figure_2.jpeg)

**Modular reduction** 

Operation = arithmetic operations in field

 Plaintext RNS to reduce evaluation key size - Evaluation key dominates request size to maintain anonymity

![](_page_42_Figure_7.jpeg)

## **Enhanced Visual Search** Results

	Technique	Queries per second	Communication per query
<b>Our Results*</b>	HE+DP+ Anonymization Network	>25,000	0.56 MB
Tiptoe [1]	Additive HE	909	~17.4 MB

\*Includes overhead due to fake queries For 3.2 million MSMARCO dataset , assuming 10-thousand cores

[1] Alexandra Henzinger, et al., "Private Web Search with Tiptoe," Symposium on Operating Systems Principles, 2023

# **Our Solution**

# Apple's privacy commitment

Client queries remain private, even from Apple

# Practical system requirements

Meets client constraints

- Small storage
- High quality experience

![](_page_44_Picture_7.jpeg)

## Meets server constraints

- High queries per second
- Small communication

![](_page_44_Picture_11.jpeg)

# **Open source Server HE implementation Server Side: Apple Swift Homomorphic Encryption** Auditable 🗸 Novel optimizations Feedback welcome via Github 🗸

https://github.com/apple/swift-homomorphic-encryption

**Device Side HE implementation: Corecrypto** 

Auditable 🗸

Novel optimizations

https://developer.apple.com/security/#corecrypto

![](_page_45_Picture_7.jpeg)

# Summary of Apple's Deployment of Homomorphic Encryption at Scale

- Enhancing on-device experiences with information from the server while maintaining one of the strongest notions of privacy
- Efficient Homomorphic Encryption for several features running on over a billion devices
- Uniquely combined with other state-of-the-art privacy technologies:
  - Differential privacy
  - Anonymization network
  - Privacy pass

# **Public documentation and resources**

#### swift-homomorphic-encryption apple

https://github.com/apple/swift-homomorphic-encryption

<u>SMS and Call Reporting</u> / Getting up-to-date calling and blocking information for your app

#### Article

#### Getting up-to-date calling and blocking information for your app

Implement the Live Caller ID Lookup app extension to provide call-blocking and identity services.

https://developer.apple.com/documentation/identitylookup/ getting-up-to-date-calling-and-blocking-information-for-yourapp

![](_page_47_Picture_8.jpeg)

Highlight | October 24, 2024

Privacy

**Combining Machine** Learning and **Homomorphic Encryption** in the Apple Ecosystem

![](_page_47_Picture_12.jpeg)

XZ

Using Private Nearest Neighbor Search for Enhanced Visual Search for photos

#### https://machinelearning.apple.com/research/homomorphicencryption

## **Scalable Private Search with Wally**

Hilal Asi, Fabian Boemer, Nicholas Genise, Muhammad Haris Mughees, Tabitha Ogilvie, Rehan Rishi, Guy N. Rothblum, Kunal Talwar, Karl Tarbe, Ruiyu Zhu, Marco Zuliani

#### https://machinelearning.apple.com/research/wally-search

![](_page_47_Picture_19.jpeg)

![](_page_48_Picture_0.jpeg)

![](_page_48_Picture_1.jpeg)

TM and © 2025 Apple Inc. All rights reserved.

![](_page_48_Picture_3.jpeg)