

BILL MURRAY SCARLETT JOHANSSON

Found In Translation: A Generative Language Modeling Approach to Memory Access Pattern Attacks

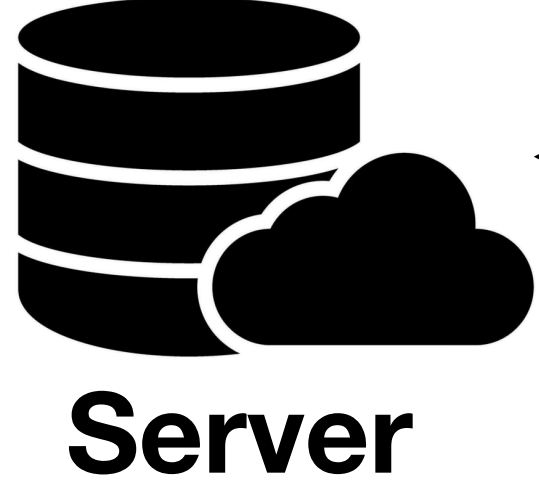
Grace Jia, Alex Wong, Anurag Khandelwal




FOCUS FEATURES PRESENTS AN AMERICAN ZOETROPE / ELEMENTAL FILMS PRODUCTION "LOST IN TRANSLATION" BILL MURRAY SCARLETT JOHANSSON GIOVANNI RIBISI ANNA FARIS FUMIHIRO HAYASHI MUSIC BY BRIAN REITZELL COSTUME DESIGNER NANCY STEINER
RESTRICTED R PARENT STRONGLY CAUTIONED
SOME SEXUAL CONTENT
PRODUCED BY ANNE ROSS K.K. BARRETT WITH SARAH FLACK DIRECTED BY LANCE ACORD PRODUCED BY CALLUM GREENE ASSOCIATED PRODUCERS MITCH GLAZER
EXECUTIVE PRODUCERS FRANCIS FORD COPPOLA FRED ROOS PRODUCED BY ROSS KATZ SOFIA COPPOLA WRITTEN AND DIRECTED BY SOFIA COPPOLA
SOUNDTRACK AVAILABLE ON emperorNorton
FOCUS FEATURES
www.lost-in-translation.com
The new film written and directed by Sofia Coppola

Secure execution on untrusted cloud

Untrusted cloud provider

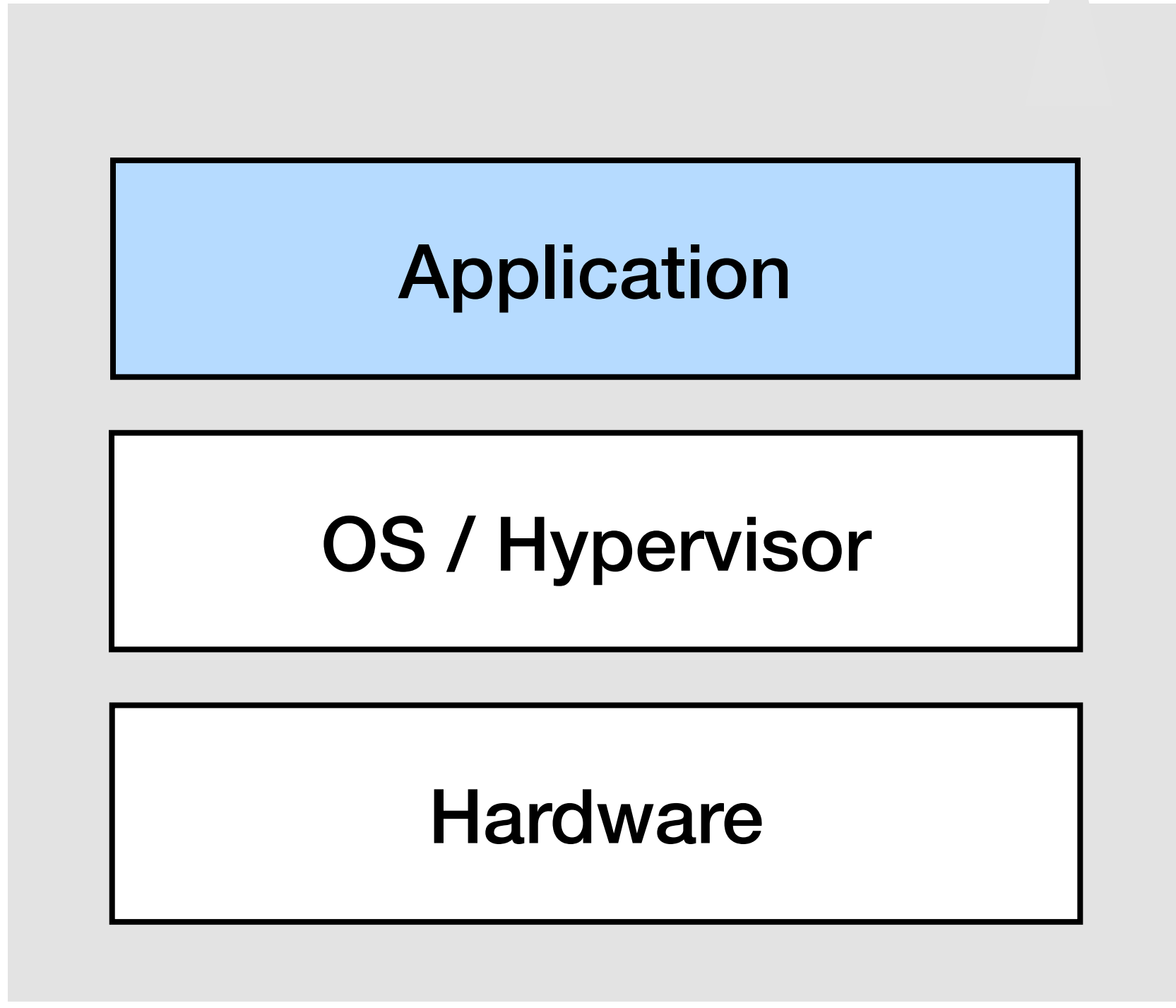


Sensitive input 

Response 



Client



Secure execution on untrusted cloud

Untrusted cloud provider



Server

Sensitive input

Response



Client

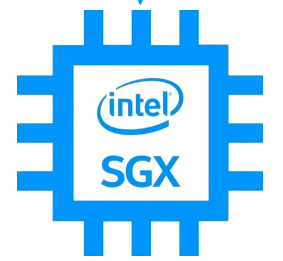
Confidential Computing Environment

Application



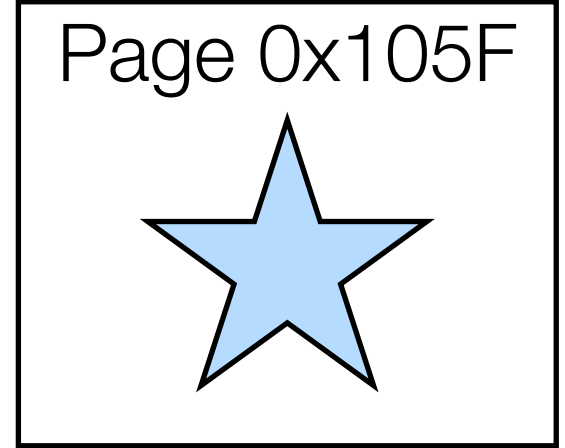
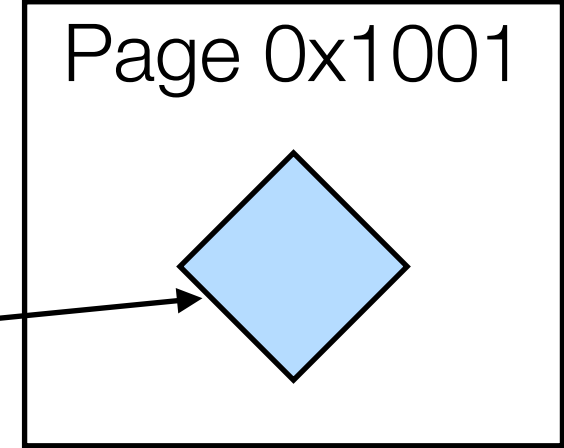
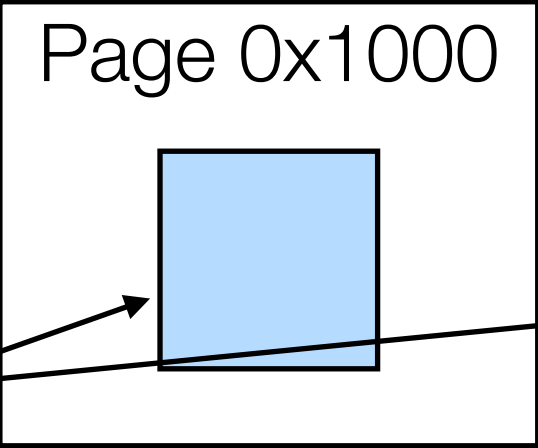
OS / Hypervisor

Hardware

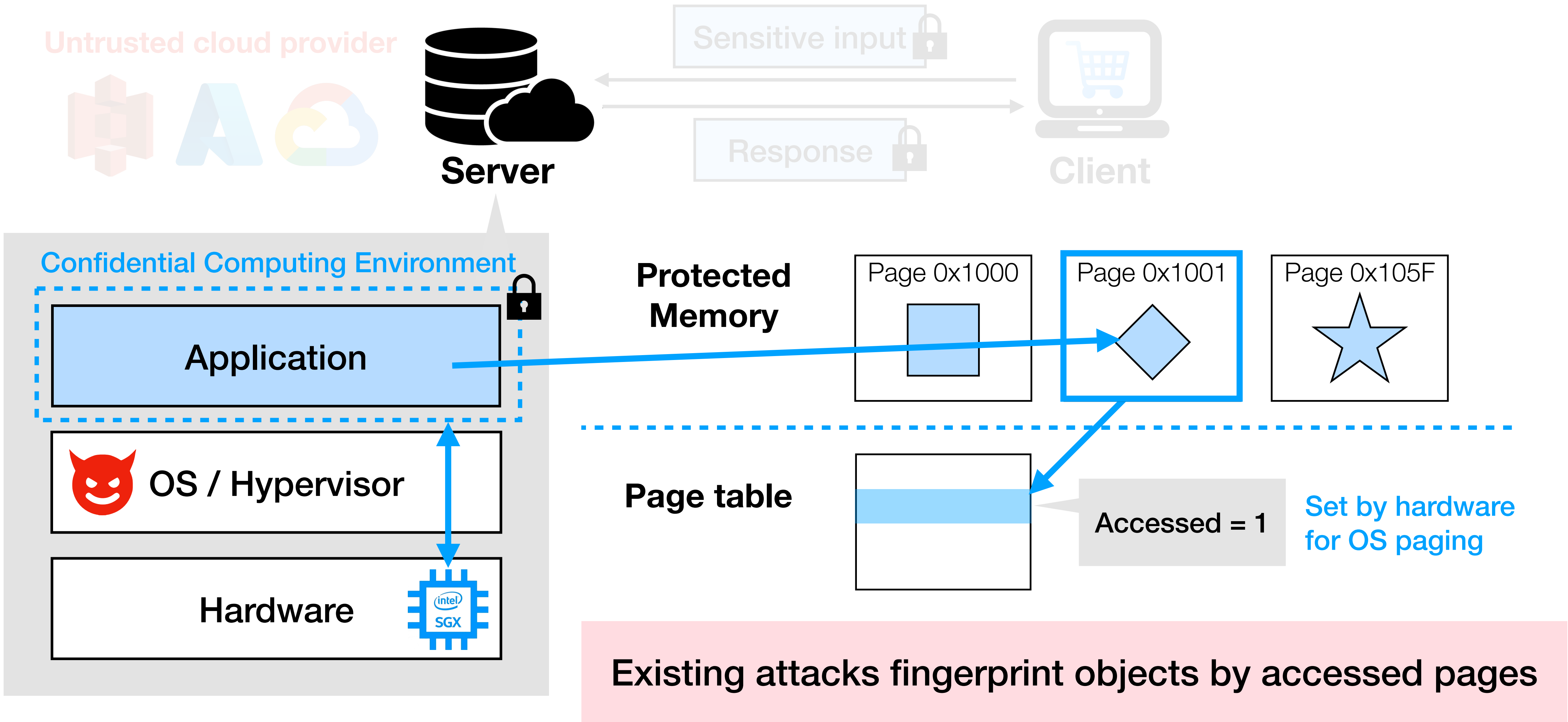


Protected Memory

Memory objects



Secure execution on untrusted cloud

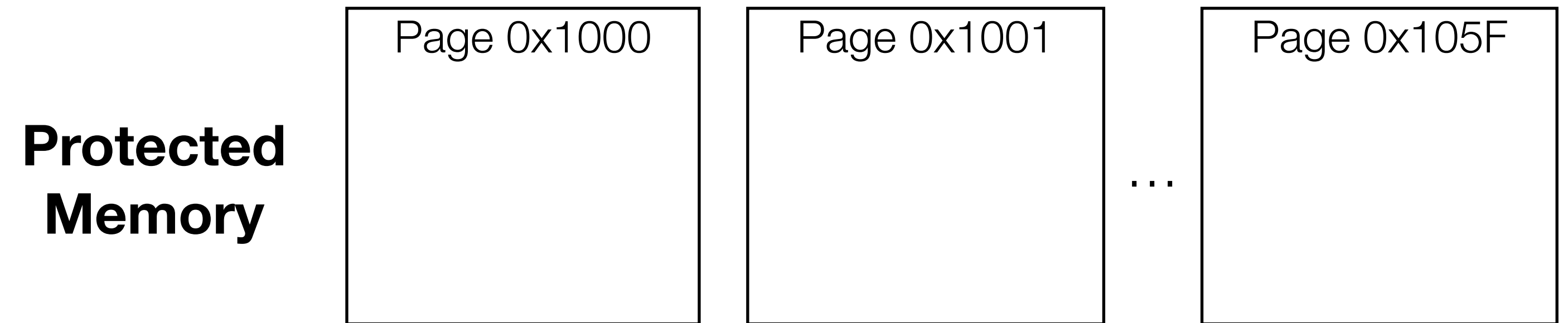


[S&P'15, AsiaCCS'16, Sec'17]

Example: Ad recommendation engine



Server



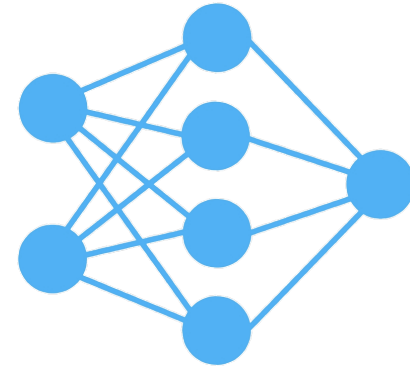
Page table

VPN	0x1000	0x1001	...	0x105F	...
Accessed?	0	0	...	0	...

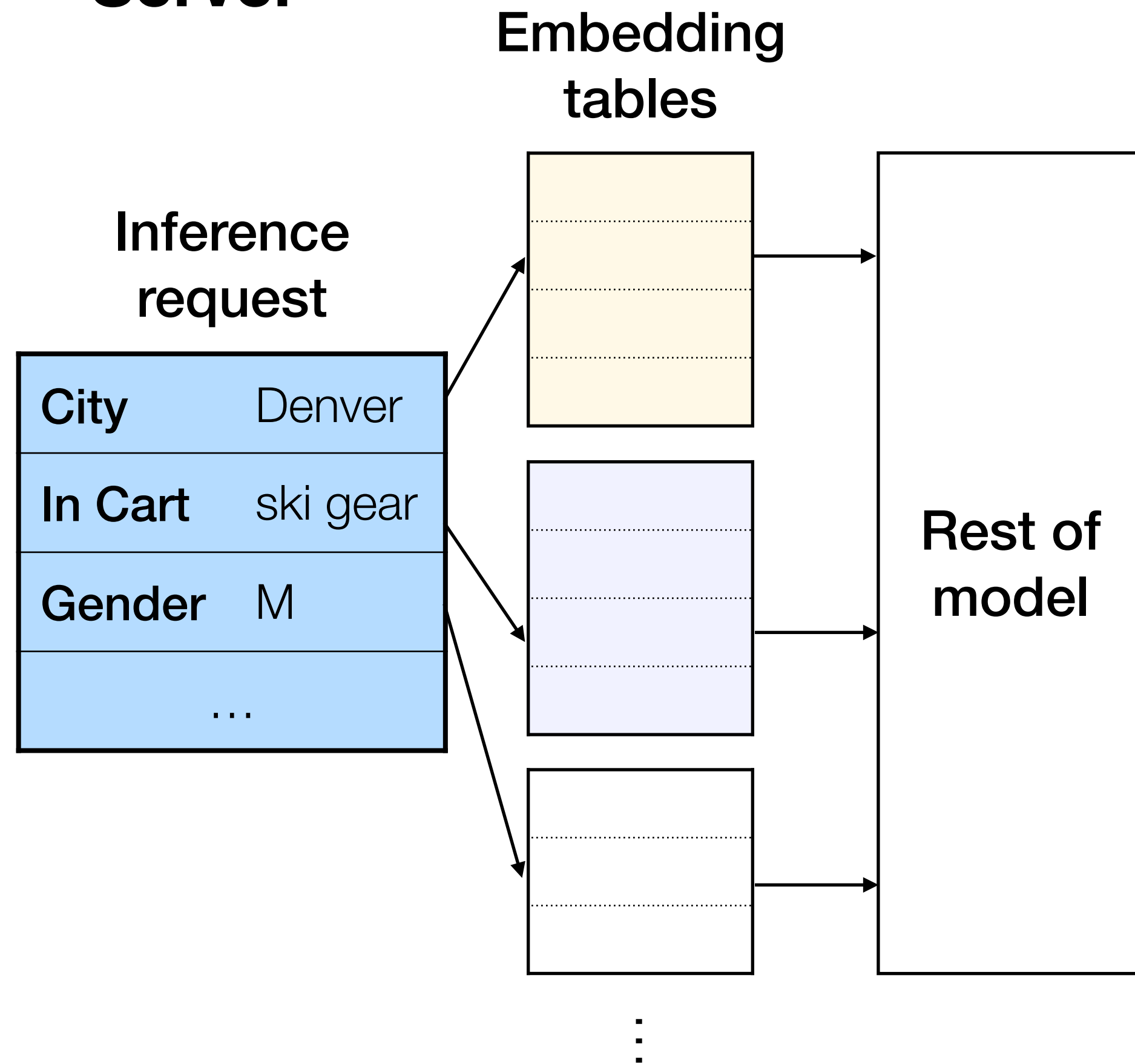
Example: Ad recommendation engine



Server

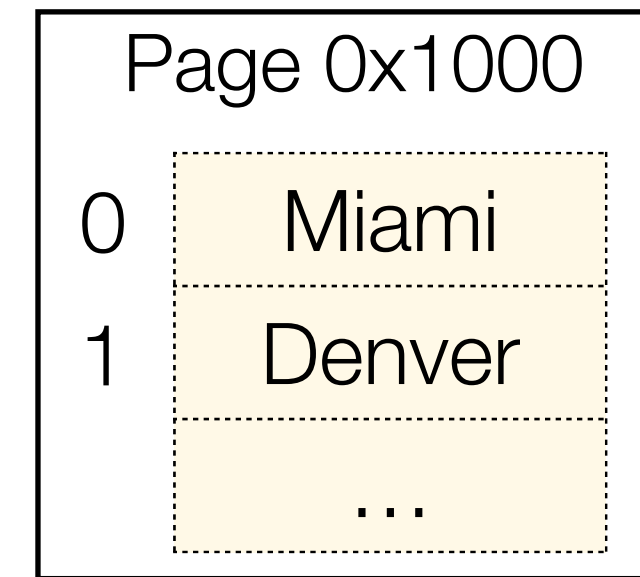


Deep learning recommendation model (DLRM)

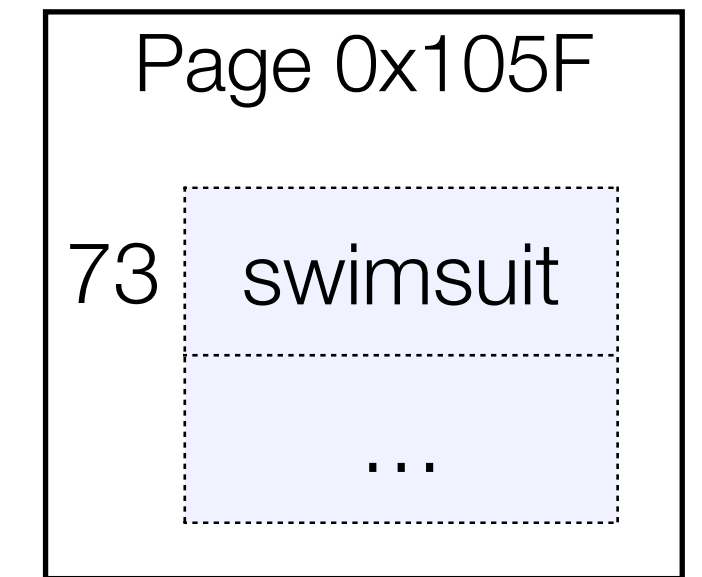
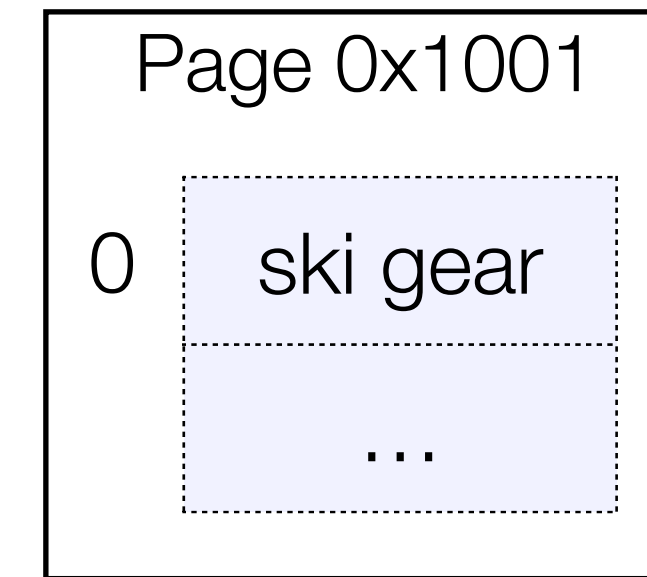


Protected Memory

City



In Cart



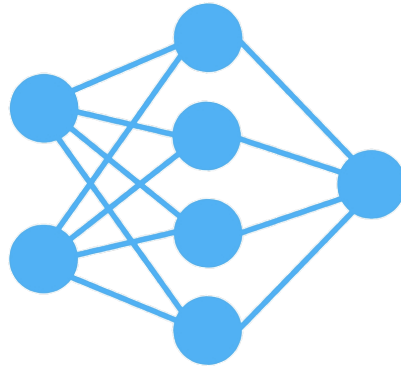
Page table

VPN	0x1000	0x1001	...	0x105F	...
Accessed?	0	0	...	0	...

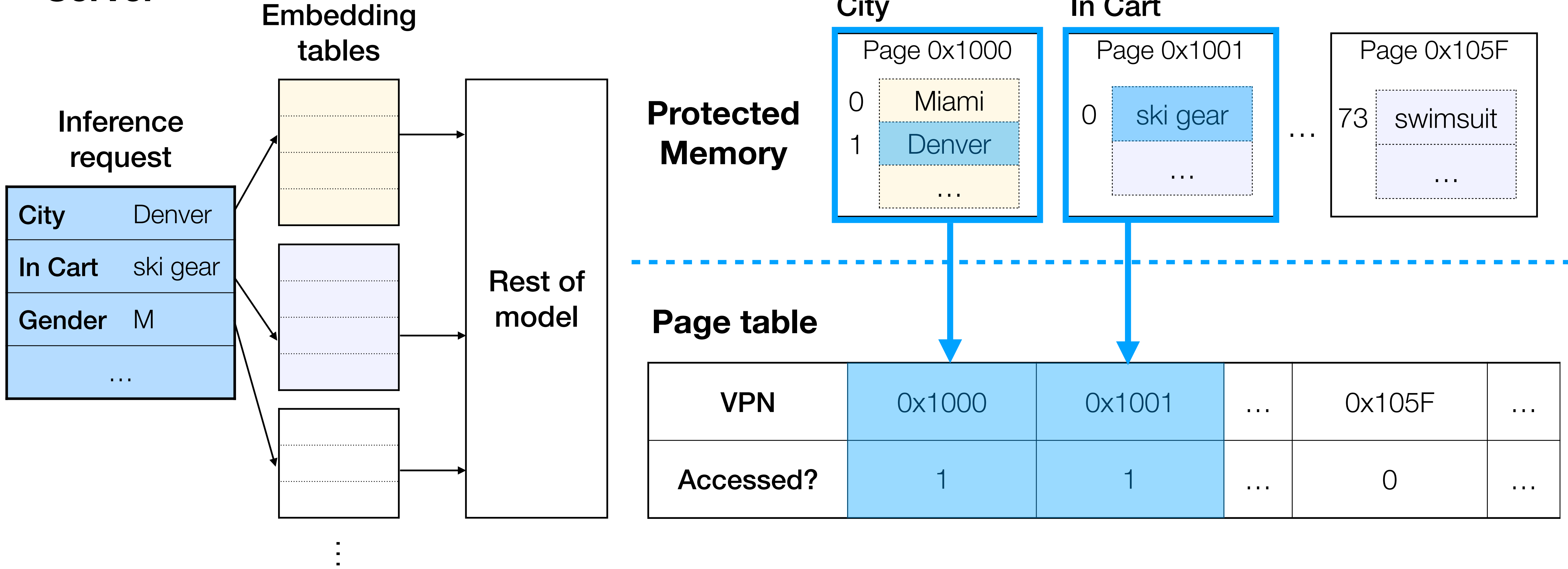
Example: Ad recommendation engine



Server

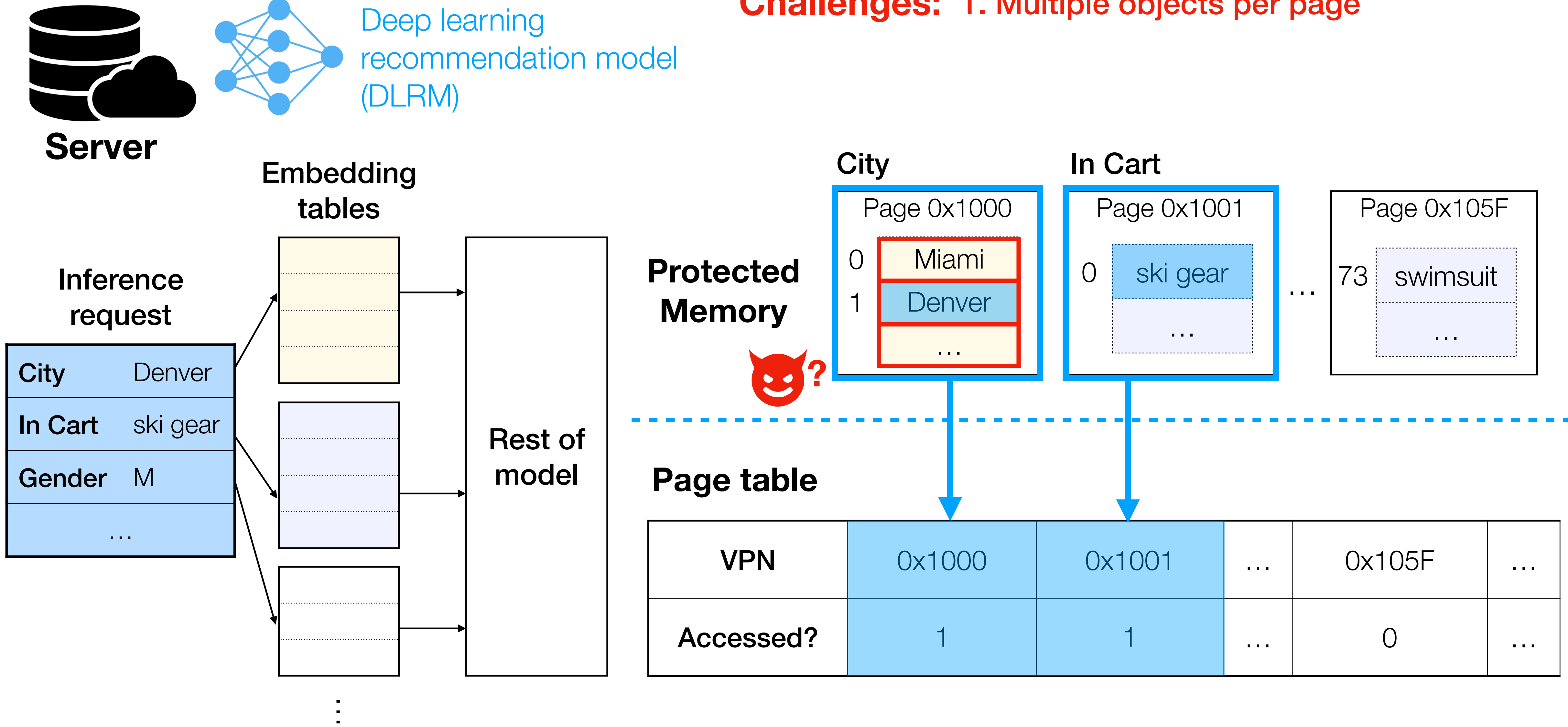


Deep learning recommendation model (DLRM)



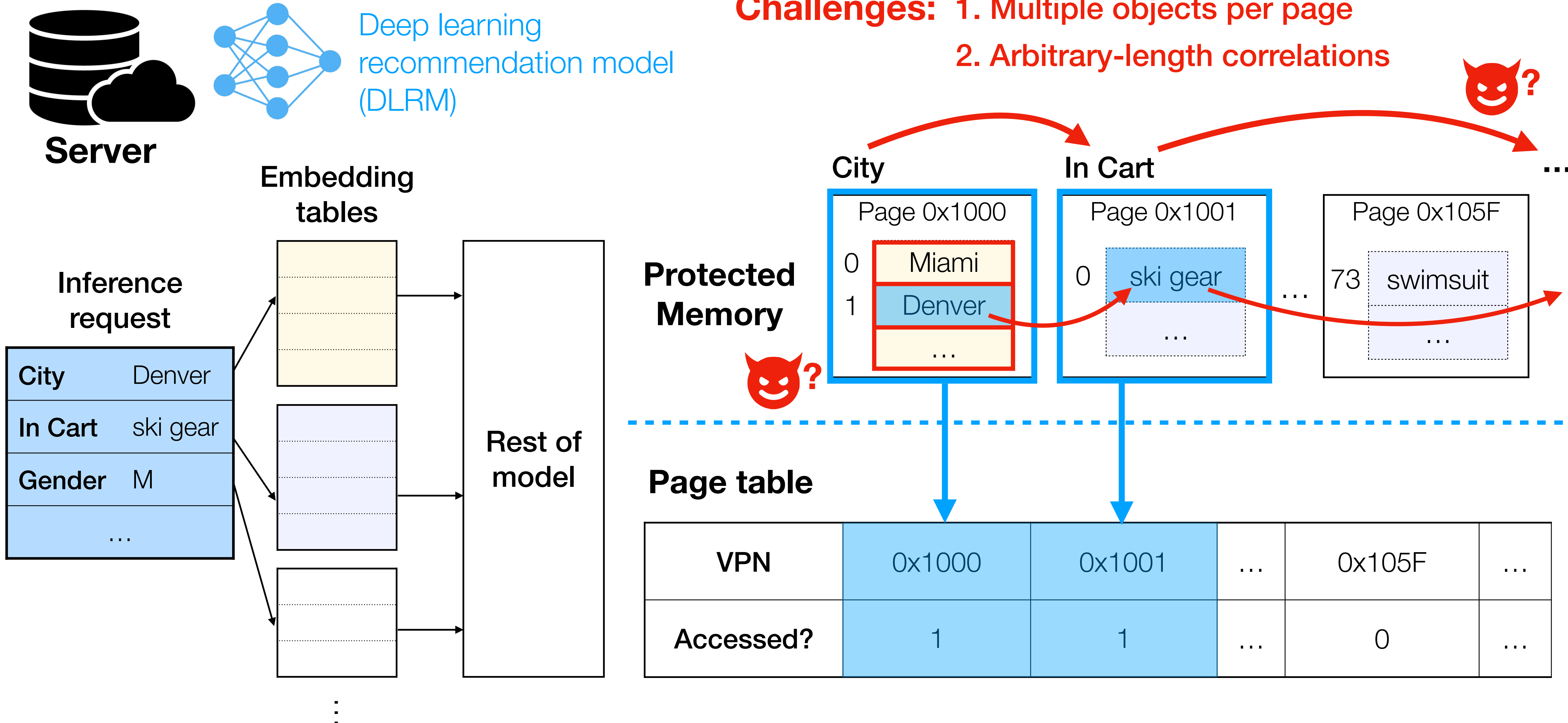
Example: Ad recommendation engine

Challenges: 1. Multiple objects per page



Example: Ad recommendation engine

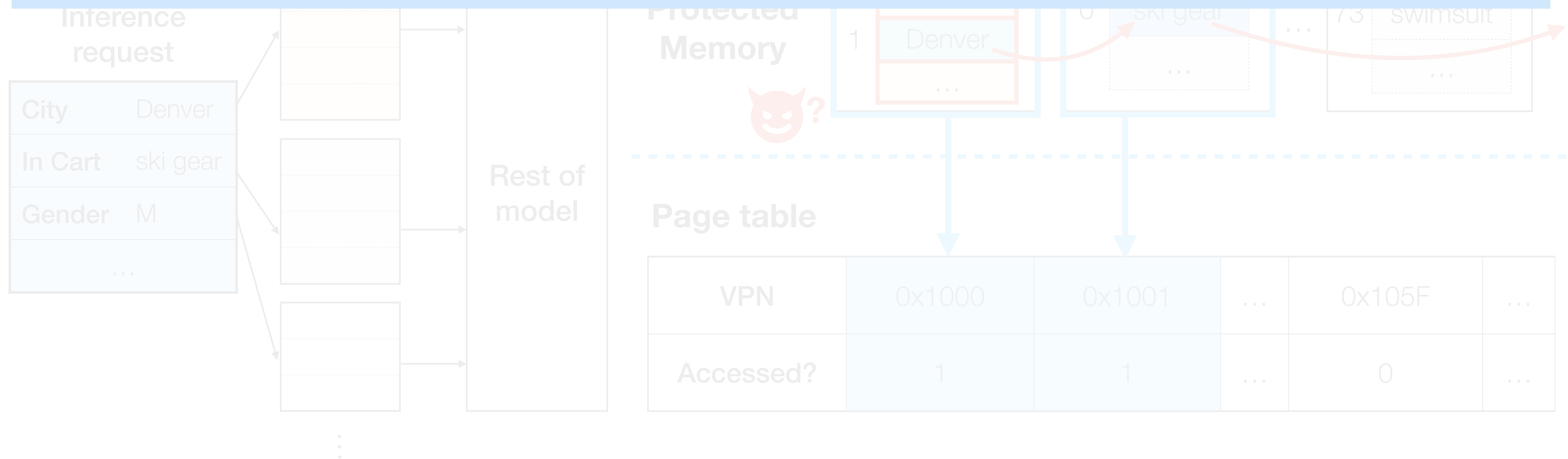
Challenges: 1. Multiple objects per page
2. Arbitrary-length correlations



Example: Ad recommendation engine

- Challenges:**
1. Multiple objects per page
 2. Arbitrary-length correlations

In the face of such challenges, is there a practical attack that can infer object-level accesses?

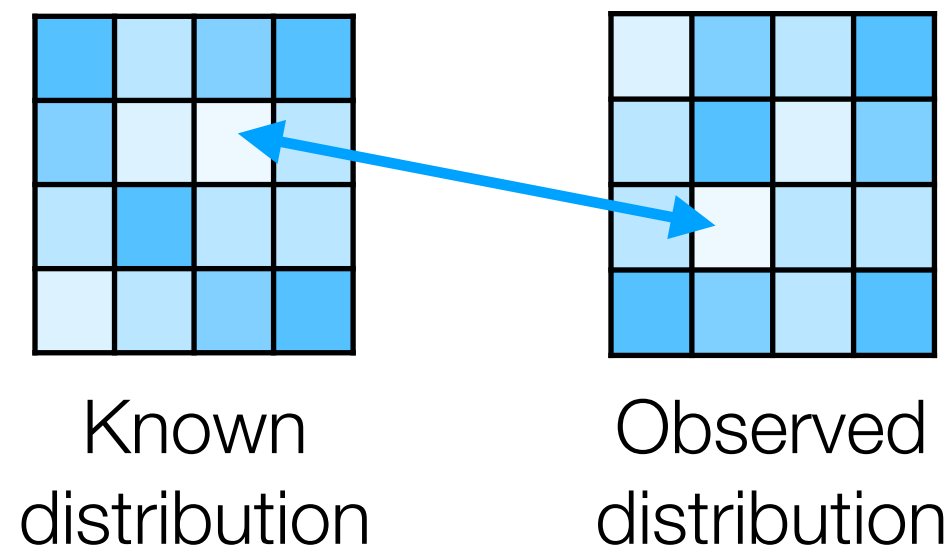


Inferring object accesses from pages

Prior attack approaches in related fields (e.g., SSE)

Markov modeling

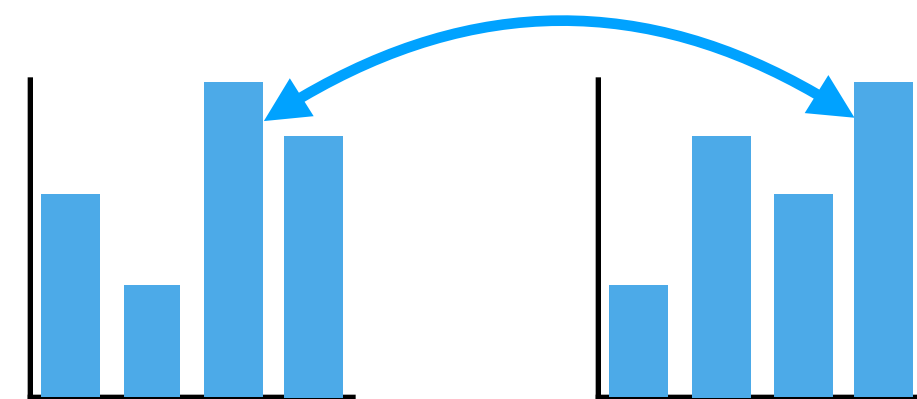
IHOP [Sec'22]



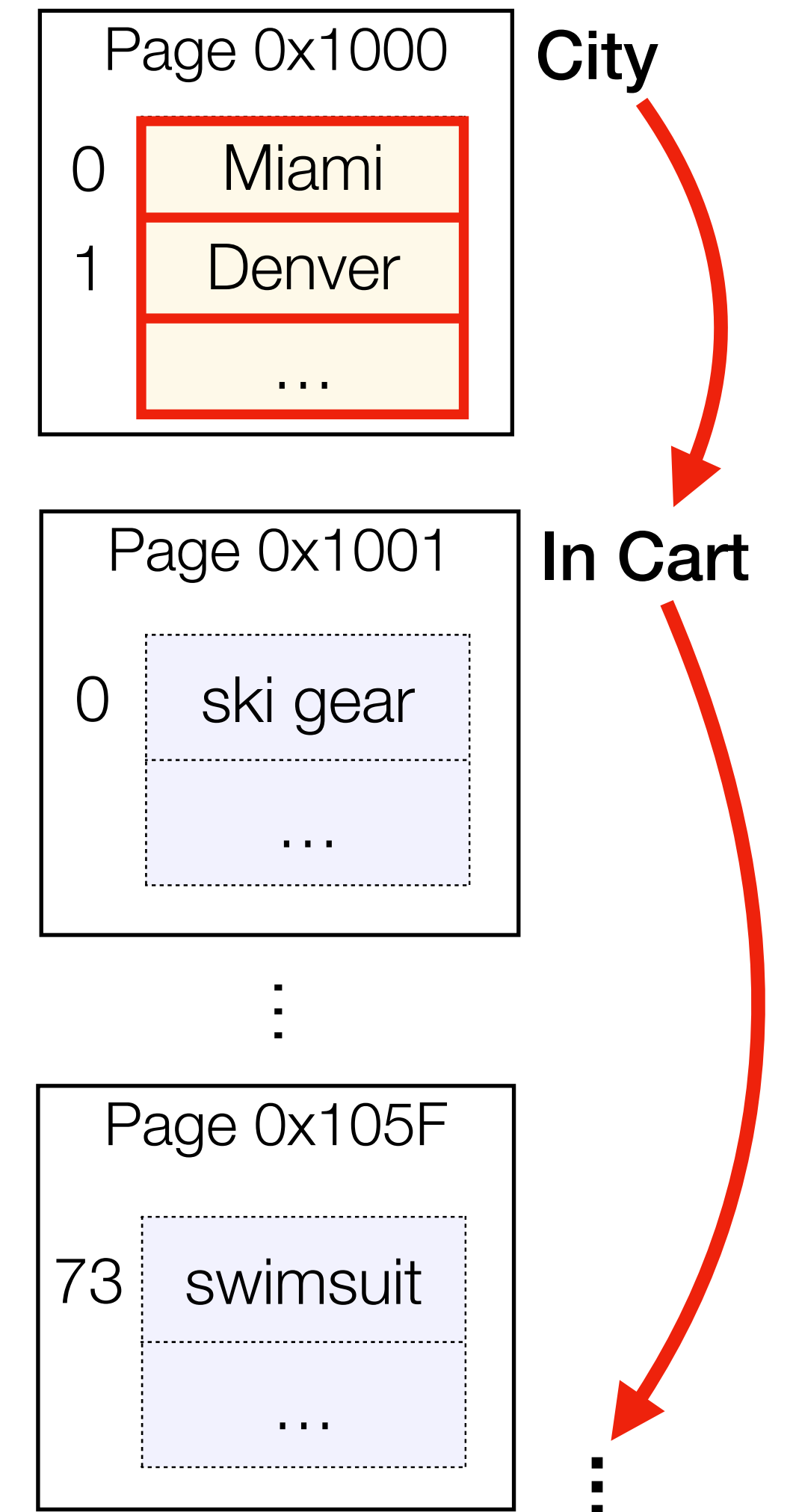
Assumes one object per page

Frequency analysis

Naive Bayes



Assumes independent accesses

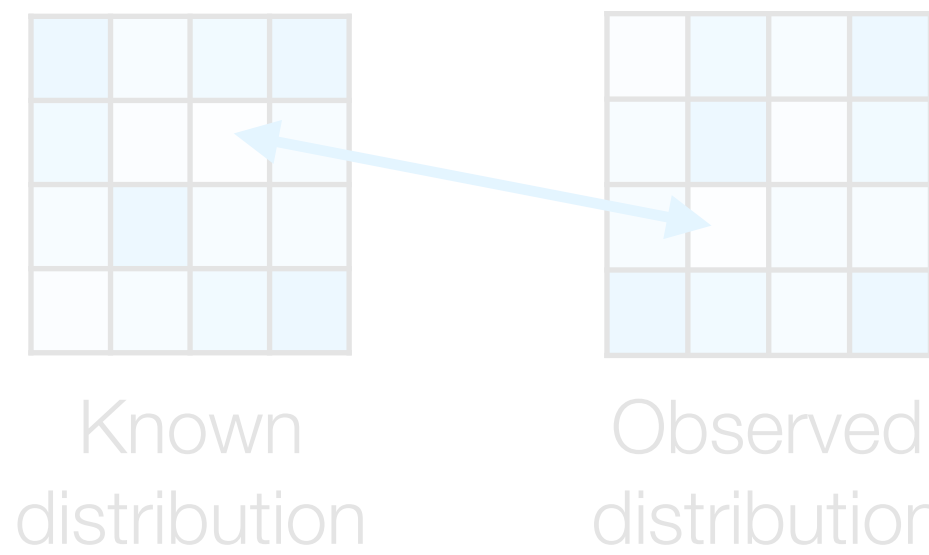


Inferring object accesses from pages

Prior attack approaches in related fields (e.g., SSE)

Markov modeling

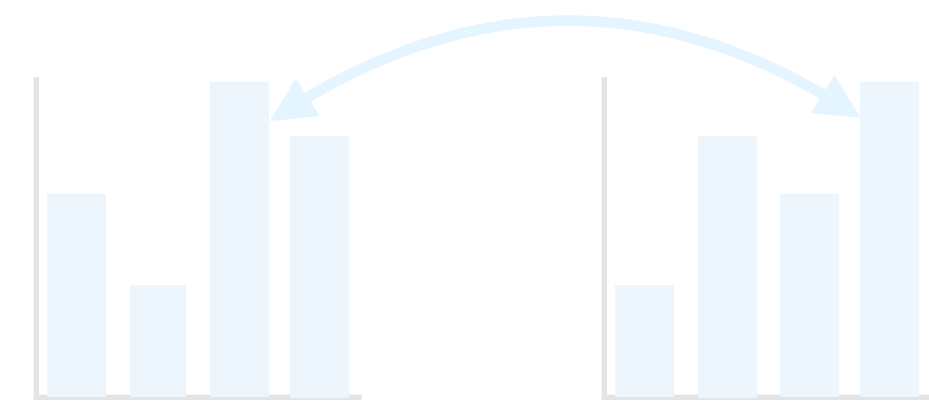
IHOP [Sec'22]



Assumes one object per page

Frequency analysis

Naive Bayes



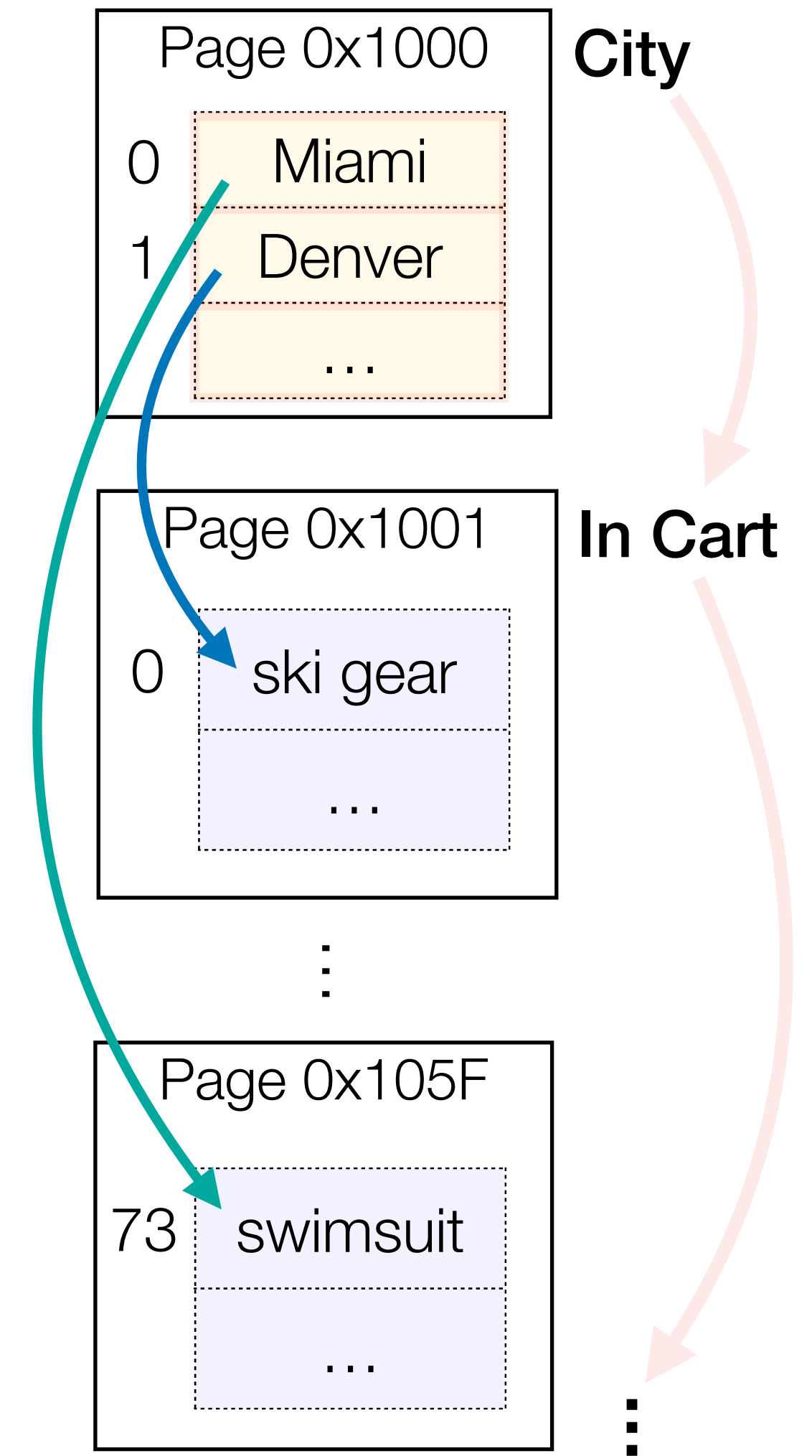
Assumes independent accesses

Our insight: Leverage correlations between objects

Adversary (OS) knows memory layout and access correlations

$0x1000, 0x1001 \implies$ Denver, ski gear

$0x1000, 0x105F \implies$ Miami, swimsuit

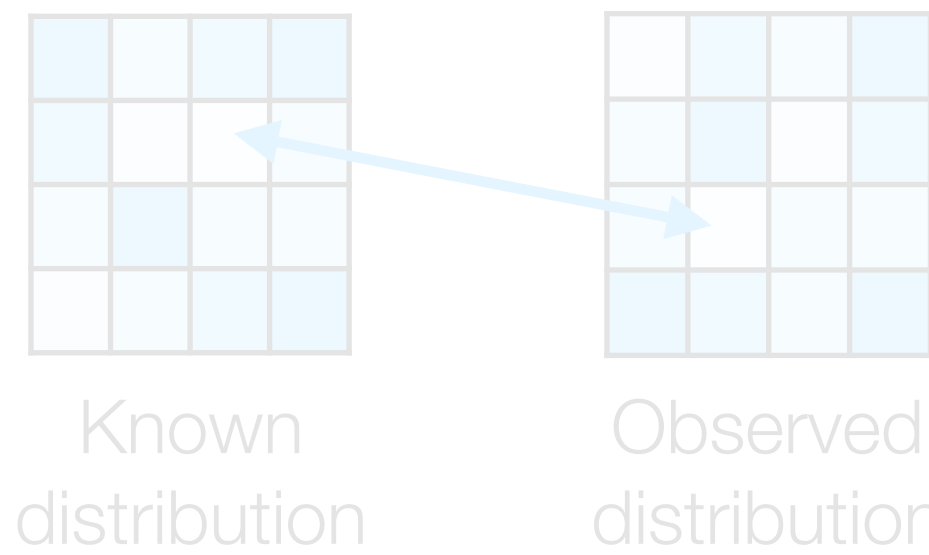


Inferring object accesses from pages

Prior attack approaches in related fields (e.g., SSE)

Markov modeling

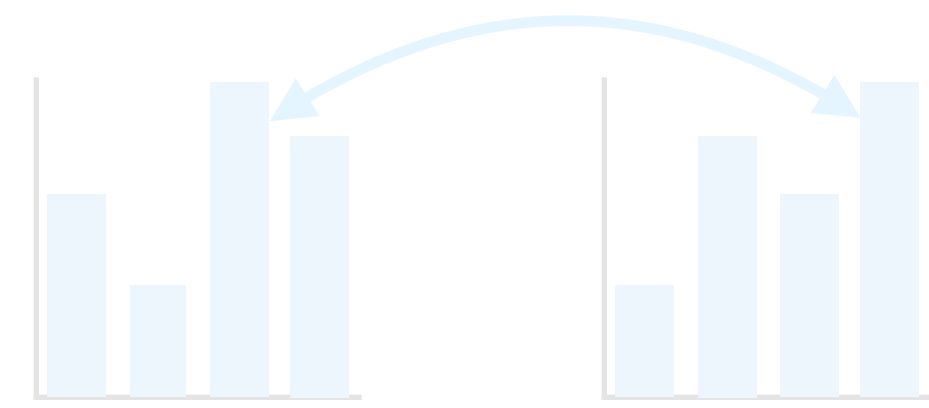
IHOP [Sec'22]



Assumes one object per page

Frequency analysis

Naive Bayes



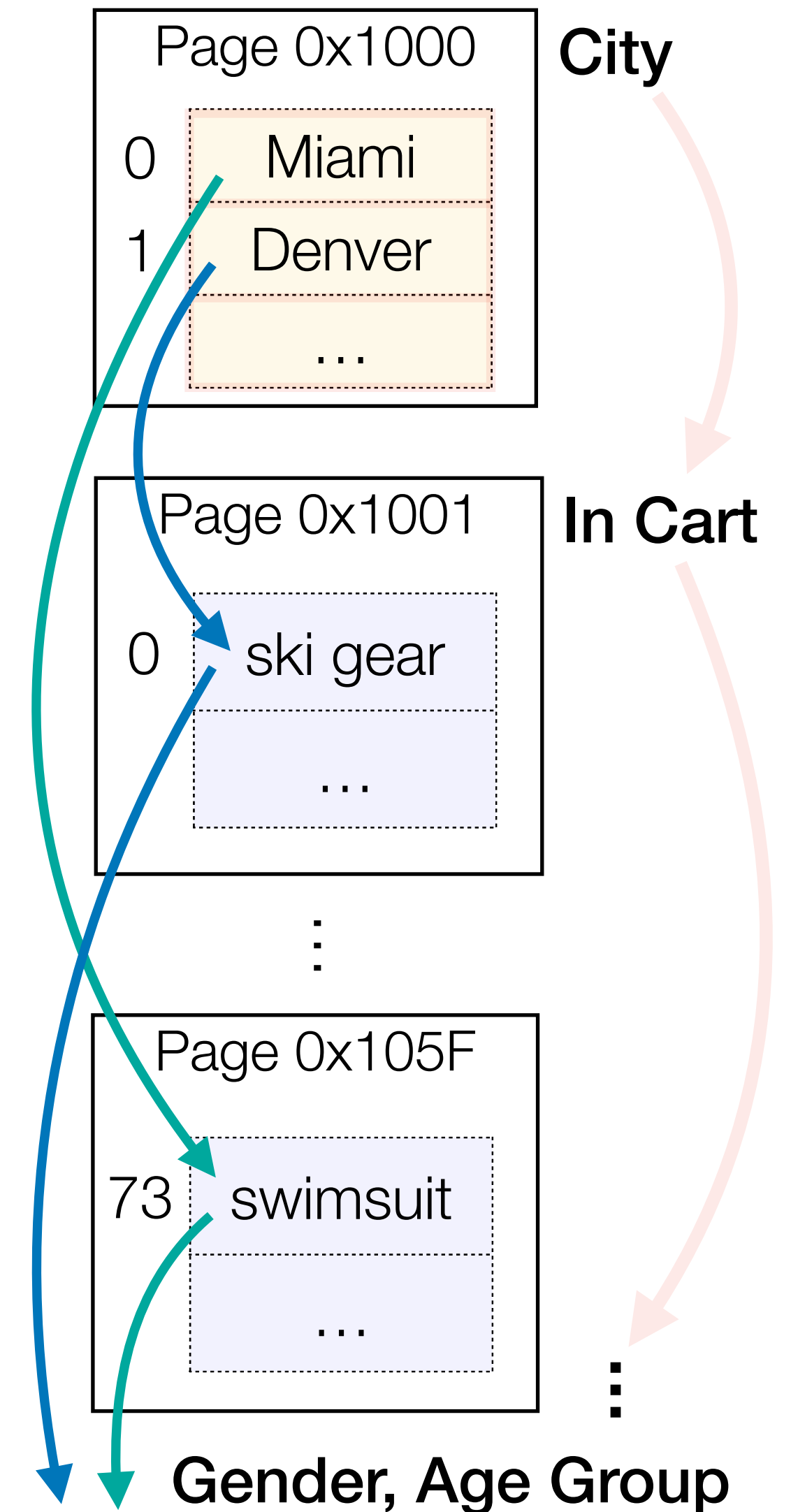
Assumes independent accesses

Our insight: Leverage correlations between objects

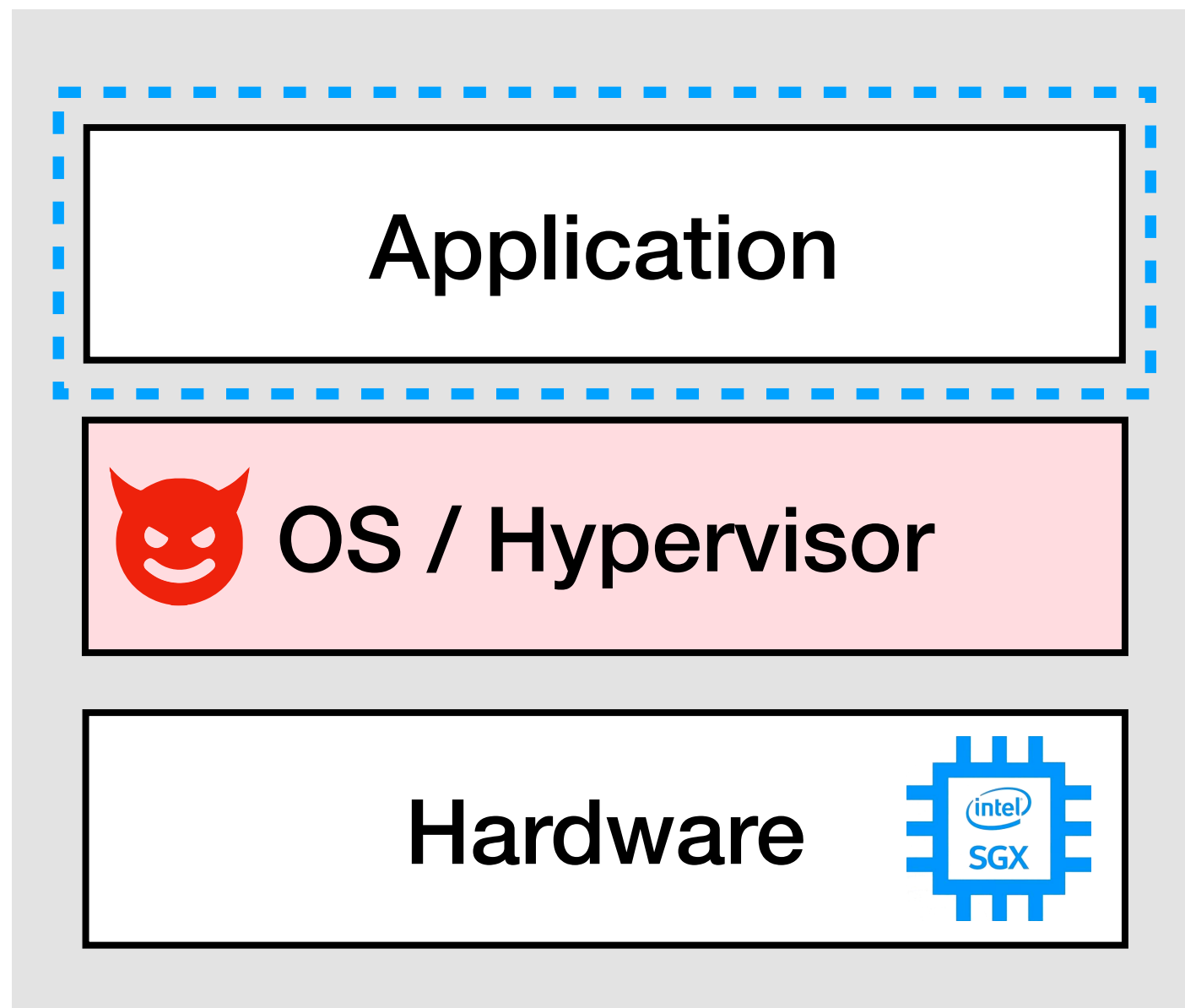
Adversary (OS) knows memory layout and access correlations

$0x1000, 0x1001 \implies$ Denver, ski gear

$0x1000, 0x105F \implies$ Miami, swimsuit



Threat Model



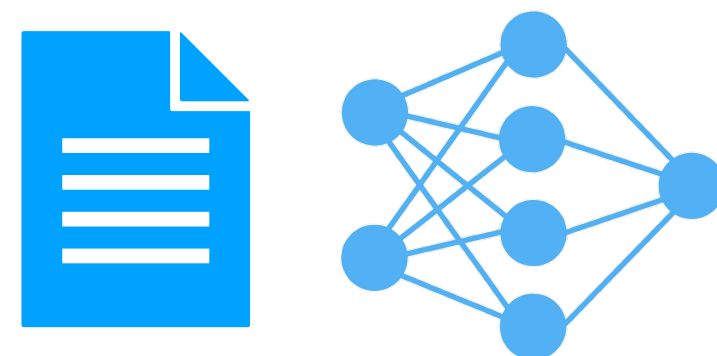
Honest-but-curious cloud adversary

- ▶ Cannot access application data or modify inputs
- ▶ Can observe application's page accesses
- ▶ Can control scheduling, memory allocations, etc.

Goal: Given page accesses, infer object accesses

Auxiliary knowledge:

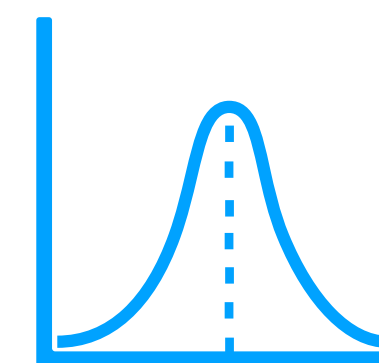
White-box access
to application



Layout of objects
per page



Distribution of
input data



Language modeling attack approach

$\Pr[\text{book} \mid \text{I, wrote, a}]$
> $\Pr[\text{sandwich} \mid \text{I, wrote, a}]$

**Parallels in syntax &
semantic rules**

$\Pr[\text{Miami} \mid \text{swimsuit}]$
> $\Pr[\text{Denver} \mid \text{swimsuit}]$

Language modeling attack approach

$\Pr[\text{book} \mid \text{I, wrote, a}]$ **Parallels in syntax & semantic rules** $\Pr[\text{Miami} \mid \text{swimsuit}]$
> $\Pr[\text{sandwich} \mid \text{I, wrote, a}]$ > $\Pr[\text{Denver} \mid \text{swimsuit}]$

***Translate* from language of page accesses to object accesses**

Learn distribution of object sequence \mathbf{y} conditioned on page sequence \mathbf{x} :

$$\Pr[y_1, \dots, y_n \mid \mathbf{x}] = \prod_{i=1}^n \Pr[y_i \mid y_1, \dots, y_{i-1}, \mathbf{x}]$$

Language modeling attack approach

$\Pr[\text{book} \mid \text{l, wrote, a}]$
> $\Pr[\text{sandwich} \mid \text{l, wrote, a}]$

Parallels in syntax & semantic rules

$\Pr[\text{Miami} \mid \text{swimsuit}]$
> $\Pr[\text{Denver} \mid \text{swimsuit}]$

Translate from language of page accesses to object accesses

Learn distribution of object sequence \mathbf{y} conditioned on page sequence \mathbf{x} :

$$\Pr[y_1, \dots, y_n \mid \mathbf{x}] = \prod_{i=1}^n \Pr[y_i \mid y_1, \dots, y_{i-1}, \mathbf{x}]$$

right hand \implies main **droite**
that is **right** \implies c'est **vrai**

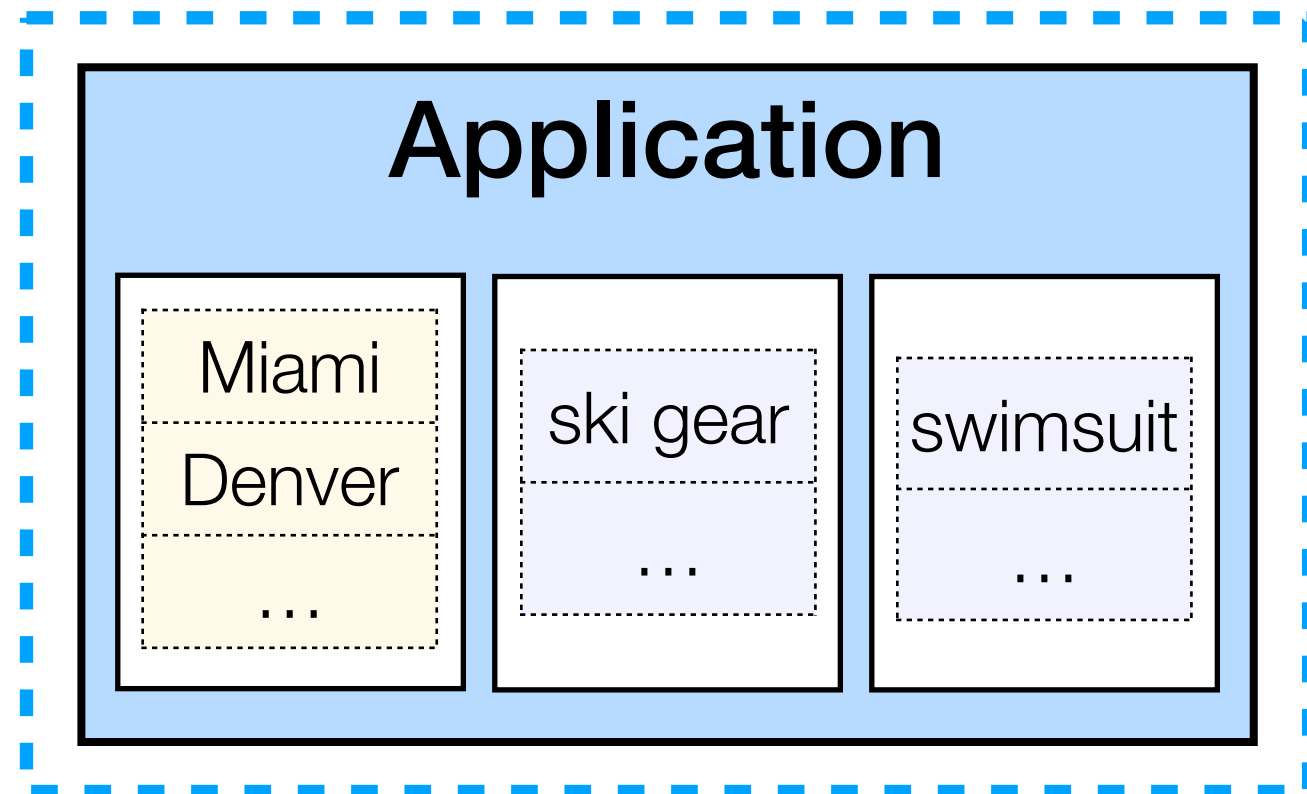
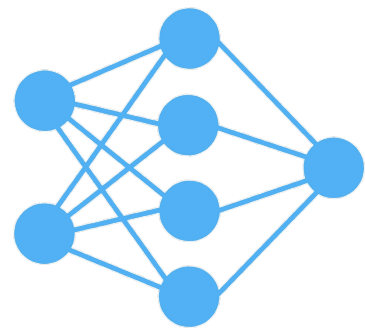
**Disambiguation
via arbitrary-
length context**

0x1000, **0x1001** \implies **Denver**, ski gear
0x1000, **0x105F** \implies **Miami**, swimsuit

Found in Translation (FiT) Attack

Initialization

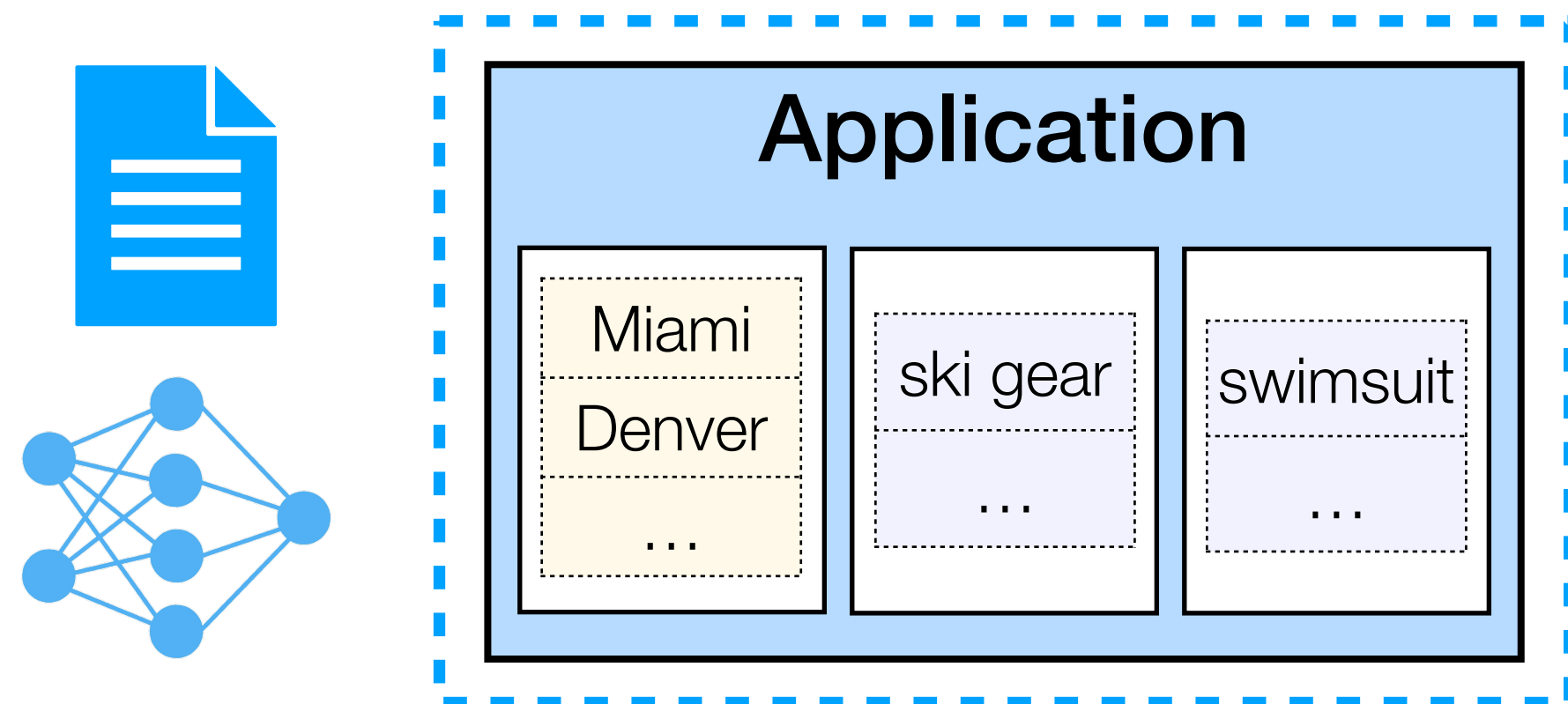
Set up monitored page table entries



Found in Translation (FiT) Attack

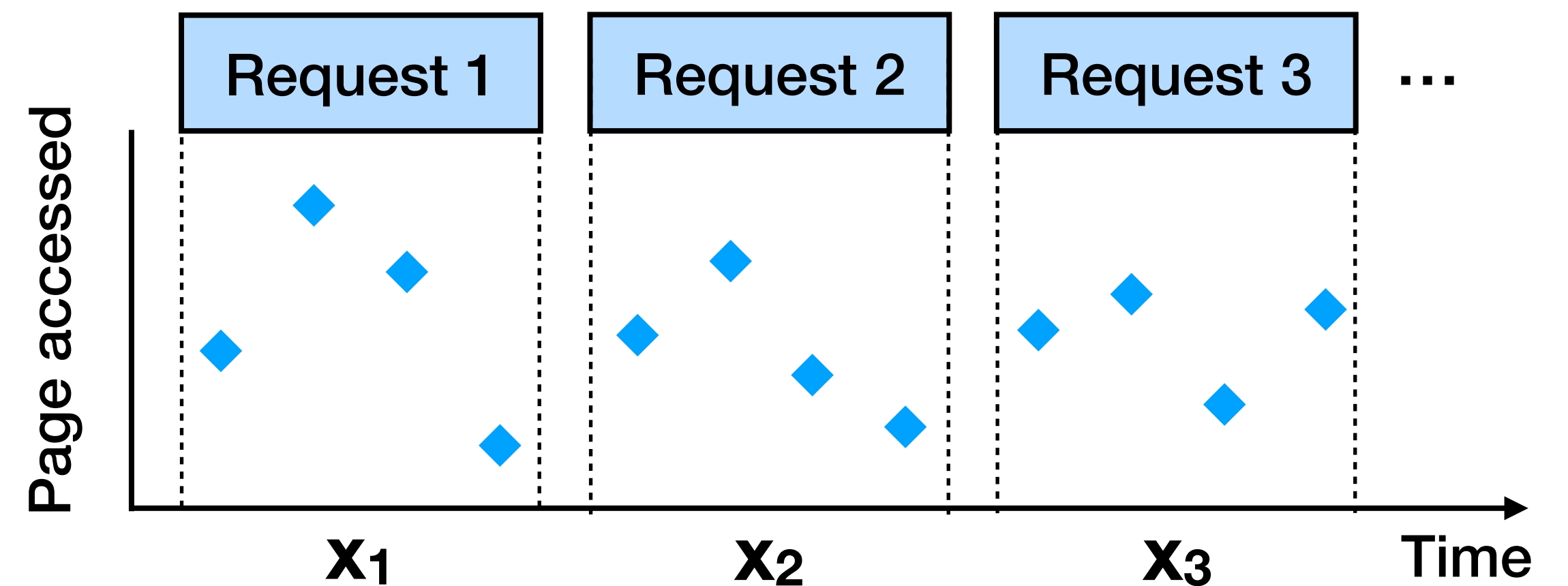
Initialization

Set up monitored page table entries



Online collection

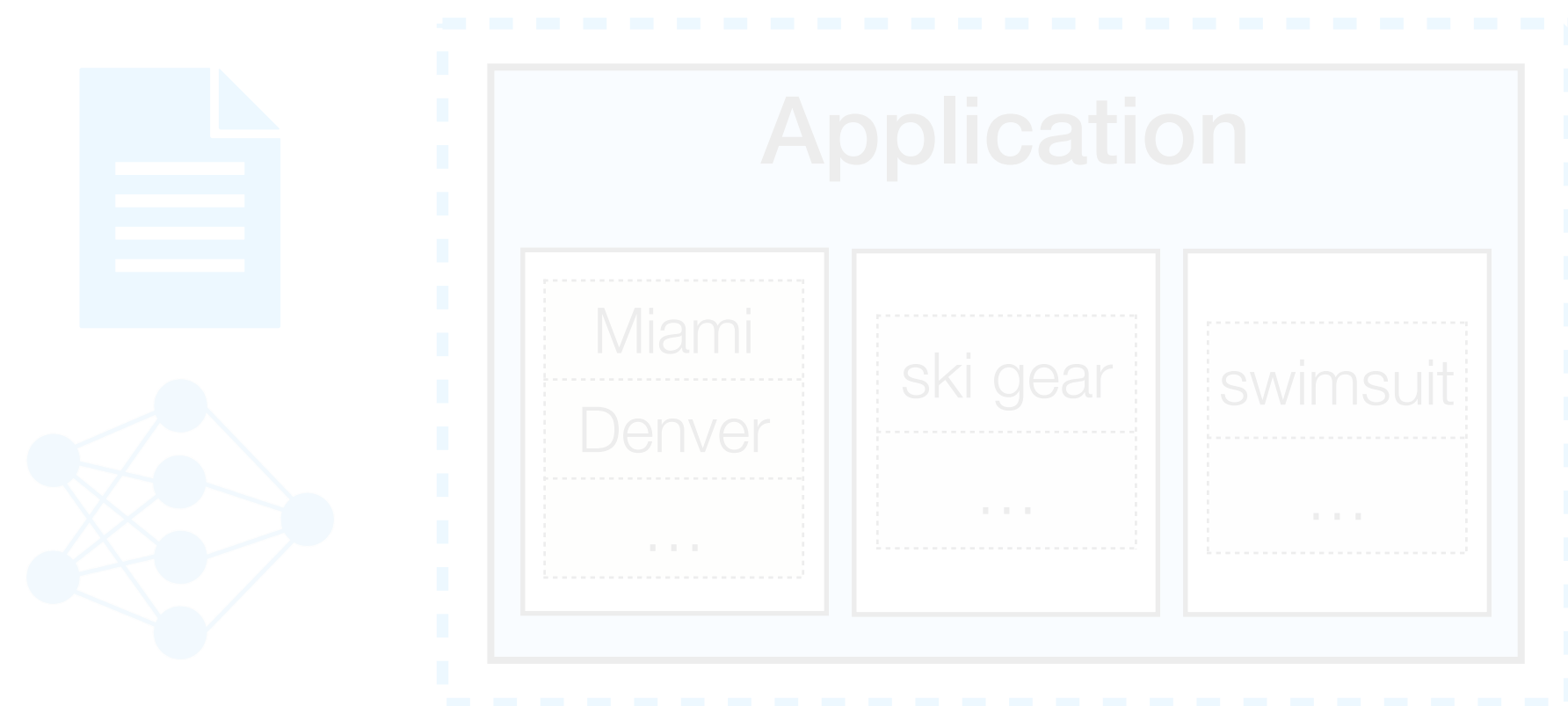
Track page access sequences \mathbf{x} via accessed bits



Found in Translation (FiT) Attack

Initialization

Set up monitored page table entries



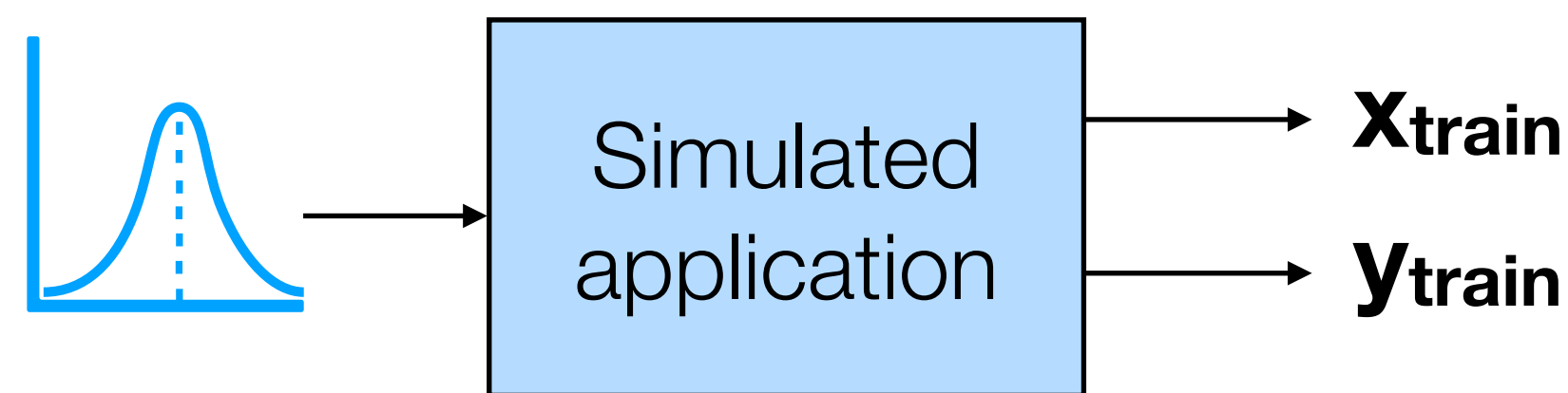
Online collection

Track page access sequences \mathbf{x} via accessed bits

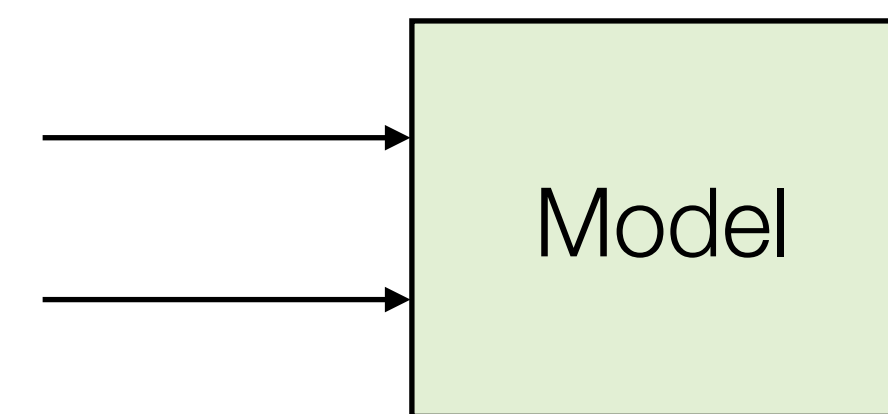


Offline analysis

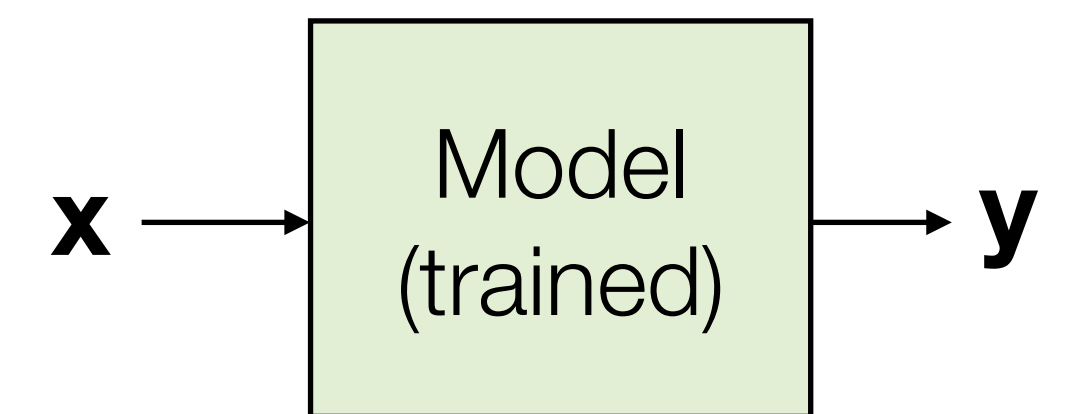
1. Generate training data



2. Training

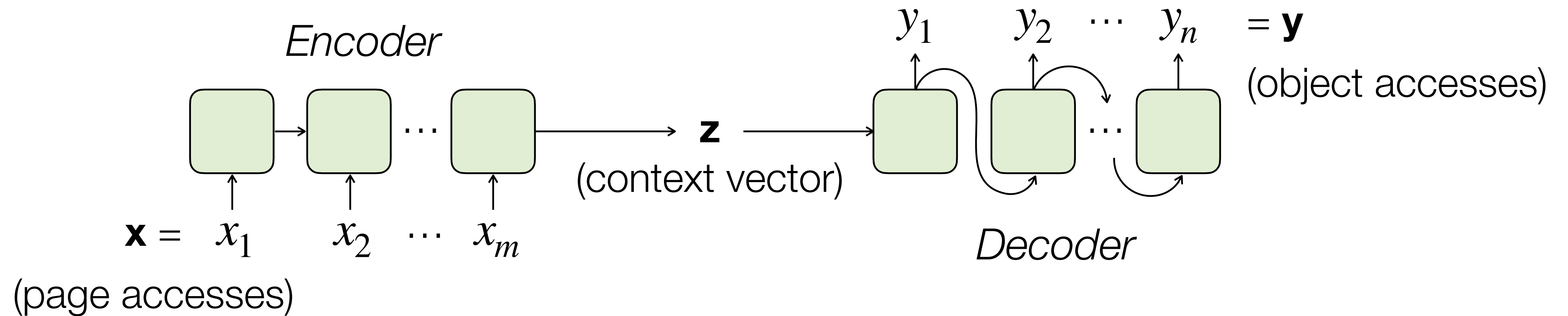


3. Inference



FIT Model Details

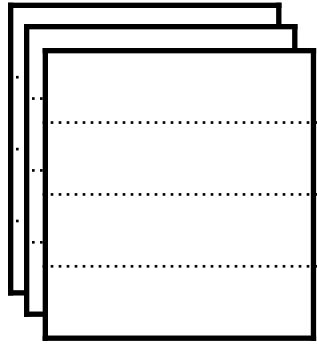
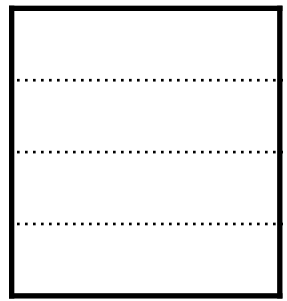
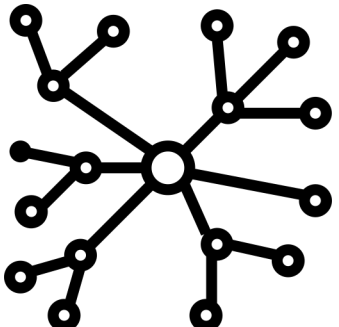
Deep generative language modeling with *recurrent encoder-decoder*



Use BERT as encoder and decoder

- ▶ Transformer processes sequences in parallel, learning global dependencies
- ▶ Leverage pre-trained weights for language prediction

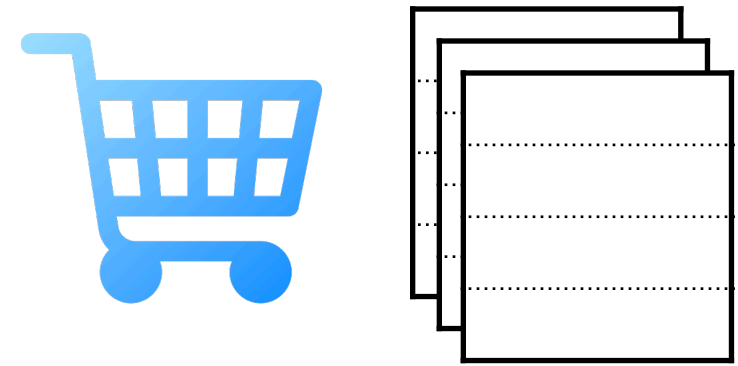
Evaluated Applications

	Privacy-sensitive use case	Access pattern	Objects	Dataset, train/test split
DLRM	User features for ad predictions	Per-feature lookup 	Embedding table entries	Criteo ad data, 1M / 100K
Medical LLM	Patient symptoms for diagnosis	Per-token lookup 	Embedding table entries	DDXPlus medical data, 500K / 50K
HNSW for semantic search	Biometrics, face images, etc.	Graph index traversal 	Nodes in graph	SIFT image search, 22.5K / 2.6K

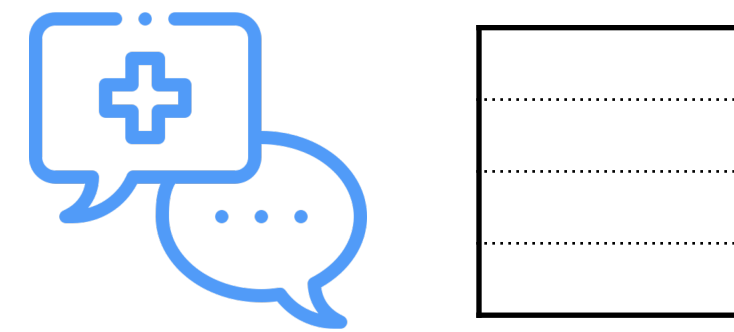
Evaluation Setup

Targeted applications

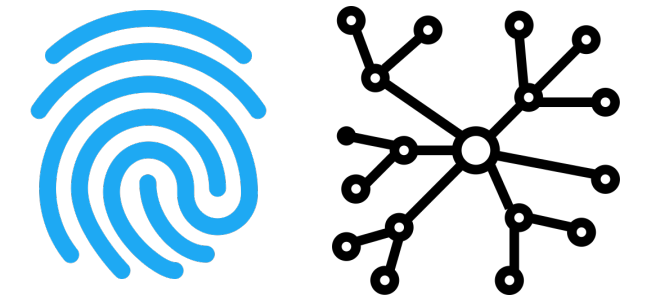
DLRM



LLM



HNSW



Executed in Nitro Enclave & Intel SGX Enclave (Ice Lake)

Compared attacks

FiT

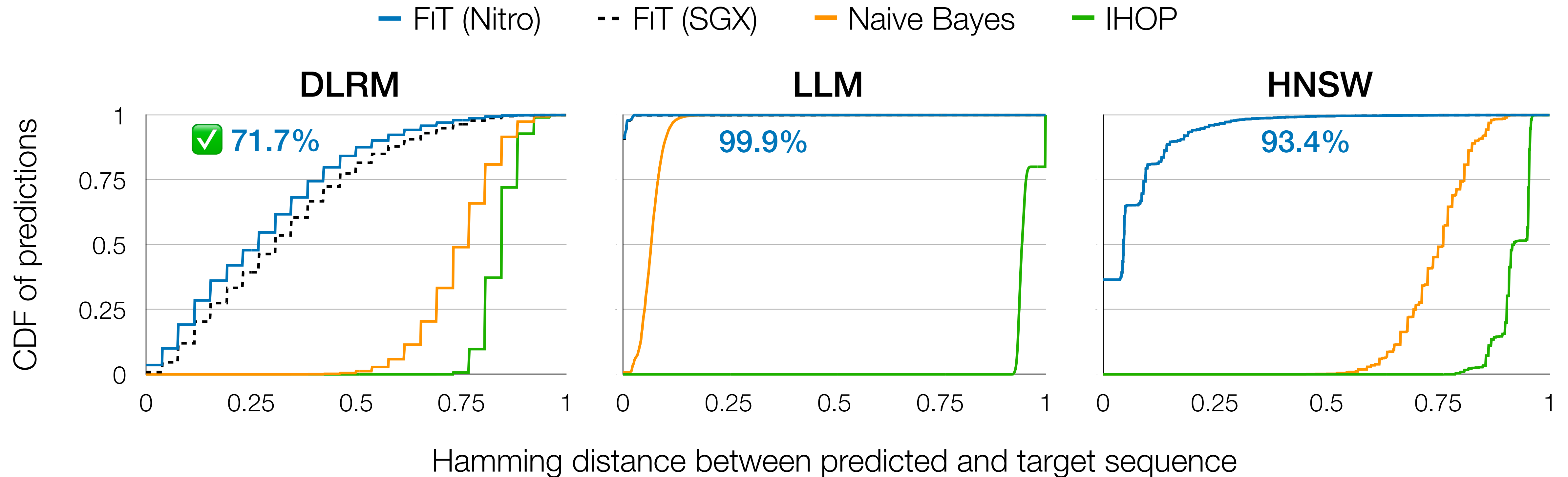
NVIDIA A40 GPUs

Naive Bayes

AMD EPYC 7302P CPUs

IHOP

Attack Efficacy

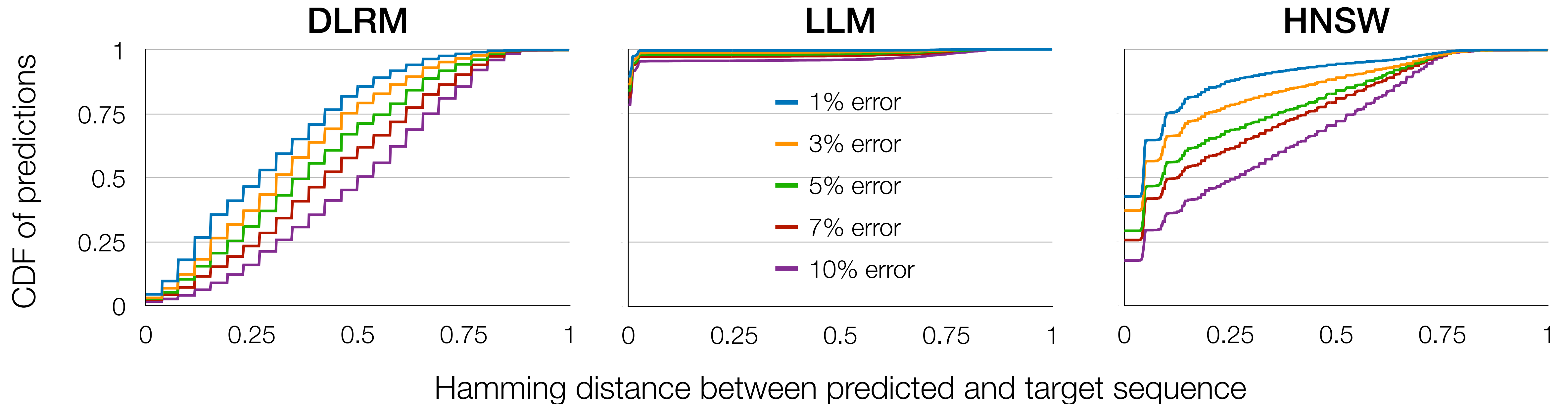


FiT achieves high accuracy by learning (i) many-to-1 object-to-page mappings
(ii) arbitrary-length correlations

FiT Sensitivity to Error

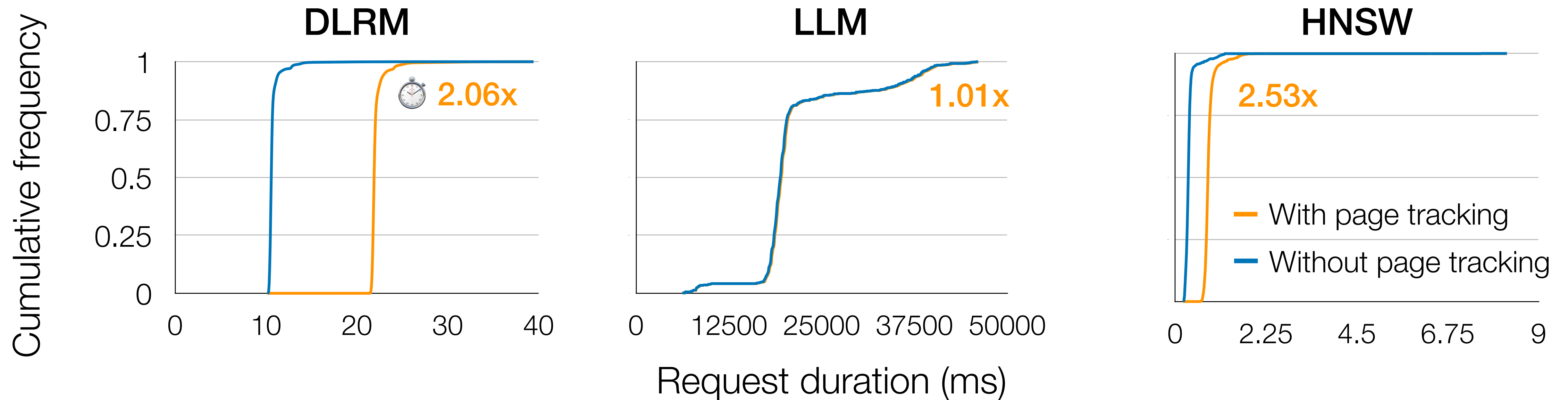
**~1% error rate in
observed page accesses**

- ▶ False positives: prefetching; accesses to co-located metadata
- ▶ False negatives: background processes clearing accessed bit



Context awareness makes FiT robust against errors in collected page sequences

FiT Latency Overheads



***Online* page tracking adds low latency – masked by network, tiered memory, etc.**

***Offline* analysis time depends on application (1 – 45 GPU hours)**

Conclusion

Found in Translation

- Leverages *language modeling* to capture both
 - (i) many-to-1 object-page mappings
 - (ii) arbitrary-length correlations
- Practical and achieves high accuracy

Future Work

- Adaptation of attack to other settings
- Sensitivity study for page sizes, memory tiering, etc.
- Correlation-hiding countermeasures

Paper

[usenix.org/conference/usenixsecurity25/
presentation/jia-grace](https://usenix.org/conference/usenixsecurity25/presentation/jia-grace)

Code

github.com/yale-nova/found-in-translation