

Encrypted Search: Intro & Basics

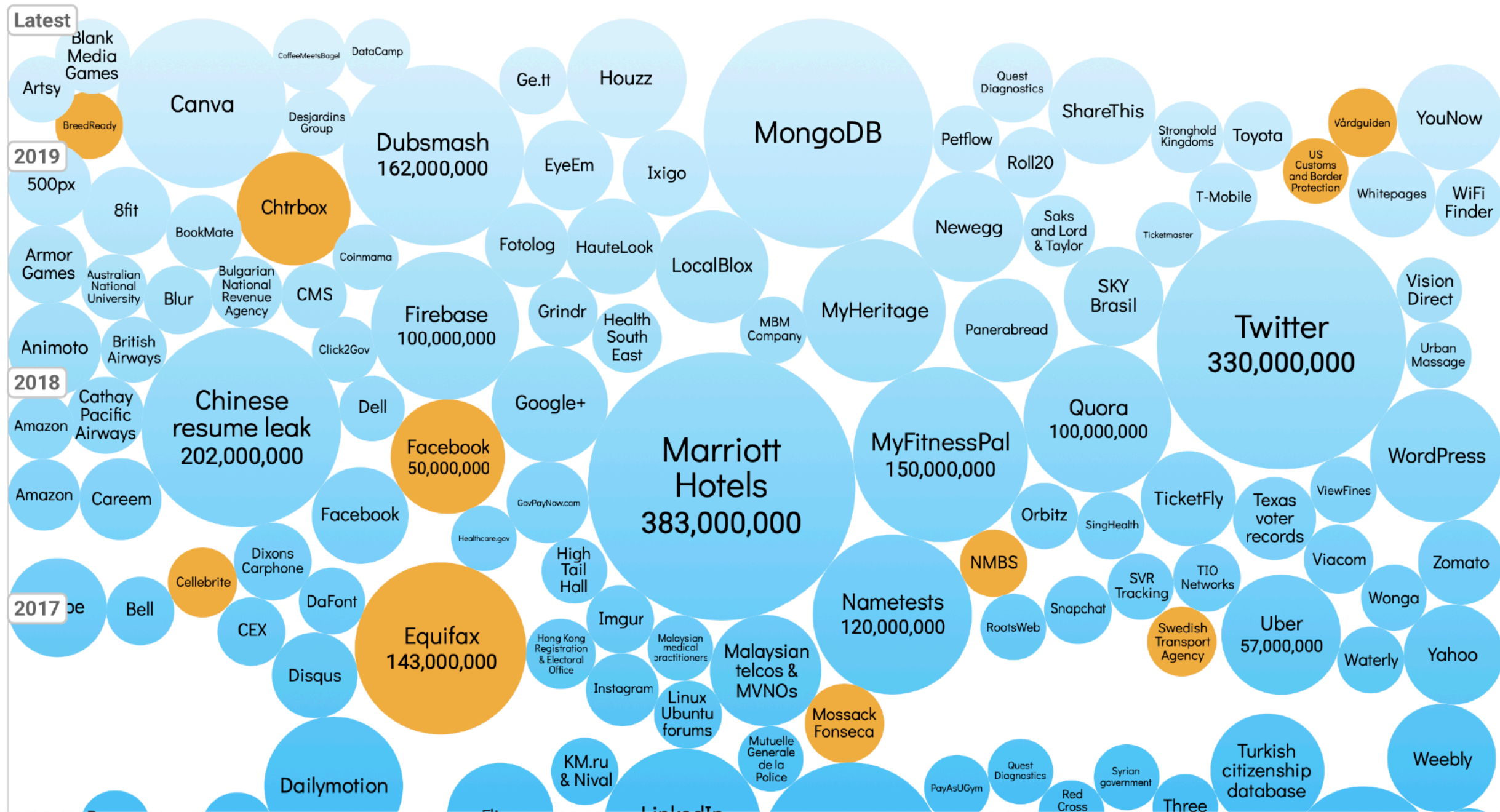
Seny Kamara



BROWN



ENCRYPTED
SYSTEMS LAB



14,717,618,286*

4%

Why so Few?



Incompetence?



Lazyness?

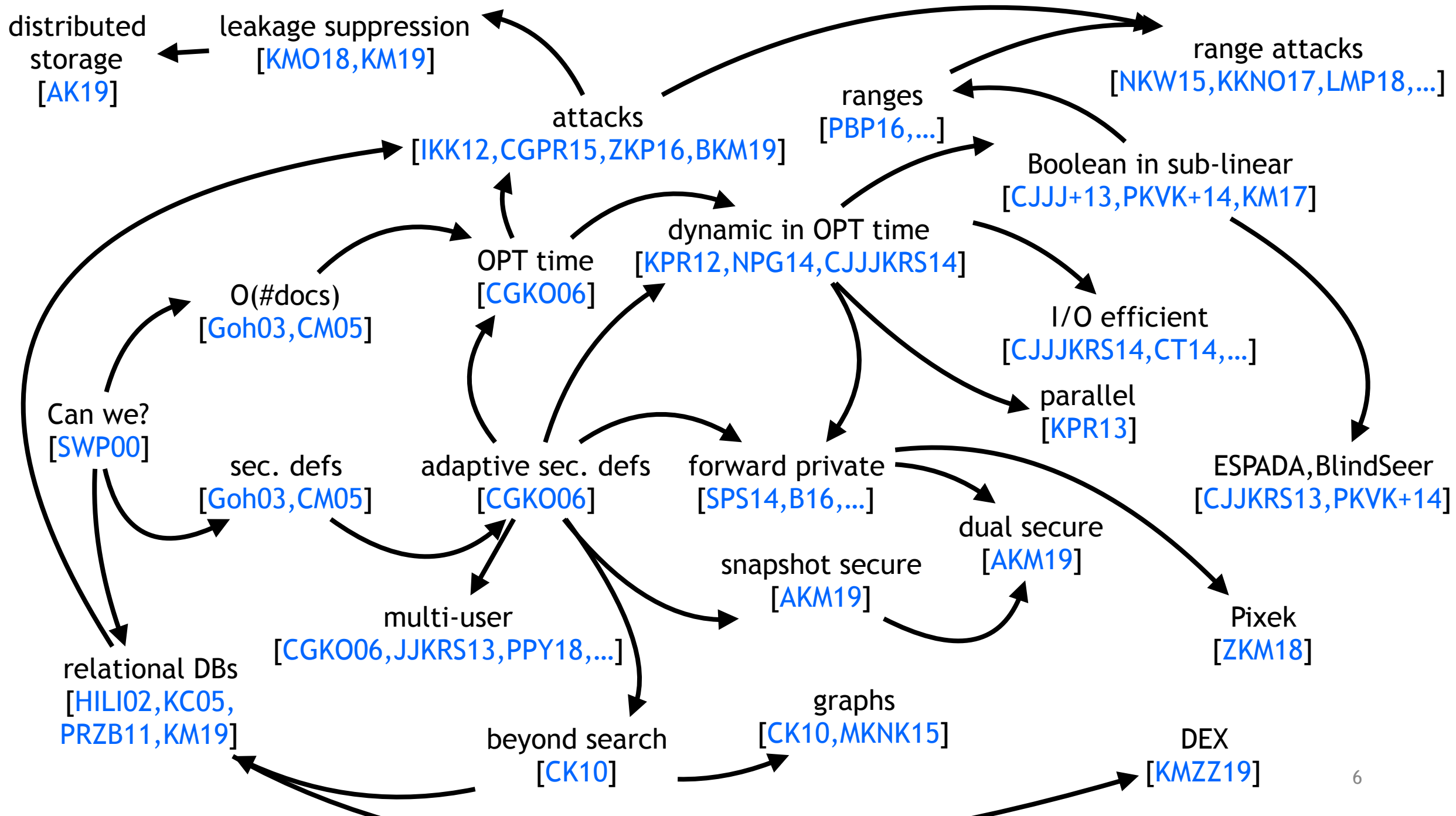


Cost?

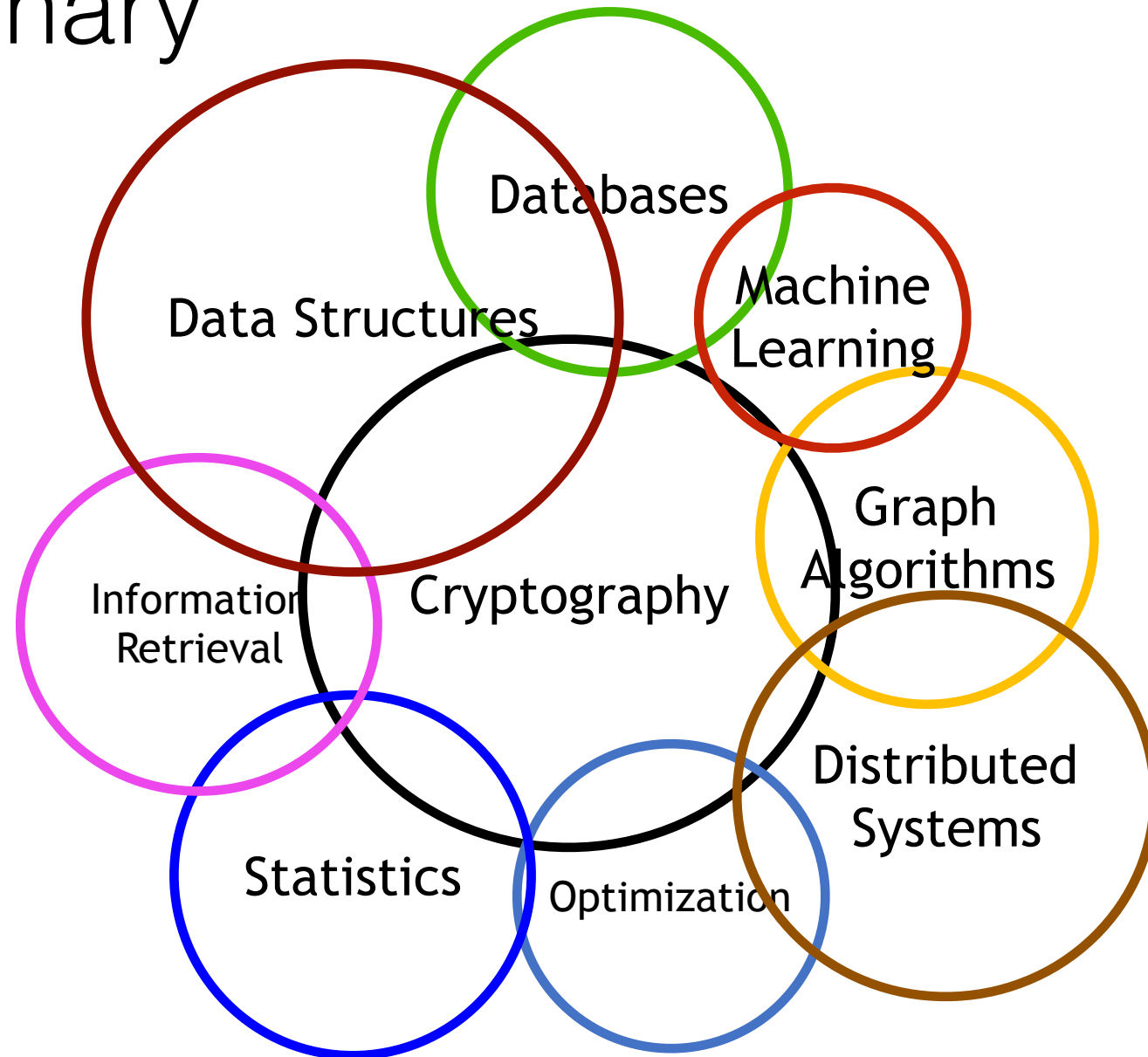
“...because it would have hurt Yahoo’s ability to index and search message data...”

— J. Bonforte in NY Times

Q: can we search on encrypted data?



Interdisciplinary

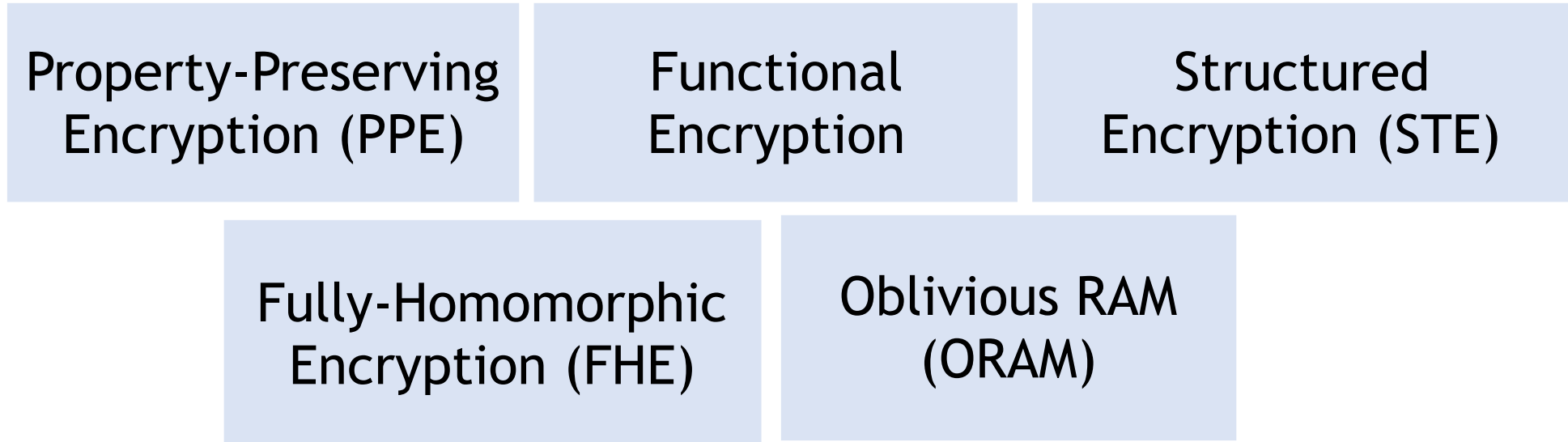


Real-World Problem

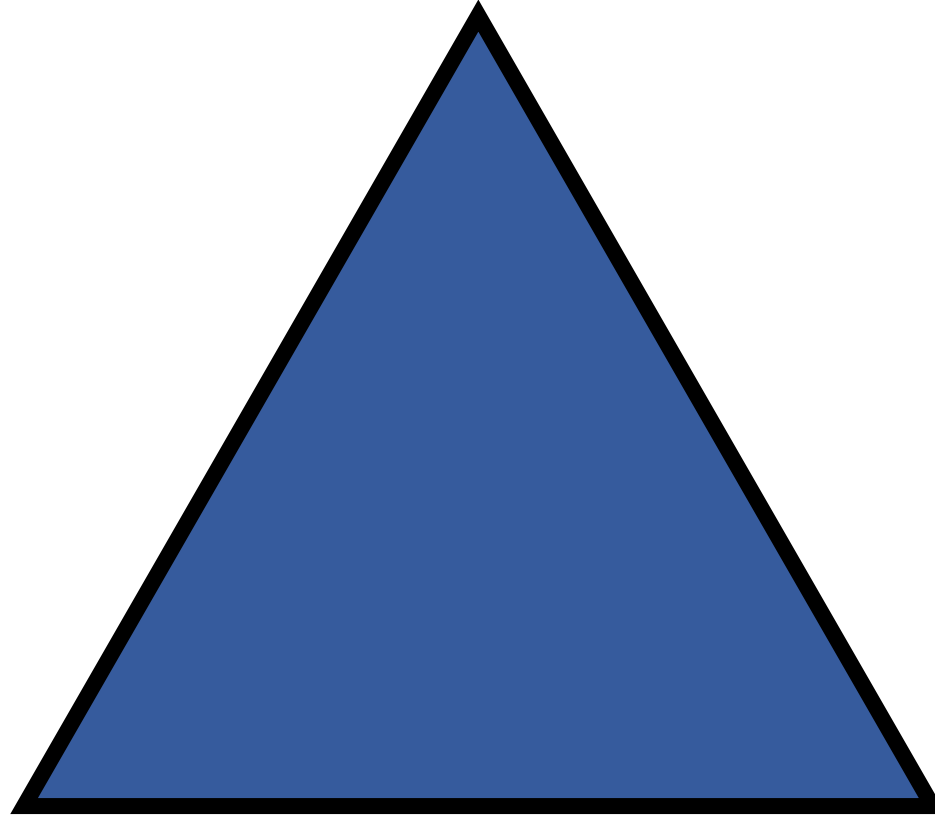


- Major companies
 - Microsoft, SAP
 - MongoDB, Cisco
 - Google Research
 - Hitachi, Fujitsu
 - more...
- Funding agencies
 - NSF
 - IARPA
 - DARPA
- Startups
 - Ciphercloud
 - Skyhigh Networks
 - Bitglass
 - Baffle
 - Cossack Labs
 - Strong Salt, Overnest
 - many many more

Encrypted Search (Building Blocks)



Efficiency



Functionality

Leakage



What is Search?

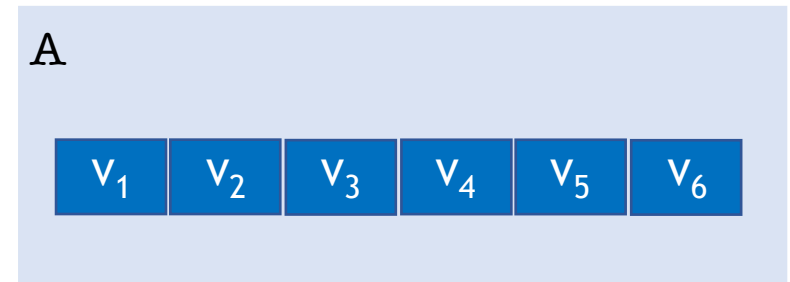
- Complexity regimes
 - linear search: $O(n)$
 - sub-linear search: $o(n)$
- Algorithmic paradigms
 - with pre-processing
 - without pre-processing
- For medium to large data
 - sub-linear search is a *requirement*; not an option

	Without Pre-Processing	With Pre-Processing
Linear	sequential scan	not interesting
Sub-Linear	read sub-set of input (errors)	data structures

Background: Data Structures

- Abstract data types
 - capture functionality
 - ex: dictionary
- Data structures
 - instantiate ADTs
 - ex: hash table, binary search tree
- As common in CS
 - we sometimes blur the distinction

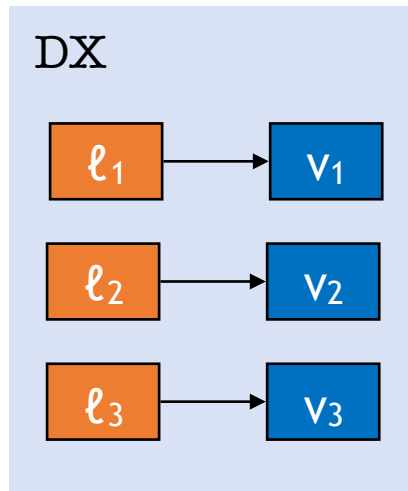
- Arrays store values



- Write: $A[i] := v_i$
- Read: $A[i]$ returns v_i

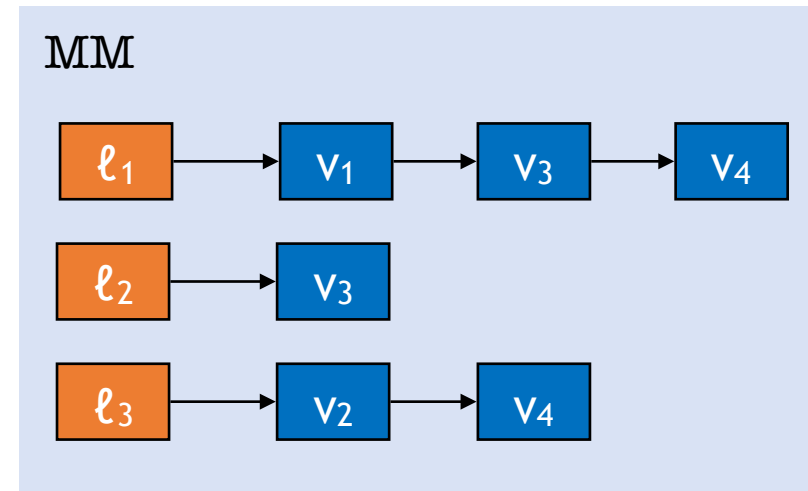
Background: Data Structures

- Dictionaries map labels to values



- Put: $DX[\ell_2] := v_2$
- Get: $DX[\ell_2]$ returns v_2

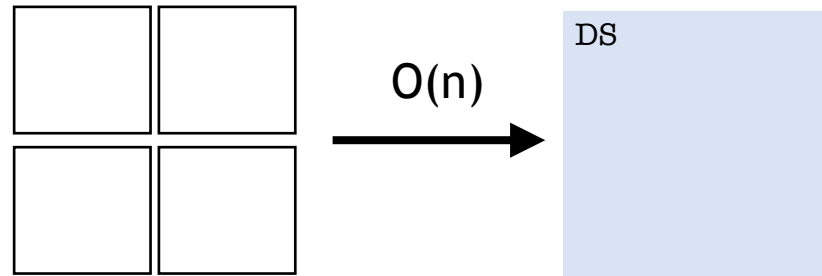
- Multi-Maps map labels to tuples



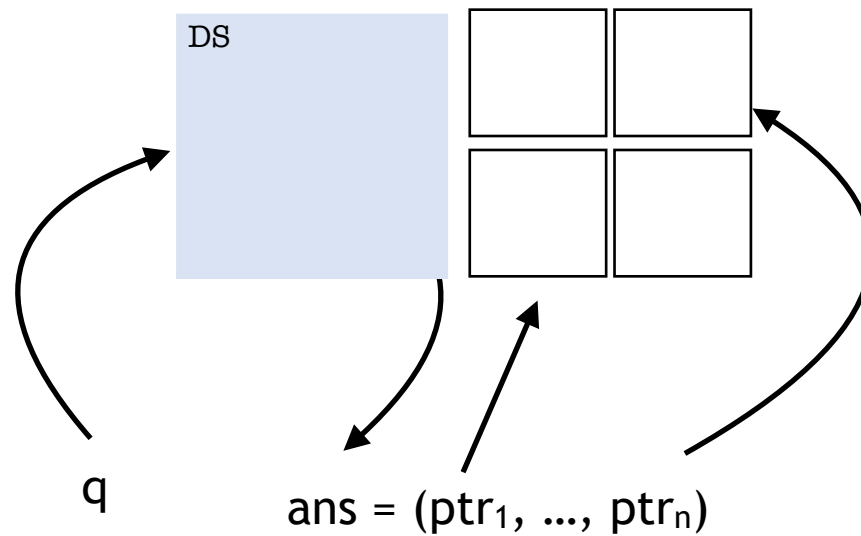
- Put: $MM[\ell_3] := (v_2, v_4)$
- Get: $MM[\ell_3]$ returns (v_2, v_4)

Keyword Search in Sub-Linear Time

Setup time

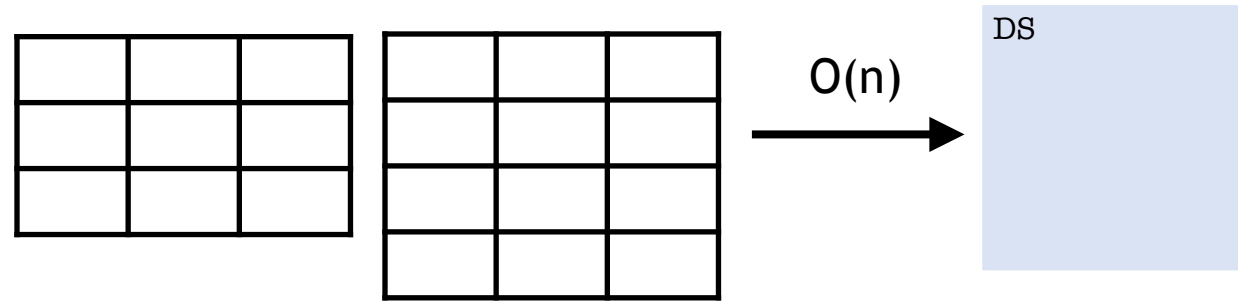


Query time

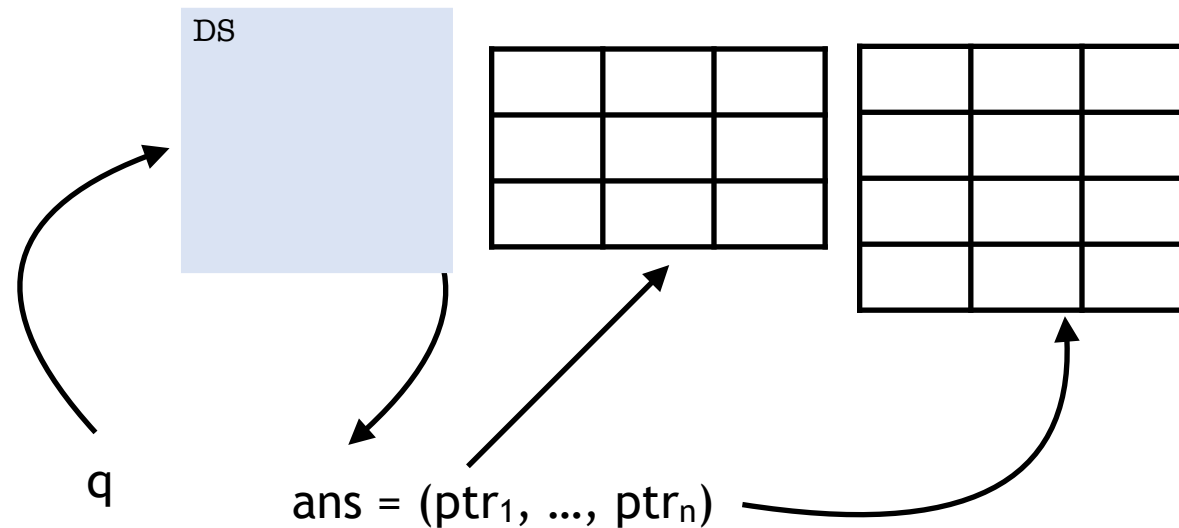


Database Queries in Sub-Linear Time

Setup time



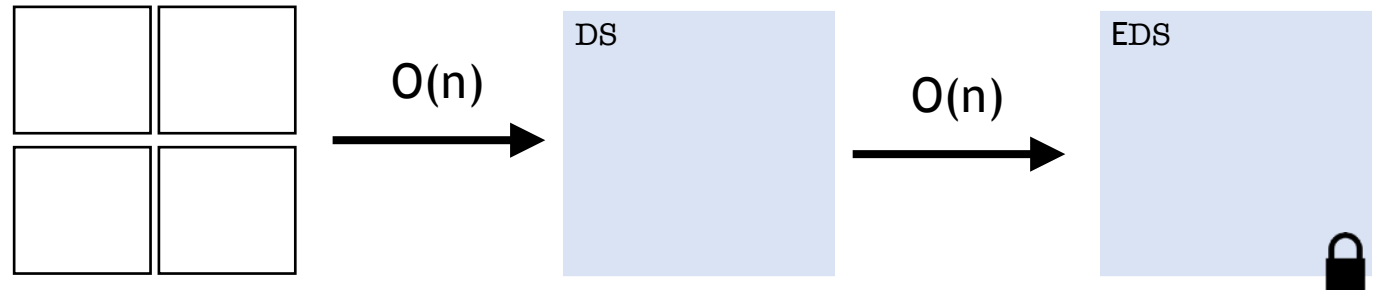
Query time



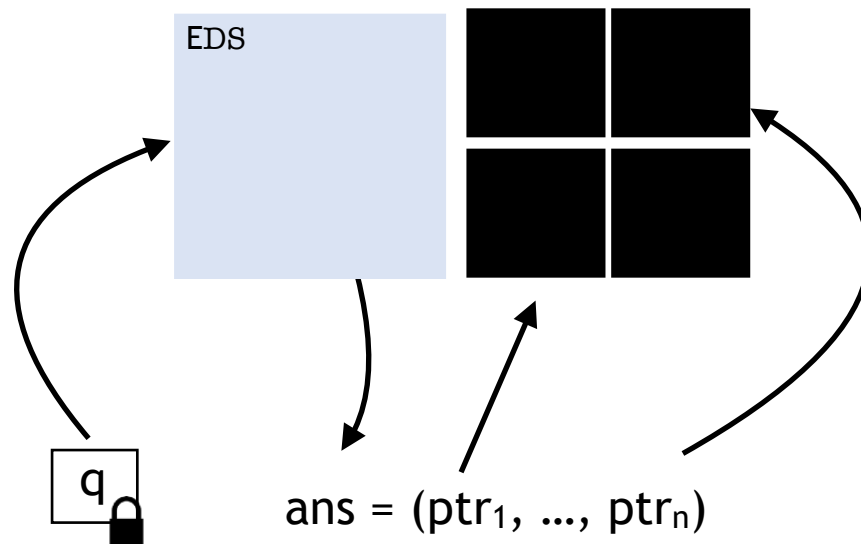
Q: how do we do sub-linear search on encrypted data?

Encrypted Keyword Search in *Sub-Linear* Time

Setup time

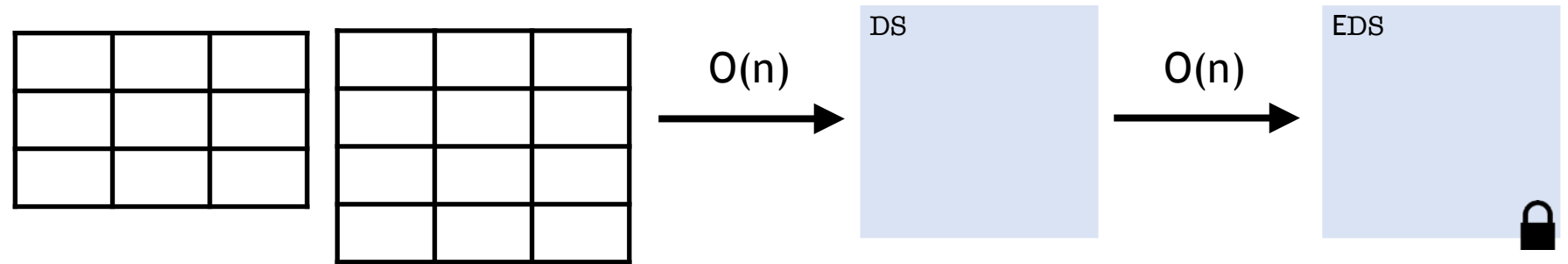


Query time

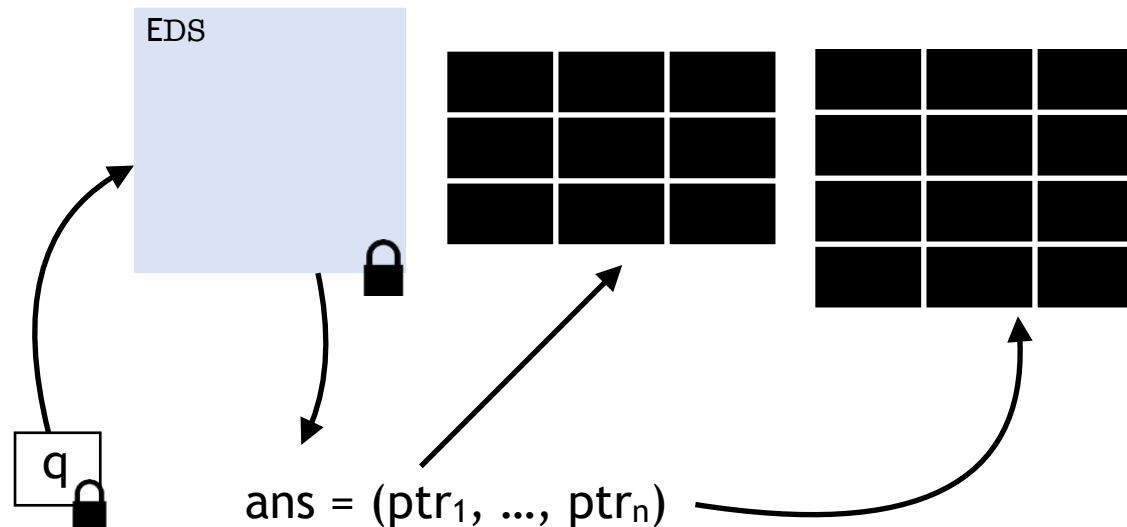


Encrypted Database Queries in Sub-Linear Time

Setup time



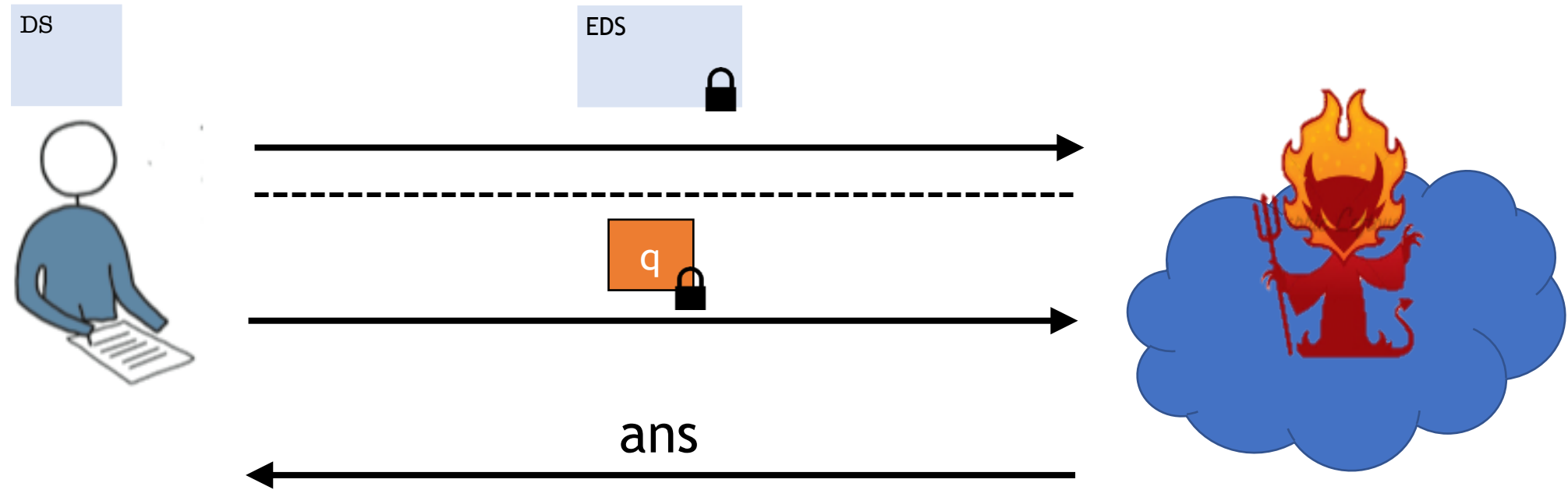
Query time



Q: how do we formalize *encrypted* data structures?

Structured Encryption

[Chase-K.10]

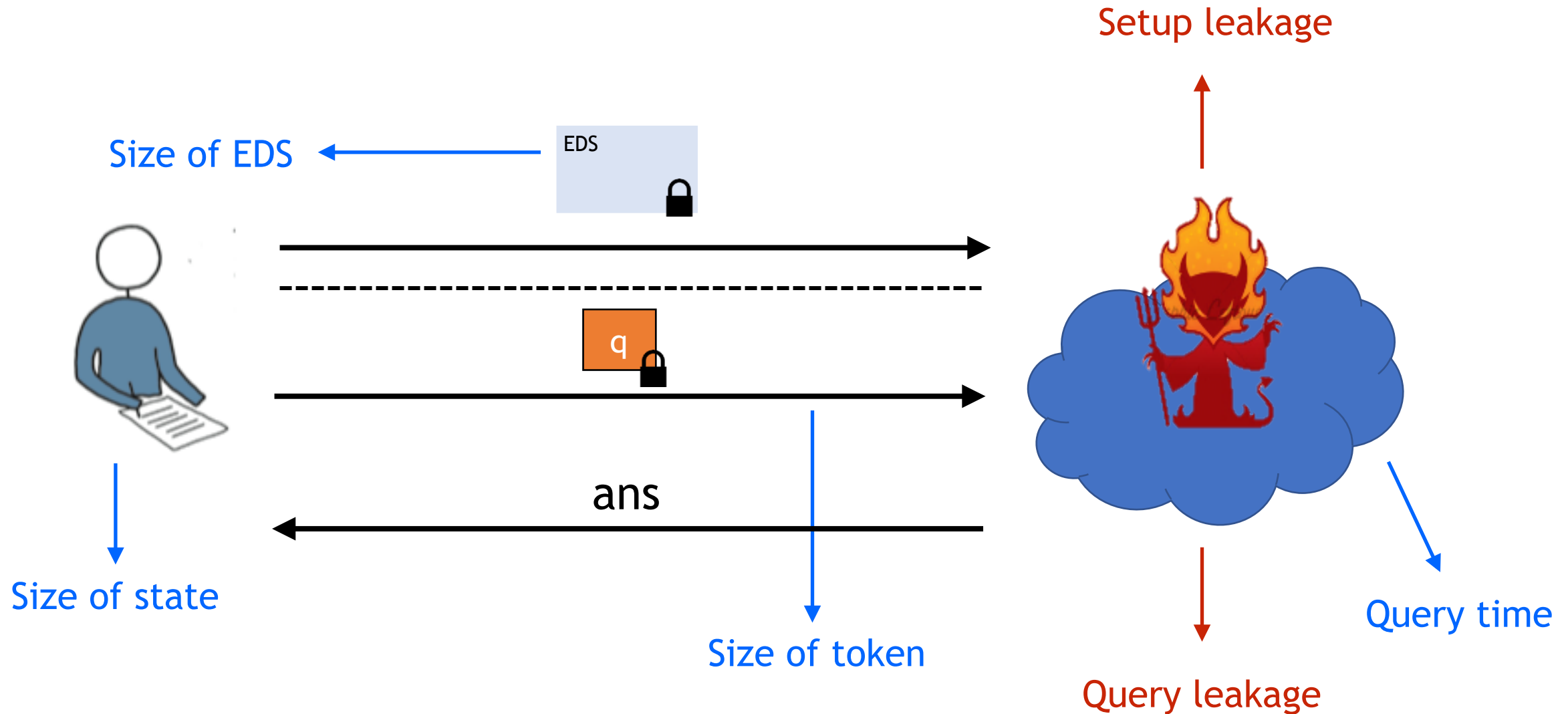


$\text{Setup}(1^k, \text{DS}) \longrightarrow (K, \text{EDS})$

$\text{Token}(K, q) \longrightarrow \text{tk}$

$\text{Query}(\text{EDS}, \text{tk}) \longrightarrow \text{ans}$

Desiderata



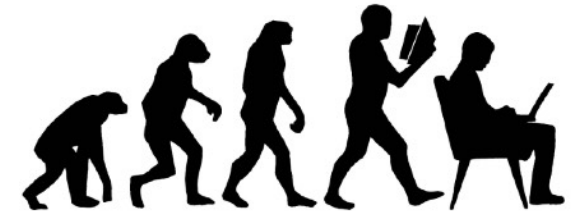
Structured Encryption

[[Chase-K.10](#)]

- Many variants of STE
 - response-revealing
 - EDS query reveals answer in plaintext
 - response-hiding
 - EDS query reveals encrypted answer
 - non-interactive queries
 - clients sends single message called a token
 - interactive queries
 - client and server execute multi-round protocol



Evolution of Structured Encryption



Efficiency

- '00 Linear in file length [SWP00]
- '03 Linear in #docs [Goh03]
- '06 Optimal [CGK006,CK10]
- '12 Optimal Dynamic [KPR12,CJJJKRS14]
- '14 I/O efficient [CT14,CJJJKRS14,ANSS16,DPP18],ASS18]

Expressiveness

- '00 Single-keyword SSE [SWP00,Goh03,CGK006,CJJJKRS14]
- '06 Multi-user SSE [CGK006,JKRS13,PPY16,HSWW18]
- '13 Boolean SSE [CJJJKRS13,PKVK+14,KM17]
- '14 Range SSE [PKVK+14,FJKNRS15]
- '18 STE-based SQL [KM18]

Security

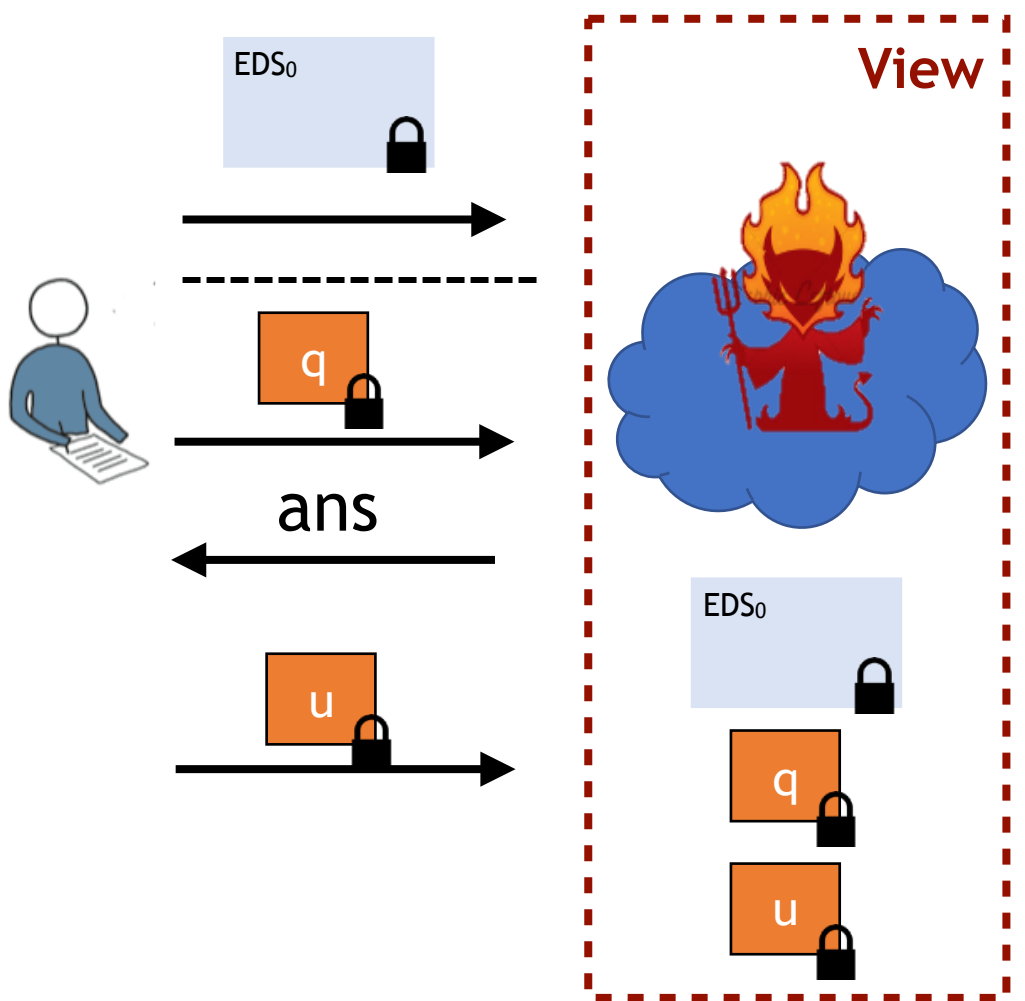
- '06 Leakage-parametrized security definitions [CGK006]
- '12 Attacks [IKK12,CGPR15,ZKP16,KMNO16,LMP18,GLMP18]
- '14 Forward/Backward Security [SPS14,Bost16,LC17,BMO17,AKM18]
- '18 Leakage Supression [KMO18,KM19]
- '19 Snapshot [AKM18]



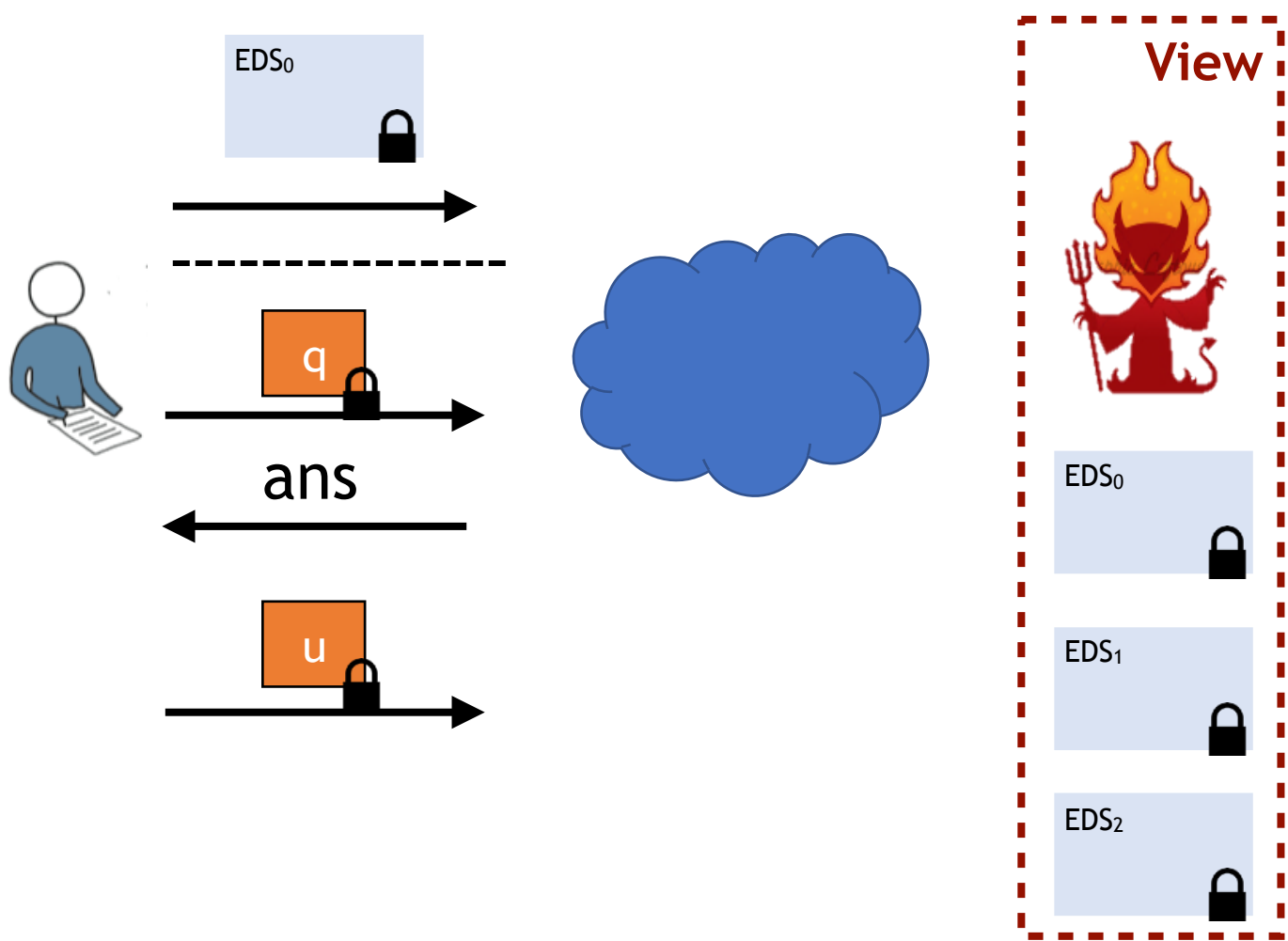
Adversarial Models

Adversarial Models

Persistent



Snapshot



Persistent (Adaptive) Security

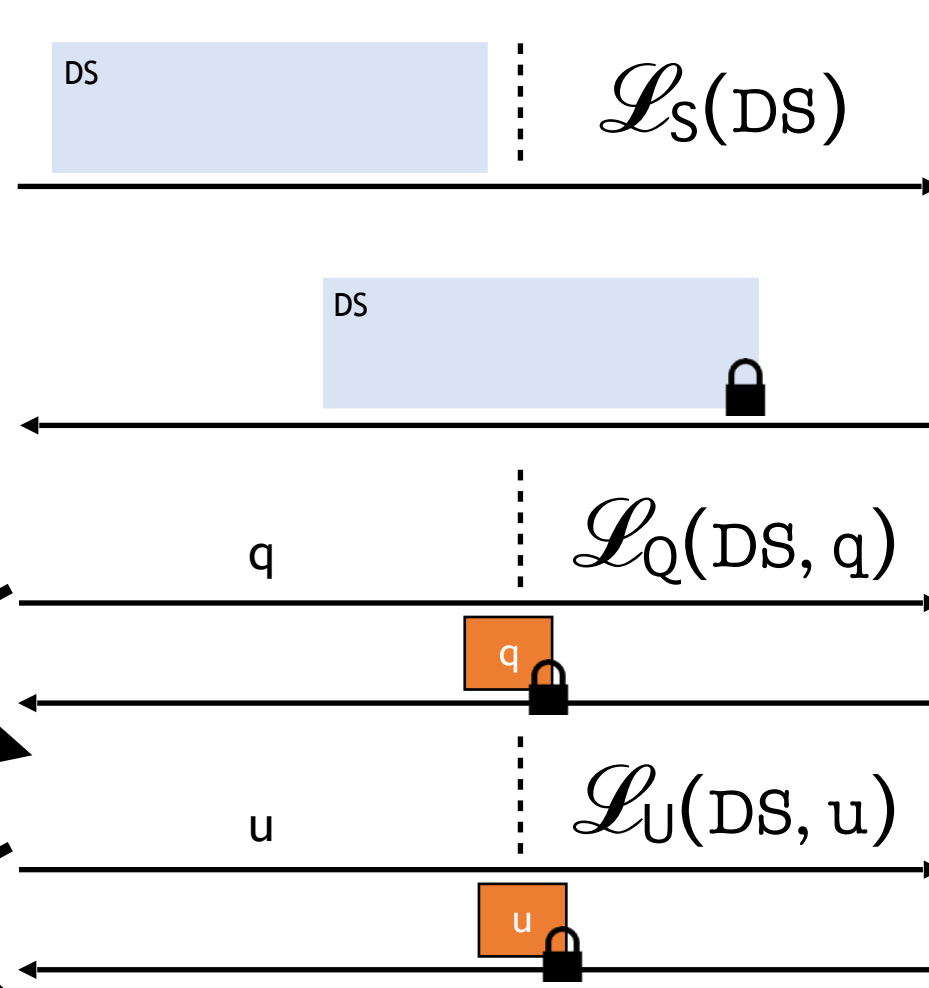
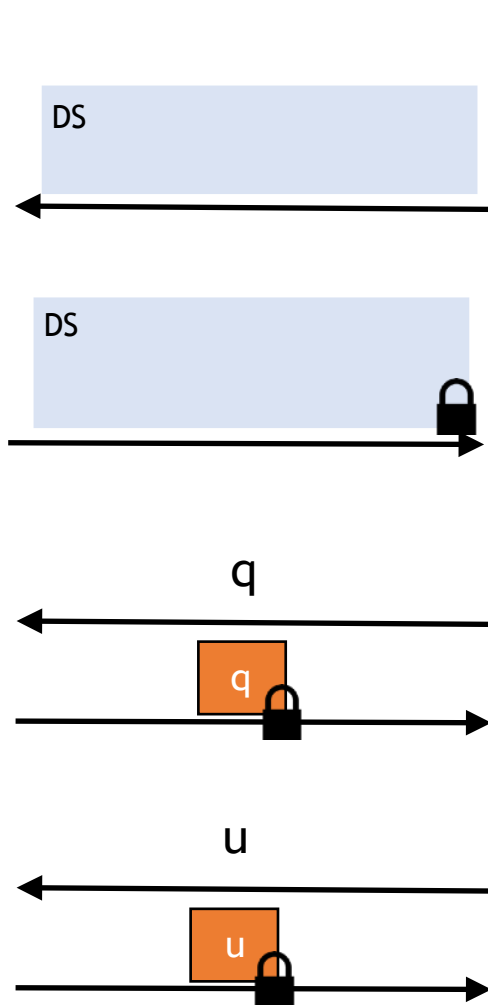
[Curtmola-Garay-K.-Ostrovsky06, Chase-K.10]

- An STE scheme is $(\mathcal{L}_S, \mathcal{L}_Q)$ -secure vs. a persistent adv. if
 - it reveals no information about the *structure* beyond \mathcal{L}_S
 - it reveals no information about the *structure* and *query* beyond \mathcal{L}_Q

Persistent (Adaptive) Security

[Curtmola-Garay-K.-Ostrovsky06, Chase-K.10]

Real



Ideal



Forward Privacy

[[Stefanov-Papamanthou-Shi14](#), [Bost16](#)]

- Informally [[SPS14](#)]
 - “Updates not correlated to previous queries”
- Formally [[Bost16](#)]
 - $\mathcal{L}_U(\text{MM}, (\ell, \mathbf{v})) = \#\mathbf{v}$

Snapshot (Adaptive) Security

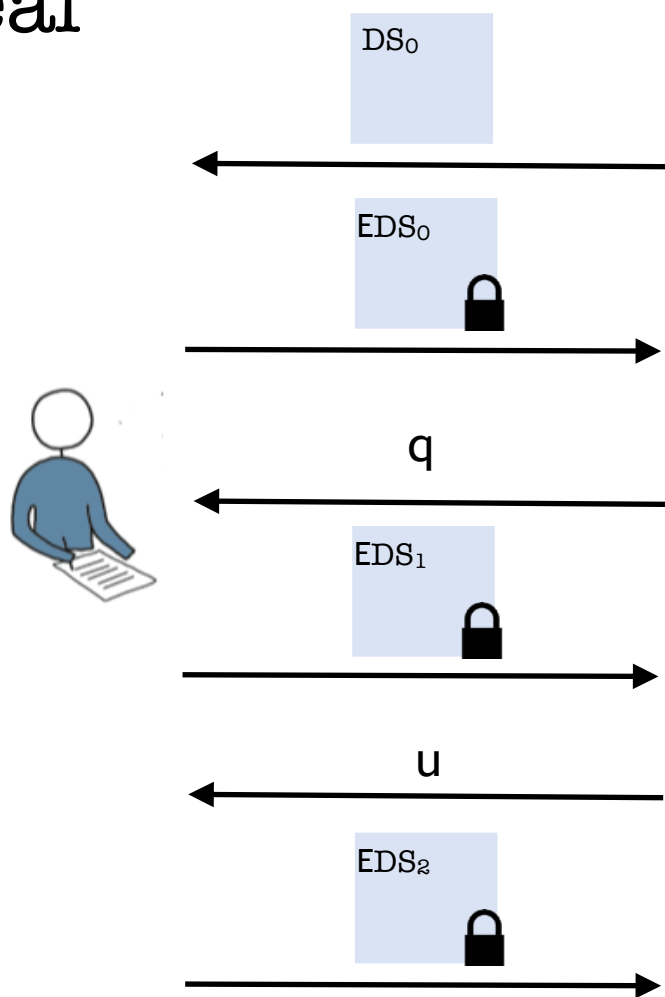
[Amjad-K.-Moataz19]

- We say that an STE scheme is \mathcal{L}_{Snp} -secure vs. a snapshot adv. if
 - it reveals no information about the *structure* beyond \mathcal{L}_{Snp}

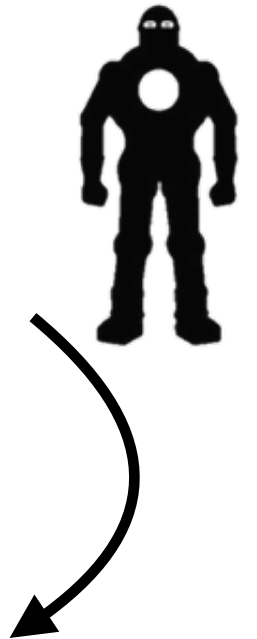
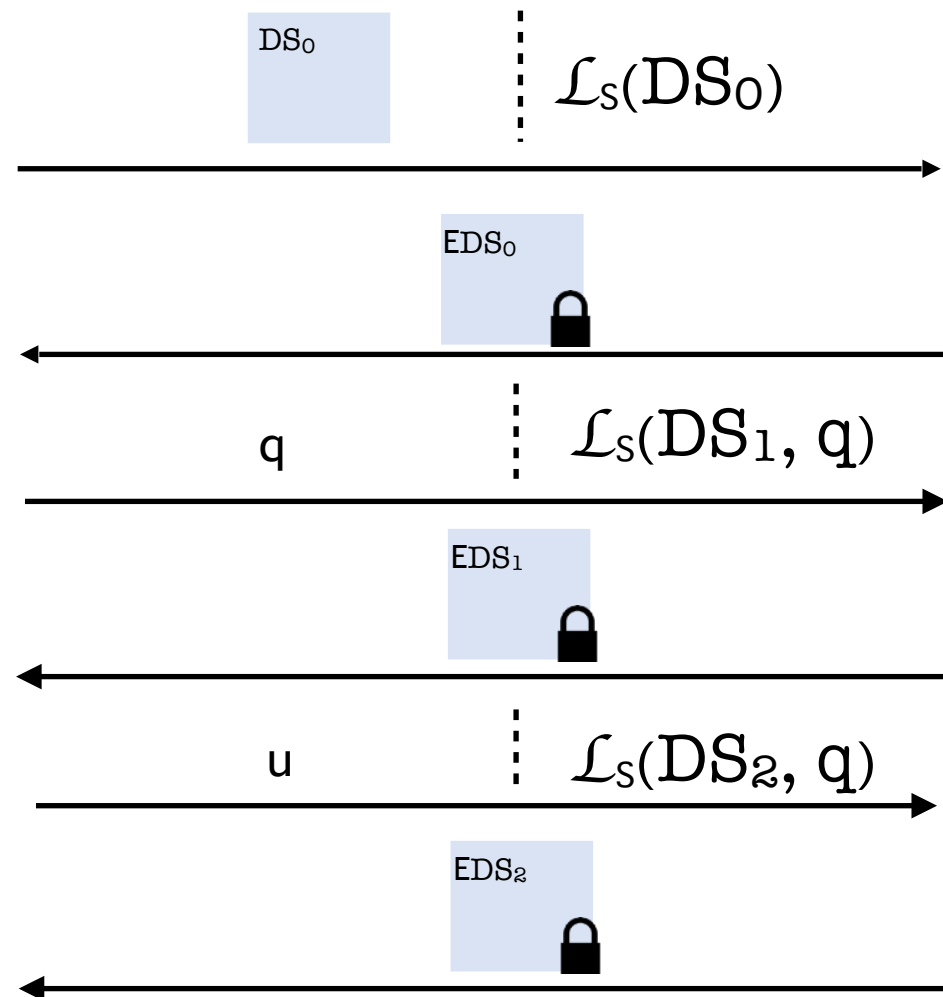
Snapshot (Adaptive) Security

[[Amjad-K.-Moataz19](#)]

Real



Ideal



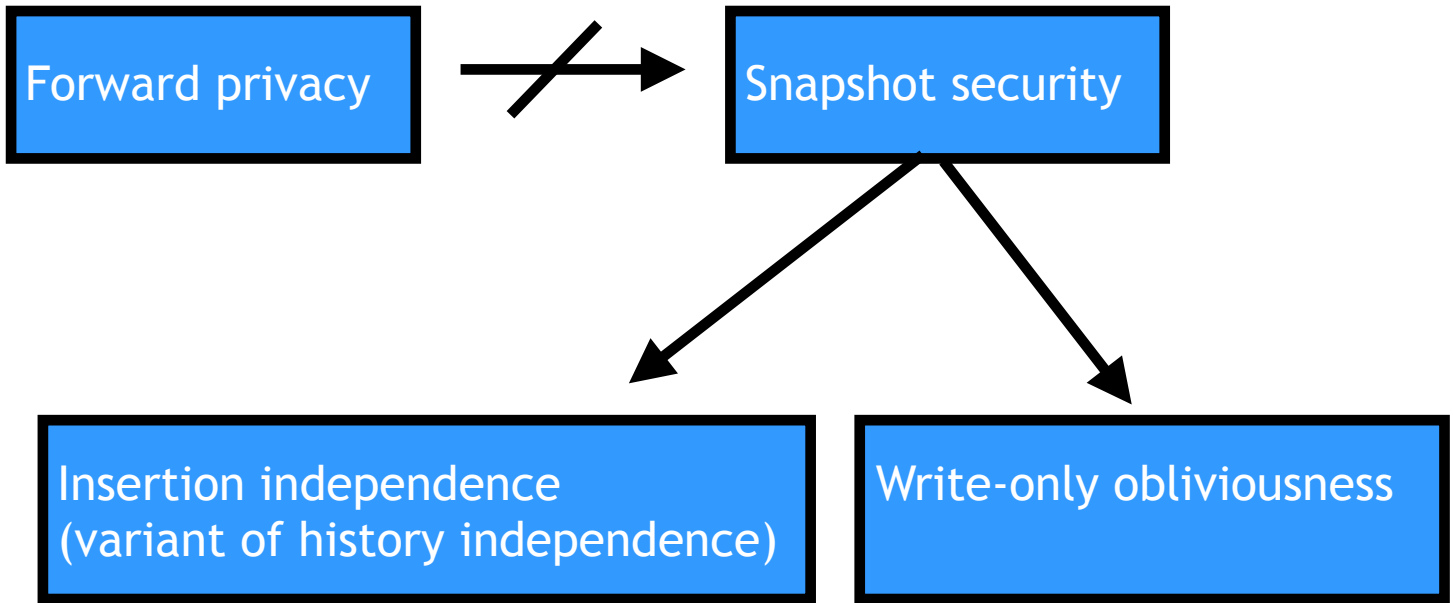
Snapshot (Adaptive) Security

[Amjad-K.-Moataz19]

Static Structures

$$\mathcal{L}_{\text{Snp}} = \mathcal{L}_{\text{S}}$$

Dynamic Structures



Q: Why do we parameterize definitions with leakage?

Leakage-Parameterized Definitions

[[Curtmola-Garay-K.-Ostrovsky, Chase-K.10](#)]

- This area is about tradeoffs
 - but traditional cryptographic definitions don't capture tradeoffs
- in 00's, different approaches were proposed to capture leakage
 - #1: limit adversary's power in the proof
 - #2: make assumptions on data (e.g., high entropy)
- Original motivations for leakage-parameterized definitions
 - Approaches #1 & #2 are misleading (sweep leakage under the rug)
 - Leakage should be made explicit and not be implicit
 - gives clear target for cryptanalysis
 - makes it (somewhat) easier to compare schemes

Q: How do we model leakage?

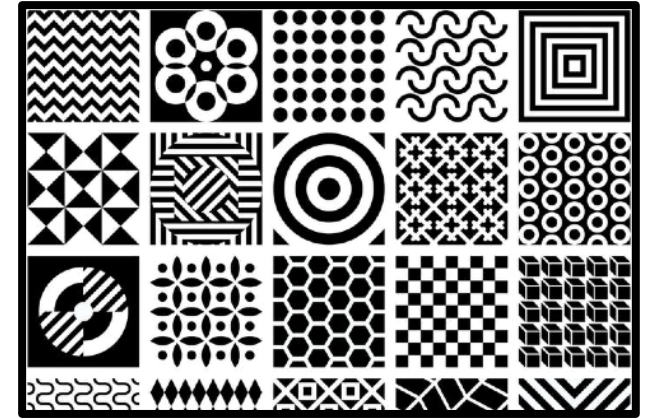
Modeling Leakage



- Each scheme has a leakage profile: $\Lambda = (\mathcal{L}_S, \mathcal{L}_Q, \mathcal{L}_U)$
 - where $\mathcal{L}_S = (\text{patt}_1, \dots, \text{patt}_n)$ is the Setup leakage
 - $\mathcal{L}_Q = (\text{patt}_1, \dots, \text{patt}_n)$ is the Query leakage
 - $\mathcal{L}_U = (\text{patt}_1, \dots, \text{patt}_n)$ is the Update leakage
- Each “operational” leakage is composed of leakage patterns
 - $(\text{patt}_1, \dots, \text{patt}_n)$

Common Leakage Patterns

[[K.-Moataz-Ohrimenko18](#)]



- **qeq**: query equality
 - a.k.a. search pattern
- **rid**: response identity
 - a.k.a. access pattern
- **qlen**: query length
- **trlen**: total resp. length
- **rlen/vol**: response length
 - a.k.a. volume pattern
- **req**: response equality
- **mqlen**: max query length
- **mrln**: max resp. length
- **srln**: sequence resp. length
- **dsize**: data size
- **usize**: update size
- **did**: data identity

Example Leakage Profiles

- The “Baseline” leakage profile for response-revealing EMMs
 - $\Lambda = (\mathcal{L}_S, \mathcal{L}_Q, \mathcal{L}_U) = (\text{dsize}, (\text{qeq}, \text{rid}), \text{usize})$
- The “Baseline” leakage profile for response-hiding EMMs
 - $\Lambda = (\mathcal{L}_S, \mathcal{L}_Q, \mathcal{L}_U) = (\text{dsize}, \text{qeq}, \text{usize})$
- Several new constructions have better leakage profiles
 - AZL and FZL [[K.-Moataz-Ohrimenko18](#)]
 - VLH and AVLH [[K.-Moataz19](#)]

Structured Encryption vs. Other Primitives

- Encrypted structures appear implicitly throughout crypto
- Oblivious RAM can be viewed as a
 - response-hiding encrypted array
 - with leakage profile $\Lambda_{\text{ORAM}} = (\mathcal{L}_S, \mathcal{L}_Q, \mathcal{L}_U) = (\text{dsize}, \perp)$
- PIR can be viewed as a
 - response-hiding encrypted array
 - with leakage profile $\Lambda_{\text{PIR}} = (\mathcal{L}_S, \mathcal{L}_Q, \mathcal{L}_U) = (\text{did}, \perp)$
- Garbled gates can be viewed as
 - response-revealing 2x2 arrays
 - $\Lambda_{\text{GG}} = (\mathcal{L}_S, \mathcal{L}_Q, \mathcal{L}_U) = (\text{dsize}, \text{qeq})$

Encrypted Multi-Maps

Encrypted Multi-Maps:

The Heart of *Sub-Linear* Encrypted Search

- EMMs are used as building block for sub-linear
 - Single keyword search [[Curtmola-Garay-K.-Ostrovsky06,...](#)]
 - Conjunctive keyword search [[Cash et al.13,...](#)]
 - Boolean keyword search [[Cash et al.13](#), [K.-Moataz17,...](#)]
 - Range queries [[Faber et al.14](#), [Demertzis et al. 16,...](#)]
 - Substring, wildcard, [[Faber et al.14,...](#)]
 - SQL databases [[K.-Moataz18,...](#)]
 - Graph databases [[Chase-K.10,...](#)]

Pidyn (Modified)

[Cash et al.14]

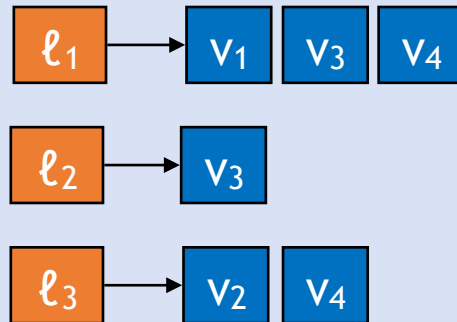
Setup

$$K_{\ell i} = F_K(w_i | 1)$$

EMM.Setup

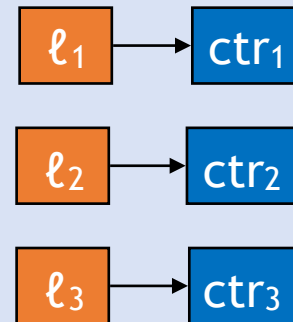
$1^k,$

MM

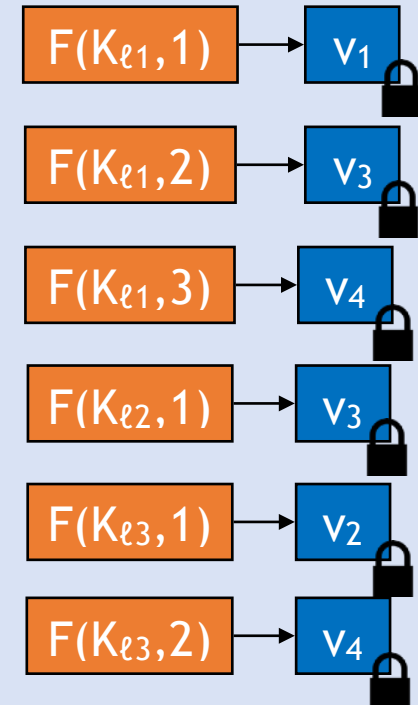


K

DX (state)

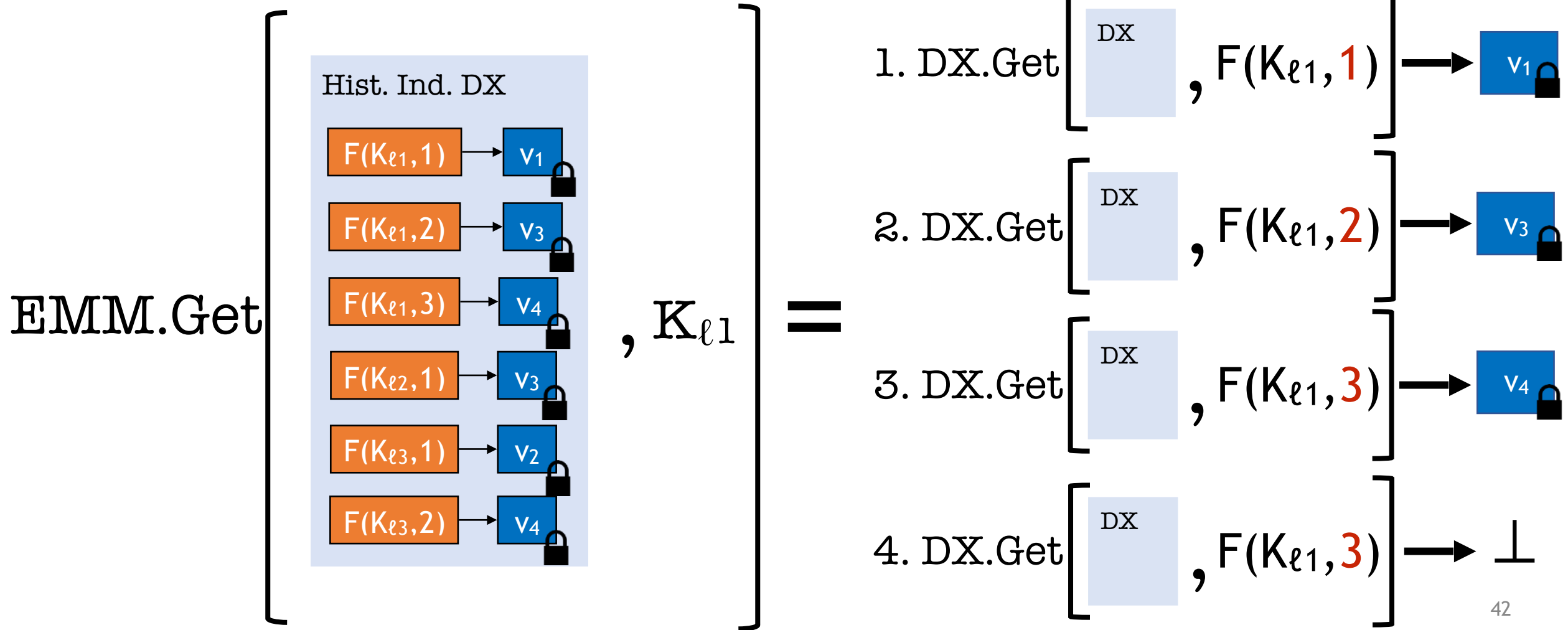


Hist. Ind. DX



Pidyn (Modified)

[Cash et al.14]

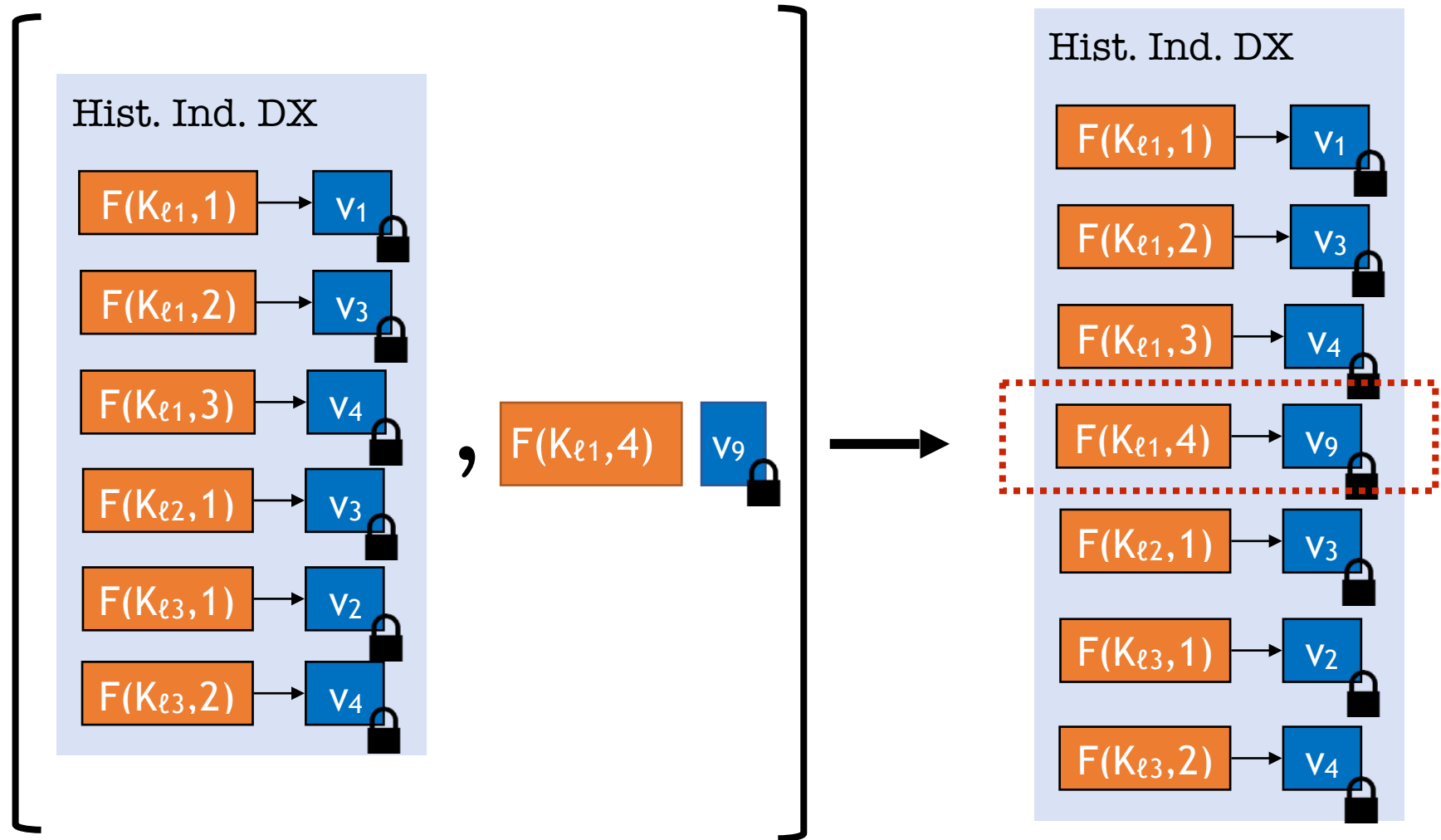


Pidyn (Modified)

[Cash et al.14]

Edit⁺

EMM.Edit⁺

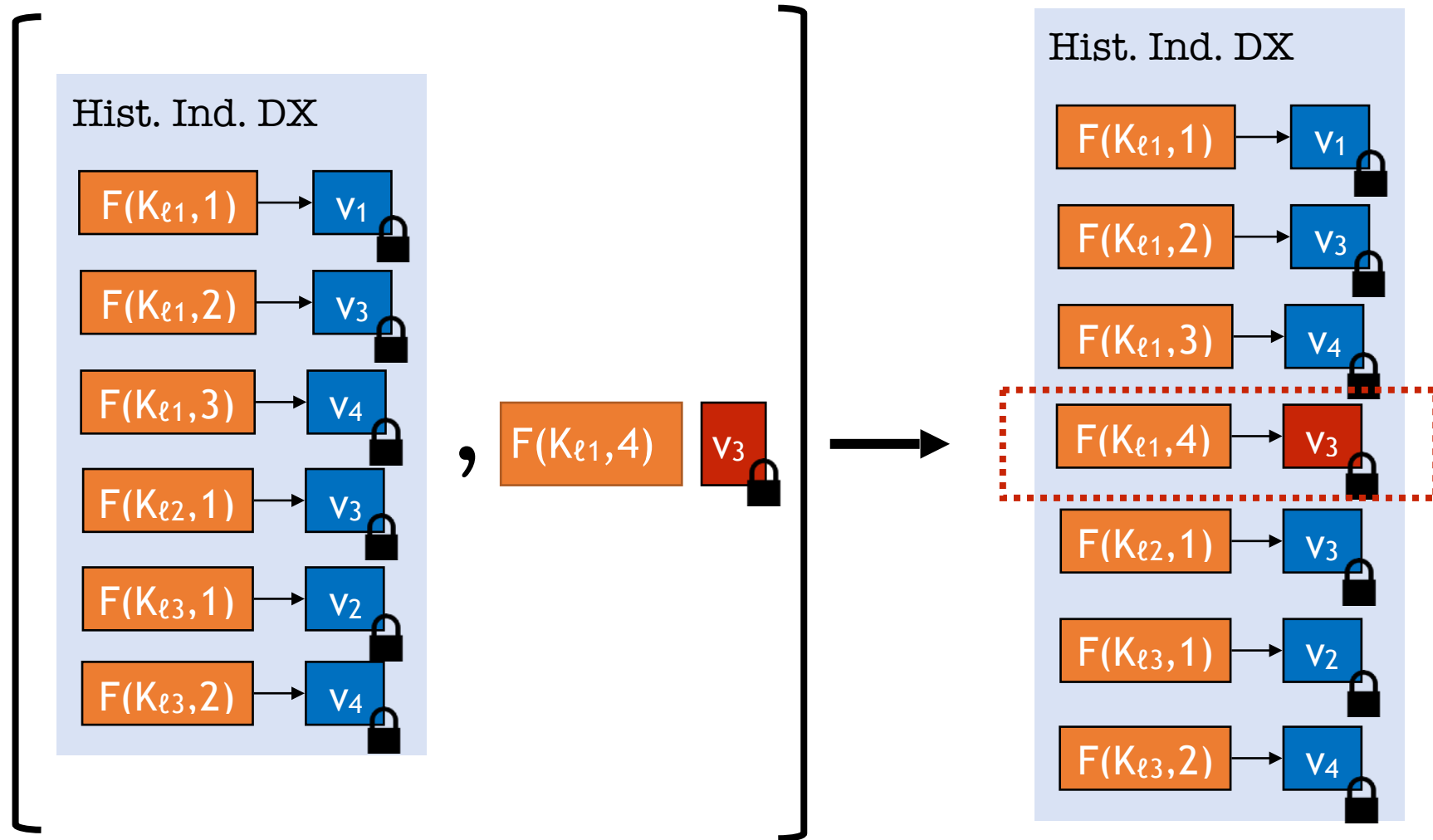


Pidyn (Modified)

[Cash et al.14]

Edit-

EMM.Edit-



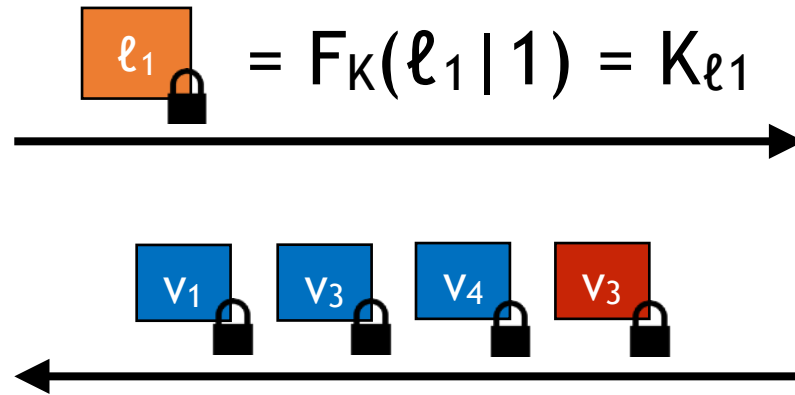
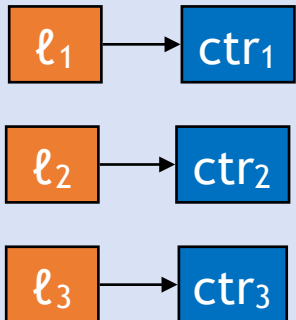
Pidyn (Modified)

[Cash et al.14]



K

DX (state)



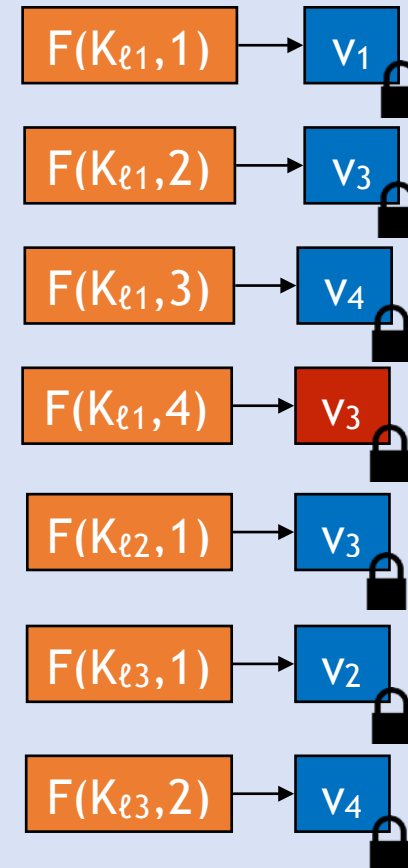
Query complexity:
 $O(\#MM[\ell] + \text{dels}_o(\ell))$

Storage complexity:
 $O(\sum_{\ell} \#MM[\ell] + \text{dels}_o(\ell))$



Get

Hist. Ind. DX



I/O Efficiency & Locality

[[Cash et al.14](#)]

- The problem with large data
 - if data is very large it gets stored on disk
- Disk seeks are very slow
 - minimize *locality*: # of non-contiguous accesses
 - minimize read efficiency: how much *additional* data is read
 - reading *contiguous* data is OK but not too long
- Pidyn has poor locality
 - **Get(ℓ)** needs $\#MM[\ell]$ non-contiguous accesses

I/O Efficiency & Locality

[[Cash et al.14](#)]

- Introduce several schemes with improved locality
 - **Pipack**: packs values in a single ciphertext
 - **Piptr**: packs pointers to values in a single ciphertext
 - this tradeoffs EMM locality for standard memory locality
 - **2Lev**: combines both techniques

Local SSE Schemes

- [[Cash-Tessaro14](#)]
 - lower bounds for “non-overlapping” schemes (improved by Asharov et al.)
- [[Asharov-Segev-Shahaf18](#)]
 - lower bound for “pad-and-split” schemes
 - $L(N)$ locality & $O(1)$ read efficiency $\implies \Omega(N \log N / \log L)$ space
 - matched by [[Demertzis-Papamanthou17](#)]
- [[Asharov-Naor-Segev-Shahaf18](#)]
 - lower bound for “statistically-ind.” schemes
 - $O(1)$ locality & $O(N)$ space $\implies \omega(1) \cdot \epsilon(n)^{-1}$ read efficiency
 - matched by [[Asharov-Segev-Shahaf18](#)]

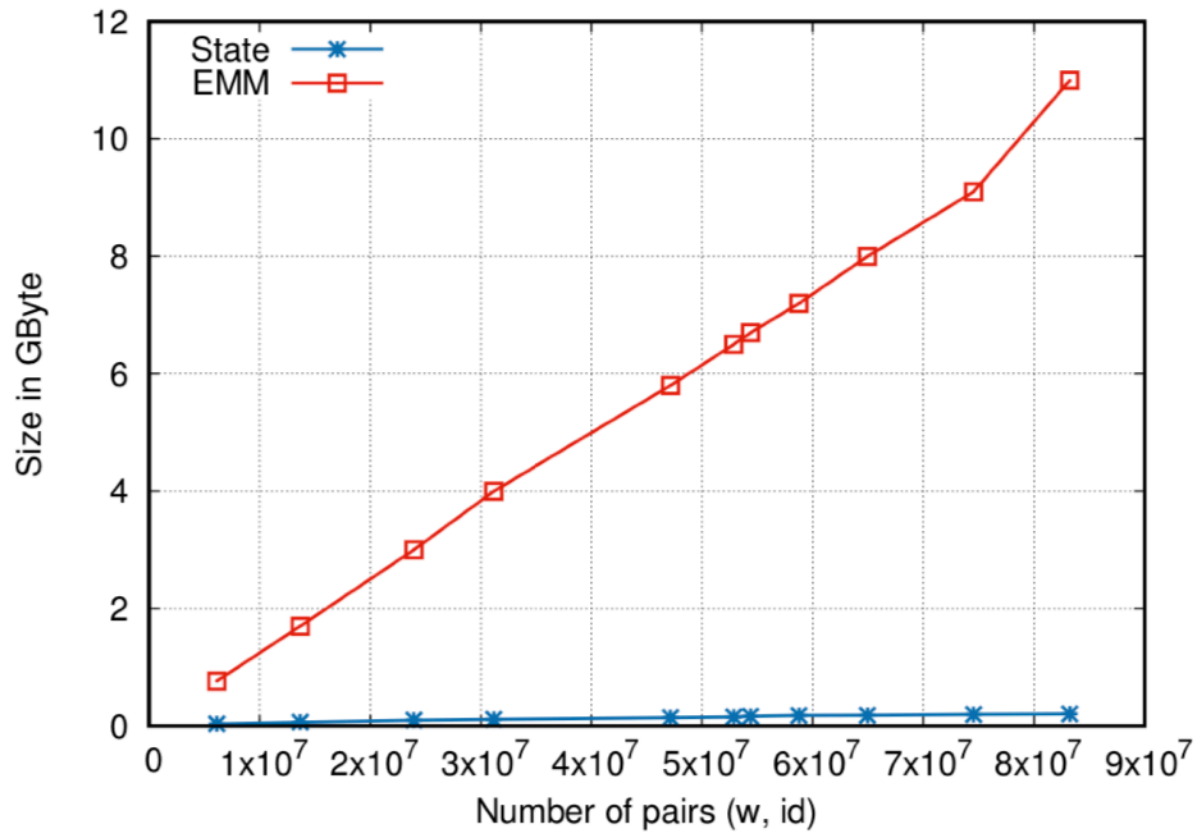
Limitations of Pidyn, Pipack, Piptr, 2Lev

- Not forward private
 - update tokens can be linked to previous search tokens
 - can be exploited using adaptive file injection attacks
- Query and storage complexity depend on total # of deletes

State-of-the-Art EMMs

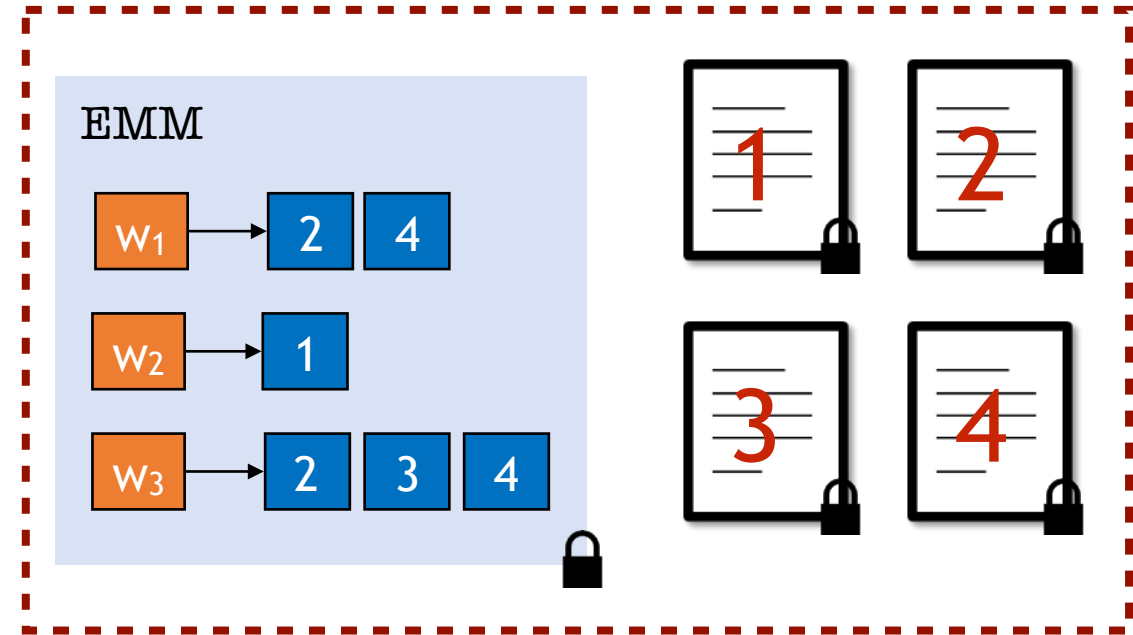
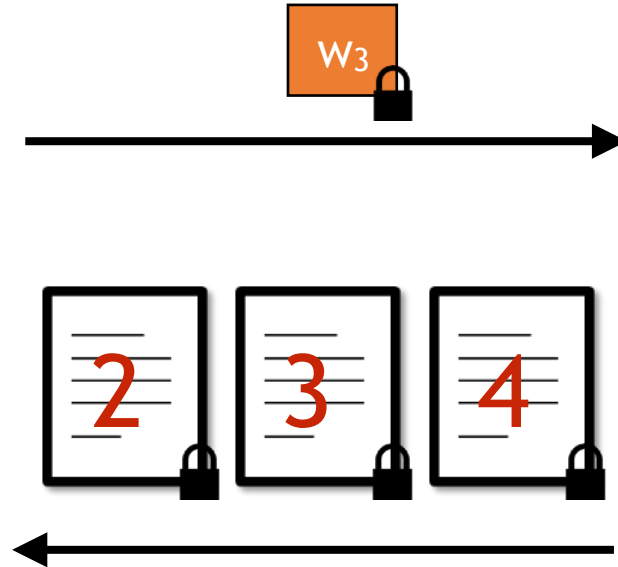
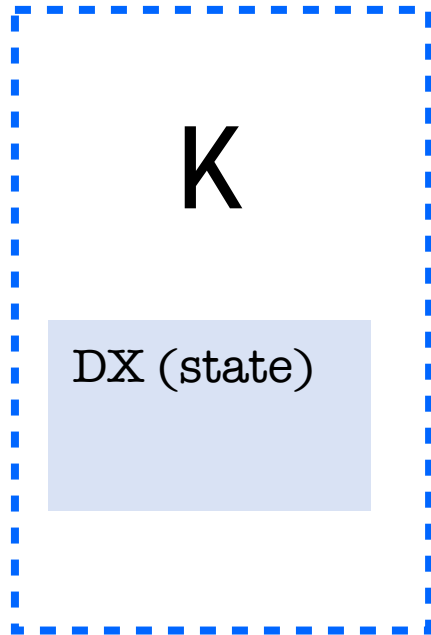
	Search	Client Storage	Forward Privacy	Snapshot
SPS'14	$O(\#MM[\ell] \cdot \text{polylog}(\#MM[\ell]))$	$O(\#\mathbb{L})$	Yes	Yes
B'16	$O(\#MM[\ell] + \text{dels}_0(w))$	$O(\#\mathbb{L})$	Yes	No
BMO'17	$O(\#MM[\ell] + \text{dels}_0(w))$	$O(\#\mathbb{L})$	Yes	No
EKPE'17	$O(\#MM[\ell] + \text{dels}_s(w))$	$O(\#\mathbb{L})$	Yes for adds No for dels	No
AKM19	$O(\#MM[\ell] + \text{dels}_r(w))$	$O(\#\mathbb{L} + ML)$	Yes	Yes

[AKM19] Client State



- EDB w/ 83 million pairs (11GB)
 - state is 210MB

Single Keyword Search from EMMs



Sub-Linear Constructions from Black-Box EMMs

- Searchable symmetric encryption [[Curtmola-Garay-K.-Ostrovsky06,...](#)]
- Graph queries [[Chase-K.10,...](#)]
- Conjunctive & disjunctive keyword search [[Cash et al. 13,](#)]
- Worst-case sub-linear disjunctive & Boolean search [[Pappas et al.14, K.-Moataz17](#)]
- Wildcard & substring search [[Faber et al.15](#)]
- Range search [[Faber et al.15,Demertzis et al.16,Podar-Boelter-Popa16](#)]
- SQL queries on relational DBs [[K.-Moataz18](#)]

Sub-Linear Constructions from Black-Box EMMs

- Why constructions based on black-box EMMs?
- Modularity
 - easy to design, understand and analyze
 - benefit from improvements in EMM efficiency
 - benefit from improvements in EMM security/leakage