

UNLYNX: A DECENTRALIZED SYSTEM FOR PRIVACY-CONSCIOUS DATA SHARING

PETS 2017

DAVID FROELICHER^{*#}, PATRICIA EGGER^{*#}, JOAO SA SOUSA^{*}, JEAN LOUIS RAISARO^{*},
ZHICONG HUANG^{*}, CHRISTIAN MOUCHET^{*}, BRYAN FORD[#], JEAN-PIERRE HUBAUX^{*}



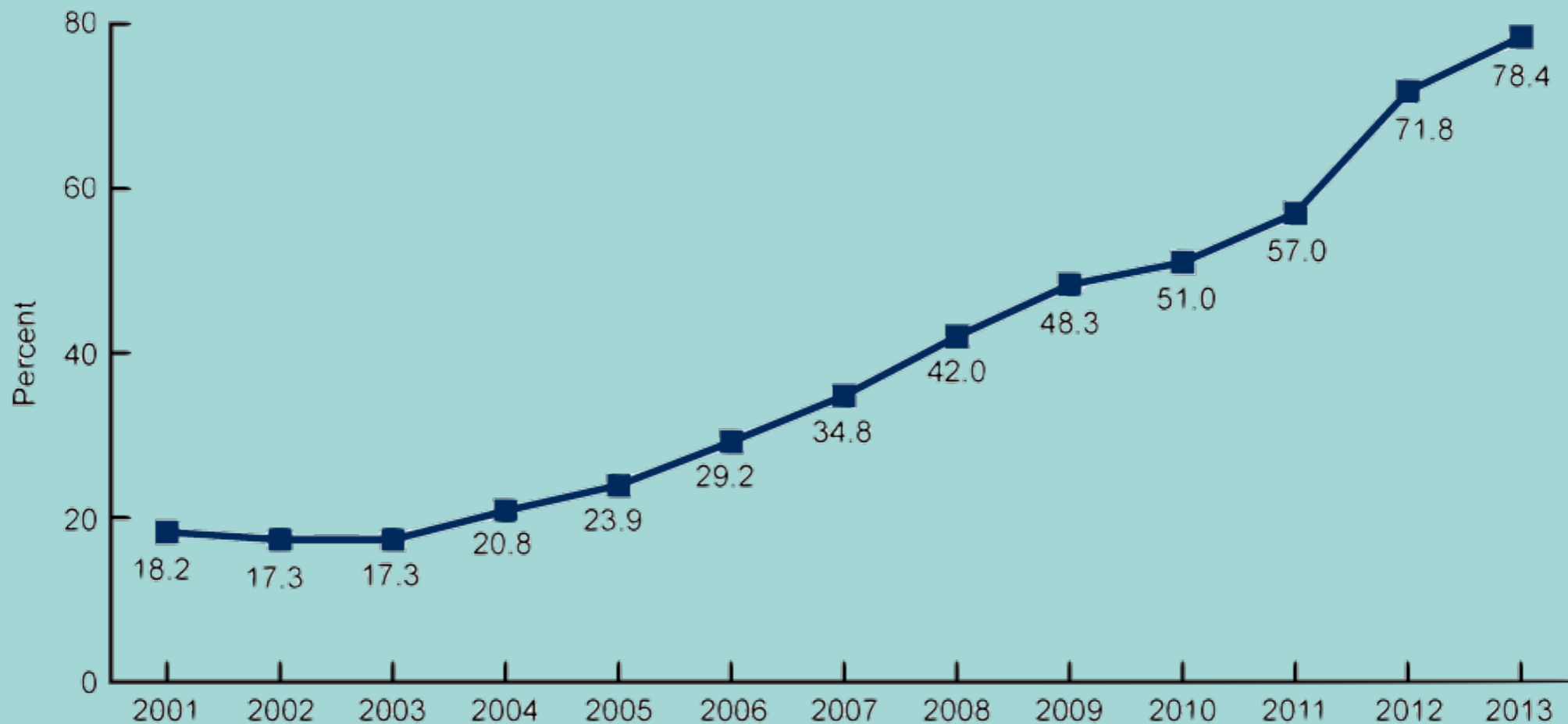
ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

^{*}LCA1 LABORATORY

[#]DeDiS LABORATORY

MORE MEDICAL DATA ARE DIGITIZED

PERCENTAGE OF OFFICE-BASED PHYSICIANS WITH ELECTRONIC MEDICAL RECORDS IN U.S.A, 2001-2013

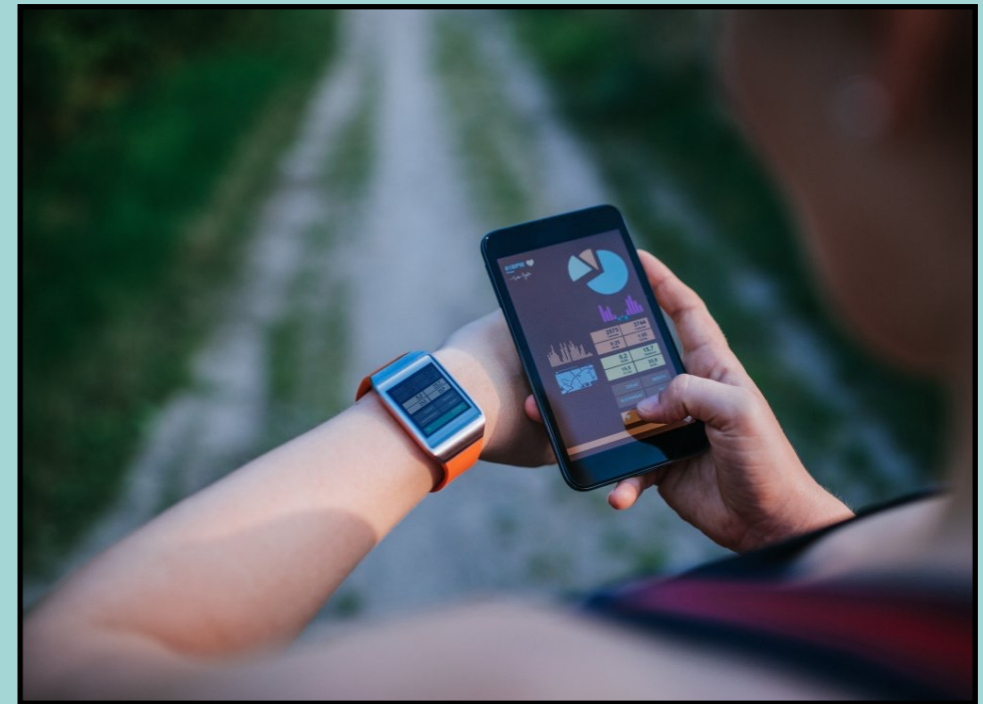


NATIONAL AMBULATORY MEDICAL CARE SURVEY (NAMCS)

MORE HEALTH DATA COLLECTED



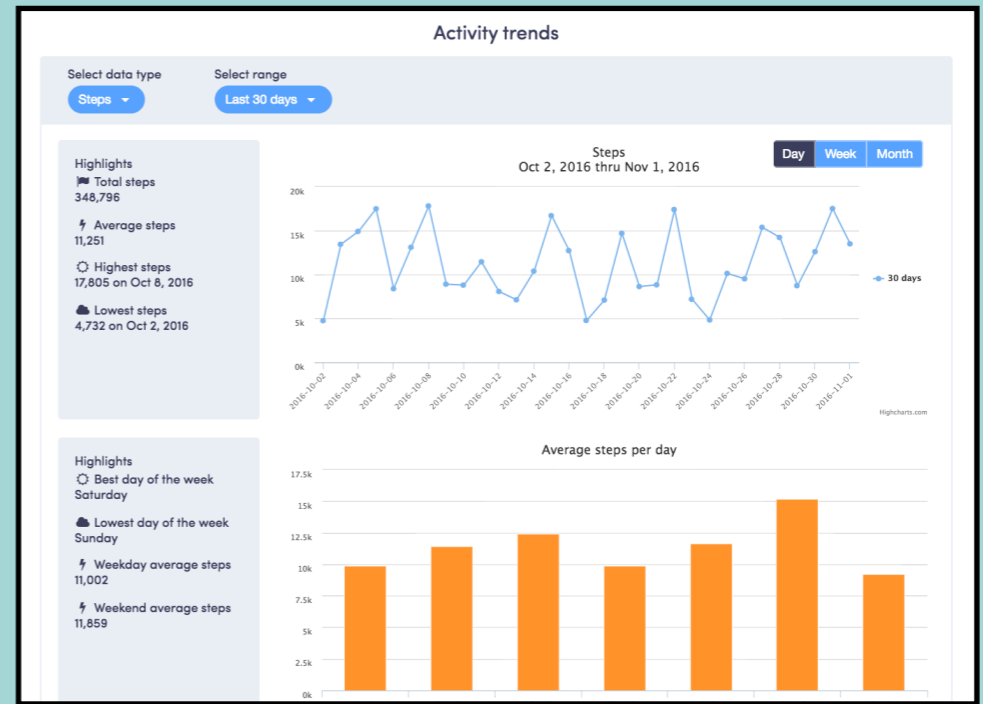
<http://blog.stridekick.com/ultimate-guide-fitness-tracker-hacks-get-most-from-fitbit/>



<http://www.consumerreports.org/cro/news/2015/06/what-you-need-to-know-about-sharing-your-medical-data/index.htm>



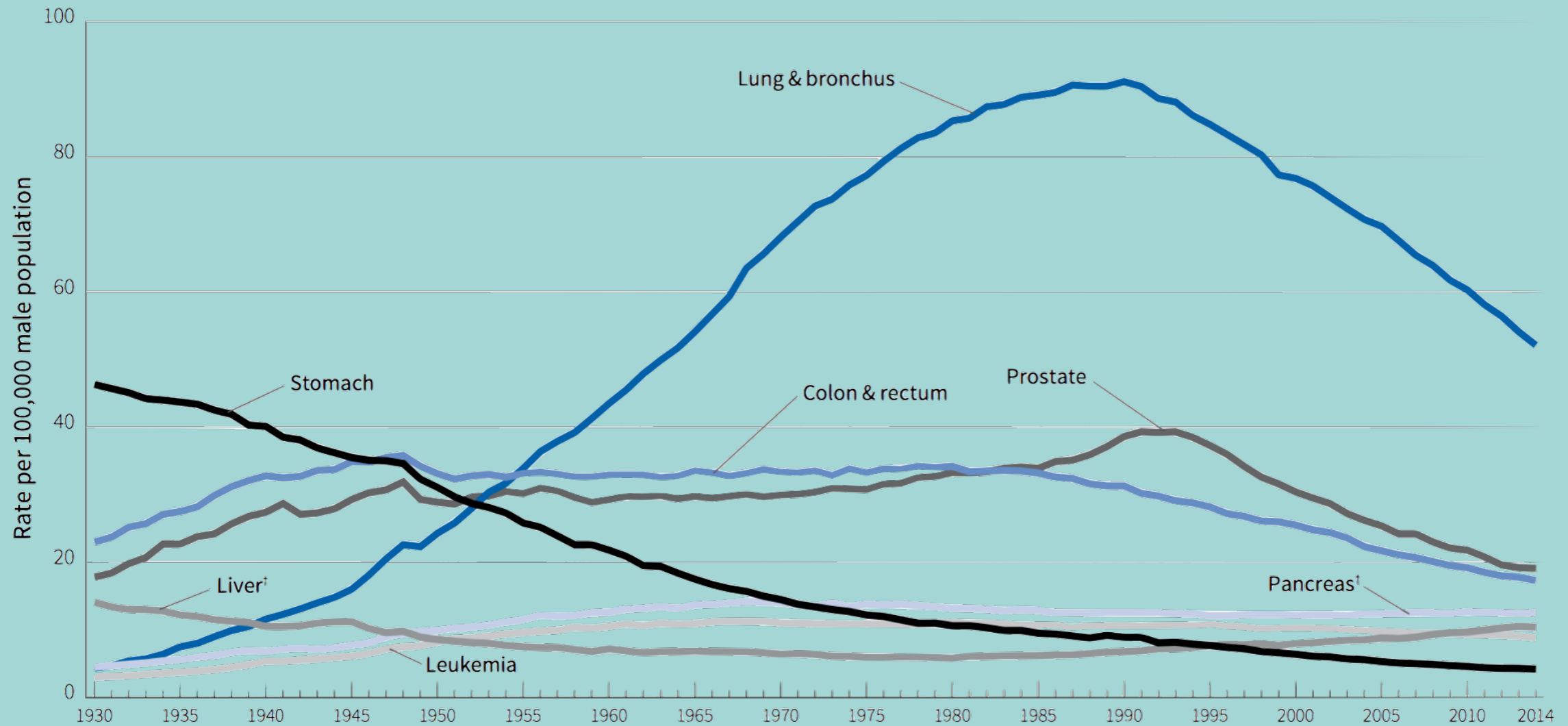
<http://time.com/collection-post/3615161/sharing-health-data/>



<http://www.designindaba.com/articles/creative-work/smart-thermometer-crowdsources-info-real-time-health-tracking>

MORE MEDICAL DATA = BETTER TREATMENTS ?

CANCER DEATH RATES* AMONG MEN, USA, 1930-2014



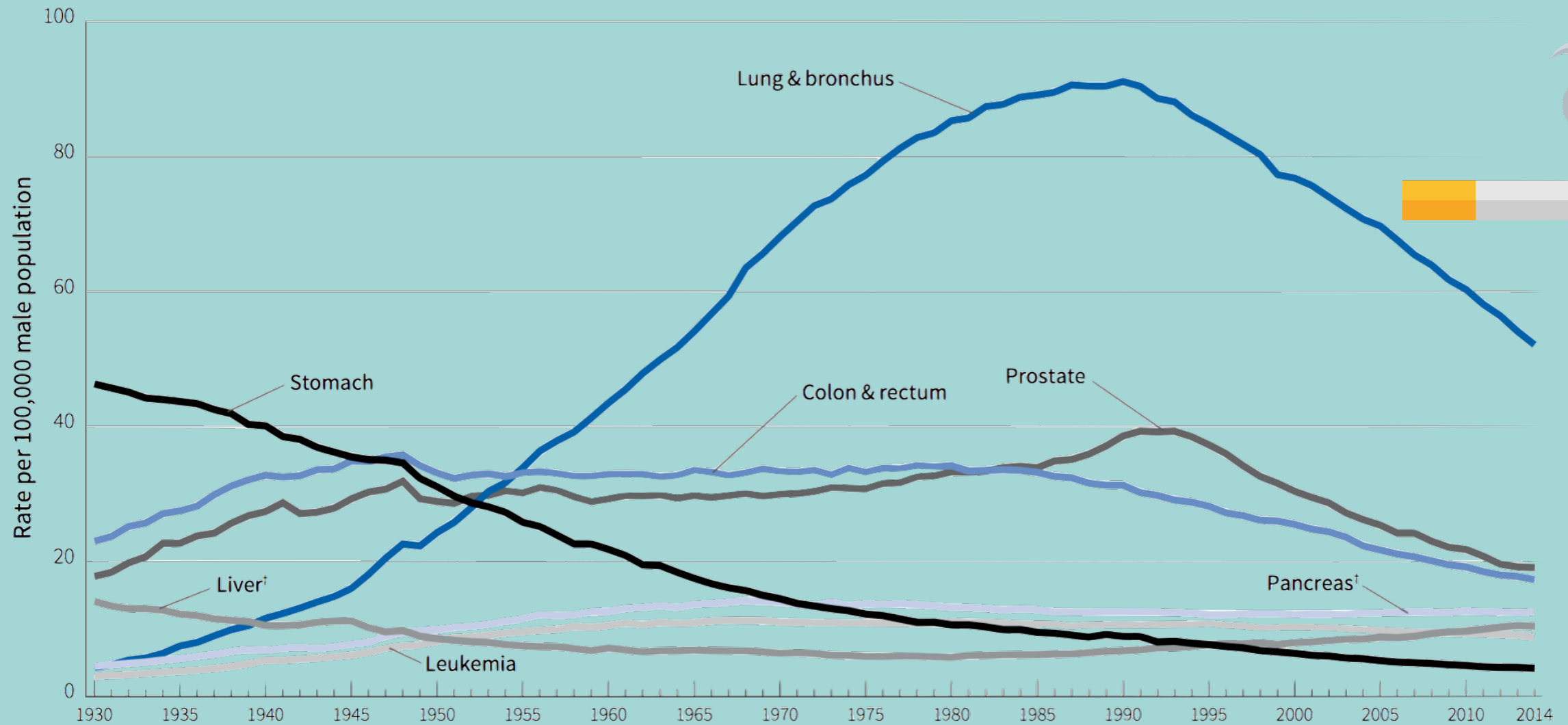
*Per 100,000, age adjusted to the 2000 US standard population. †Mortality rates for pancreatic and liver cancers are increasing.

Note: Due to changes in ICD coding, numerator information has changed over time. Rates for cancers of the liver, lung and bronchus, uterus, and colon and rectum are affected by these coding changes.

Source: US Mortality Volumes 1930 to 1959 and US Mortality Data 1960 to 2014, National Center for Health Statistics, Centers for Disease Control and Prevention.

MORE MEDICAL DATA = BETTER TREATMENTS ?

CANCER DEATH RATES* AMONG MEN, USA, 1930-2014



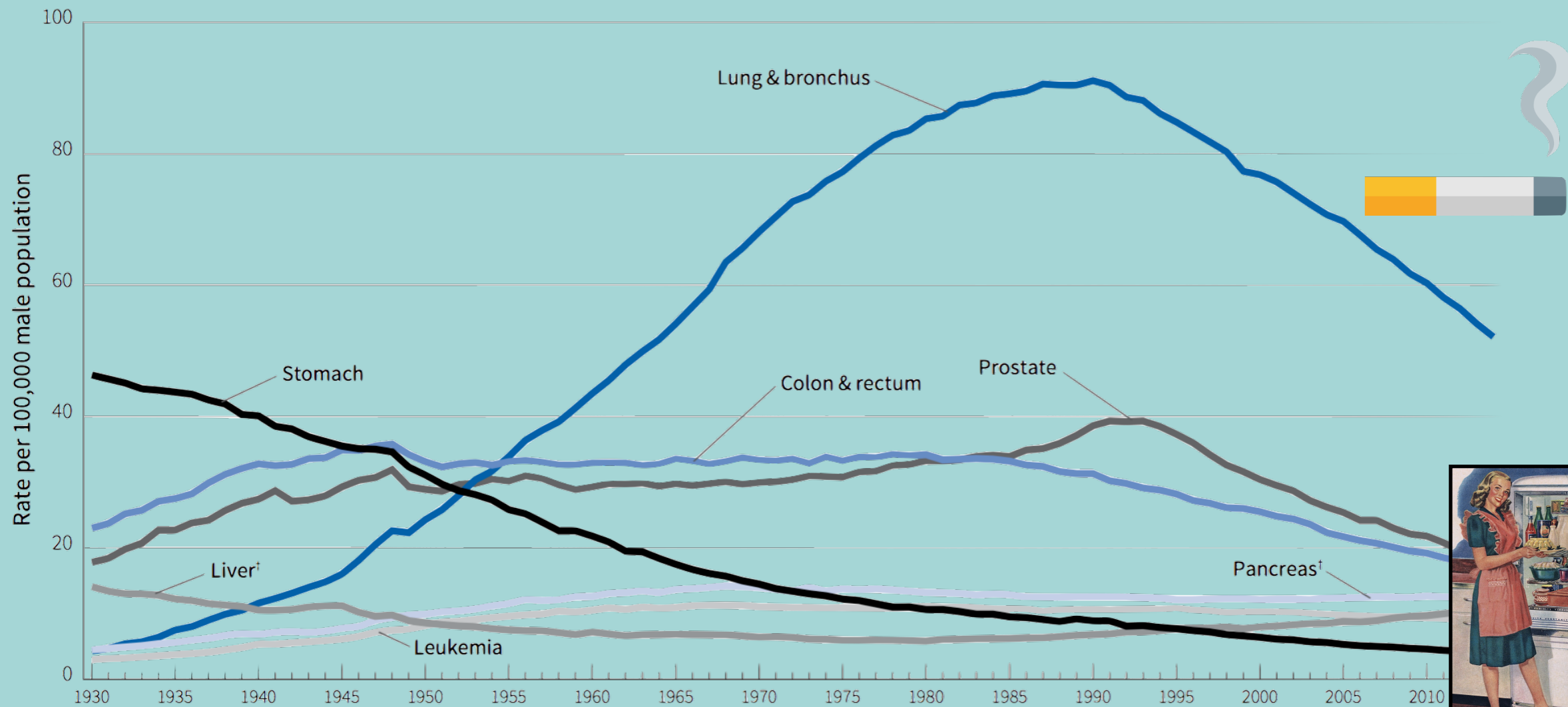
*Per 100,000, age adjusted to the 2000 US standard population. †Mortality rates for pancreatic and liver cancers are increasing.

Note: Due to changes in ICD coding, numerator information has changed over time. Rates for cancers of the liver, lung and bronchus, uterus, and colon and rectum are affected by these coding changes.

Source: US Mortality Volumes 1930 to 1959 and US Mortality Data 1960 to 2014, National Center for Health Statistics, Centers for Disease Control and Prevention.

MORE MEDICAL DATA = BETTER TREATMENTS ?

CANCER DEATH RATES* AMONG MEN, USA, 1930-2014



*Per 100,000, age adjusted to the 2000 US standard population. †Mortality rates for pancreatic and liver cancers are increasing.

Note: Due to changes in ICD coding, numerator information has changed over time. Rates for cancers of the liver, lung and bronchus, uterus, and colon and rectum are affected by these coding changes.

Source: US Mortality Volumes 1930 to 1959 and US Mortality Data 1960 to 2014, National Center for Health Statistics, Centers for Disease Control and Prevention.

SENSITIVE-DATA SHARING IS DIFFICULT



SENSITIVE-DATA SHARING IS DIFFICULT



<http://www.gmill.net/proje-Grain-Storage-Silos>

UNLYNX

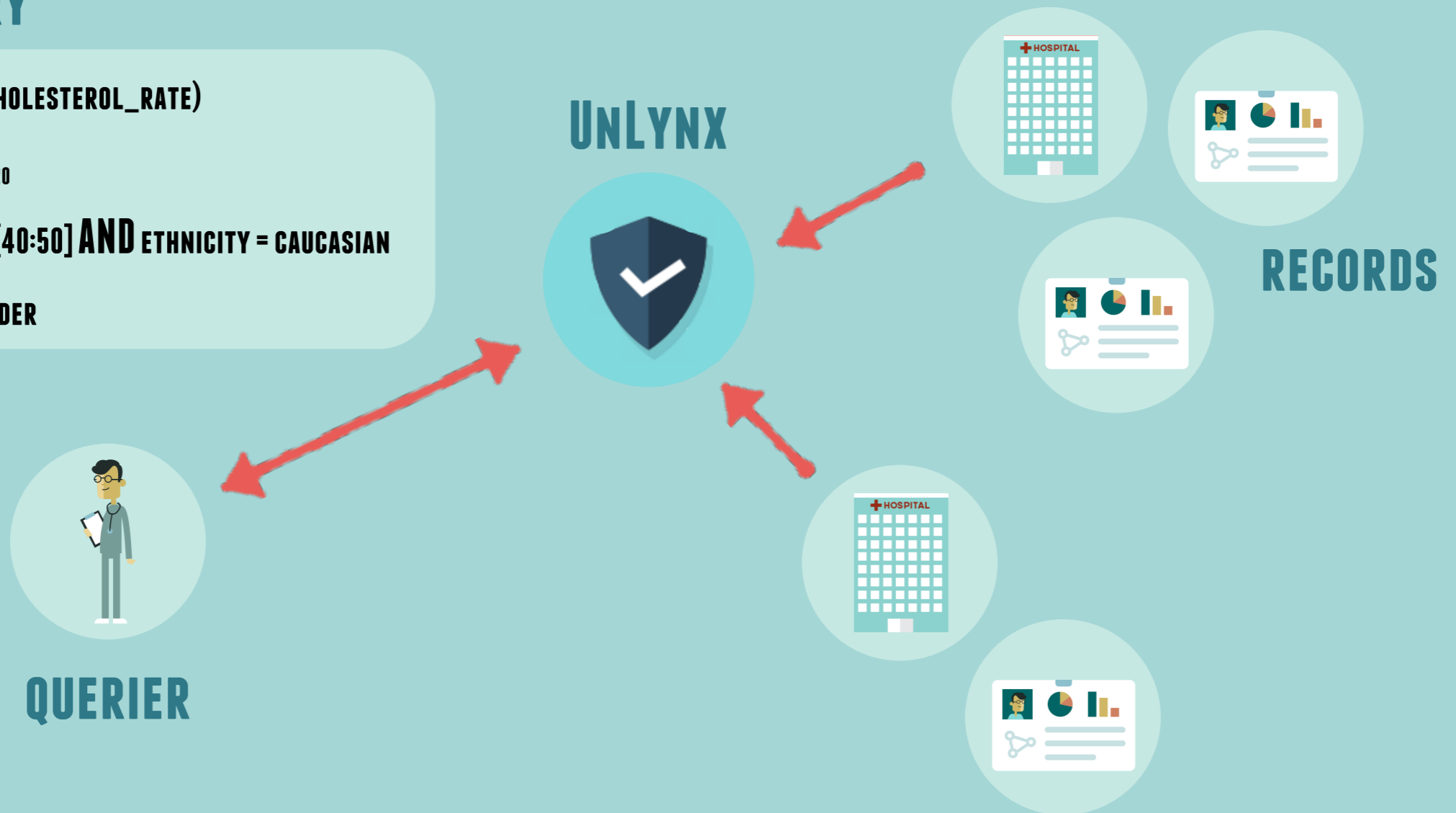
QUERY

```
SELECT AVG (CHOLESTEROL_RATE)
FROM DP1,...,DP20
WHERE AGE IN [40:50] AND ETHNICITY = CAUCASIAN
GROUP BY GENDER
```

DATA PROVIDERS

UNLYNX

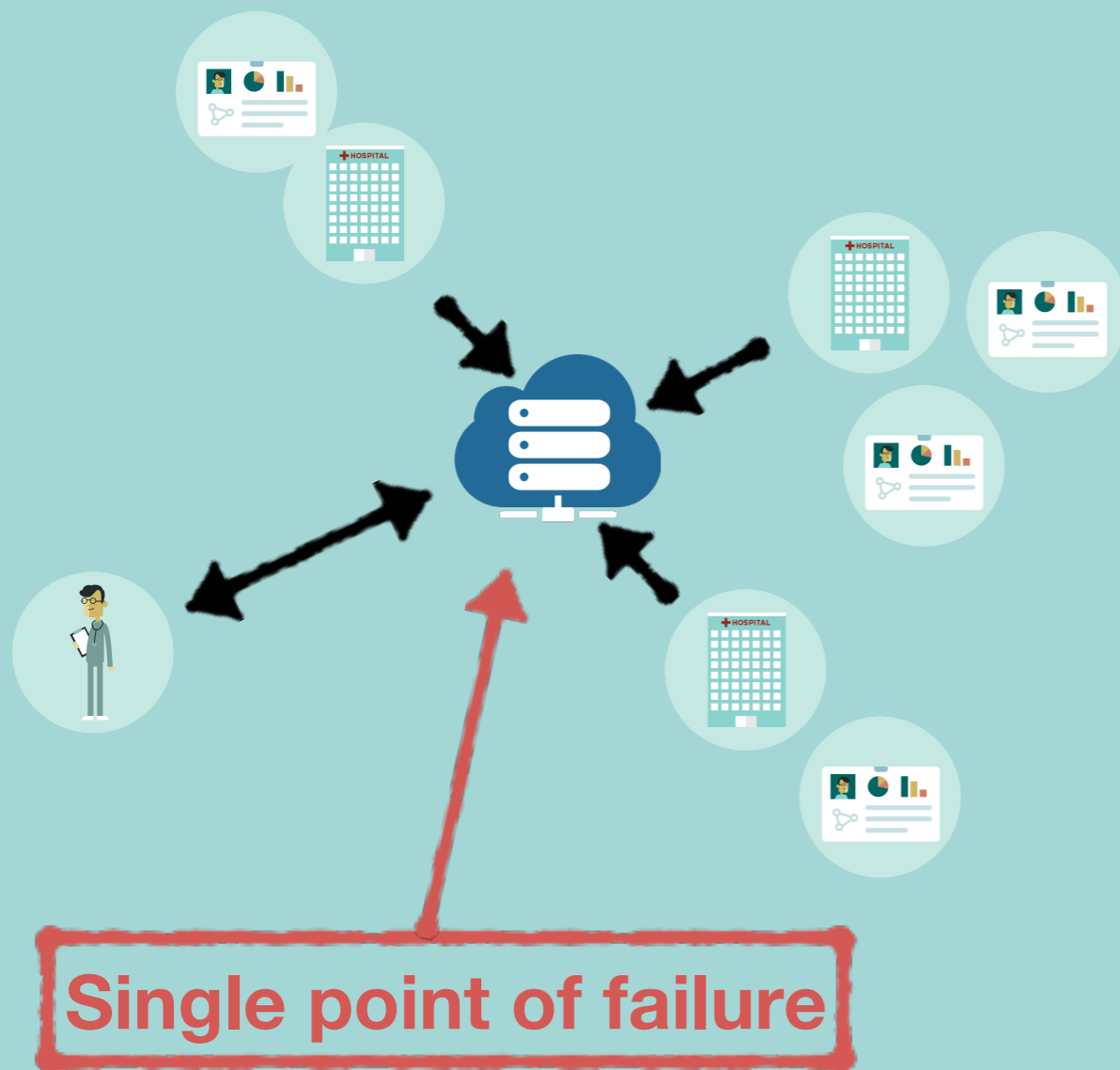
RECORDS



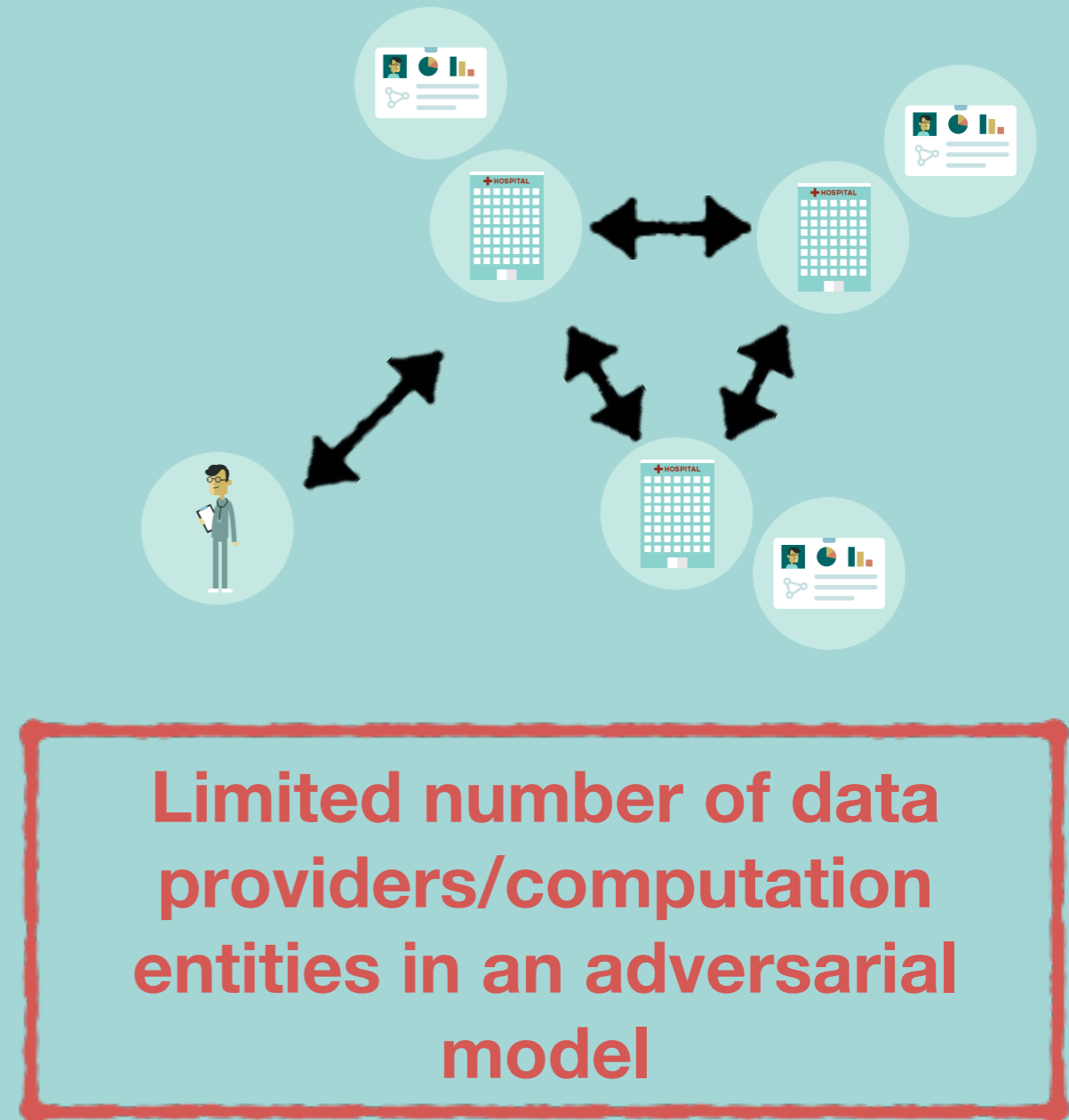
Allow statistical queries on multiple independent databases while ensuring privacy and confidentiality for data providers.

EXISTING DATA SHARING SOLUTIONS

CENTRALIZED SOLUTIONS



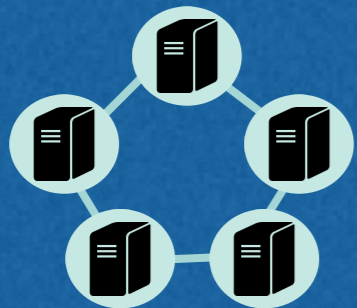
DECENTRALIZED SOLUTIONS



REQUIREMENTS



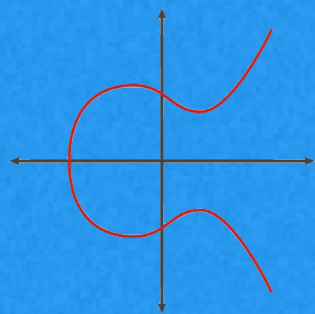
BUILDING BLOCKS



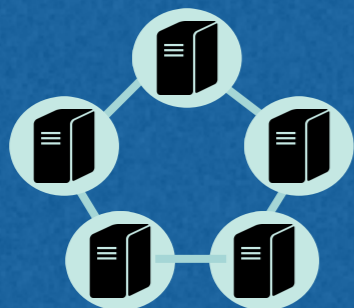
Collective Authority



BUILDING BLOCKS



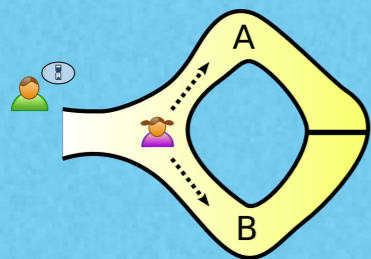
**Additively-homomorphic
ElGamal crypto scheme**



Collective Authority

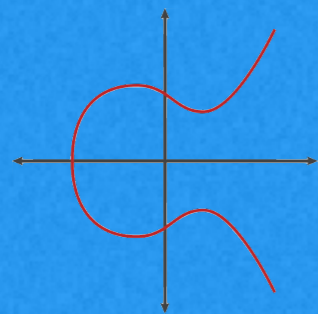


BUILDING BLOCKS

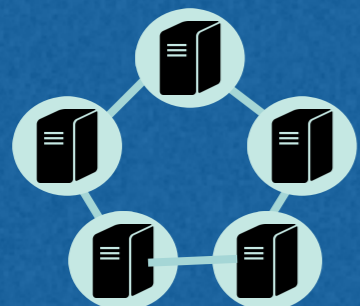


https://simple.wikipedia.org/wiki/Zero-knowledge_proof

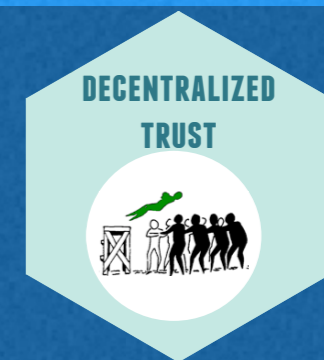
Zero-Knowledge Proofs of Correctness



Additively-homomorphic ElGamal crypto scheme



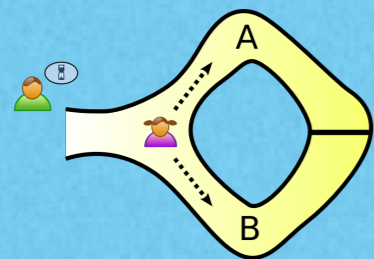
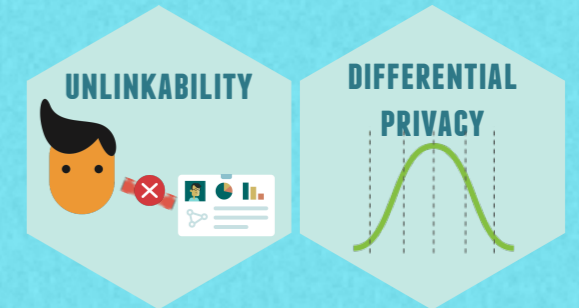
Collective Authority



BUILDING BLOCKS

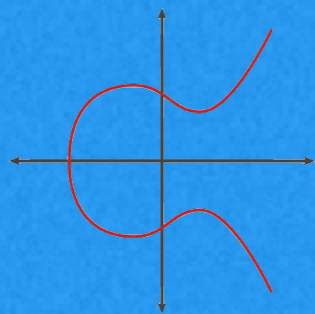


Verifiable Shuffle

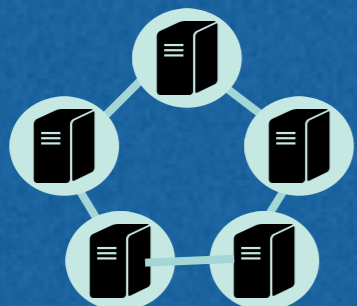


Zero-Knowledge Proofs of Correctness

https://simple.wikipedia.org/wiki/Zero-knowledge_proof



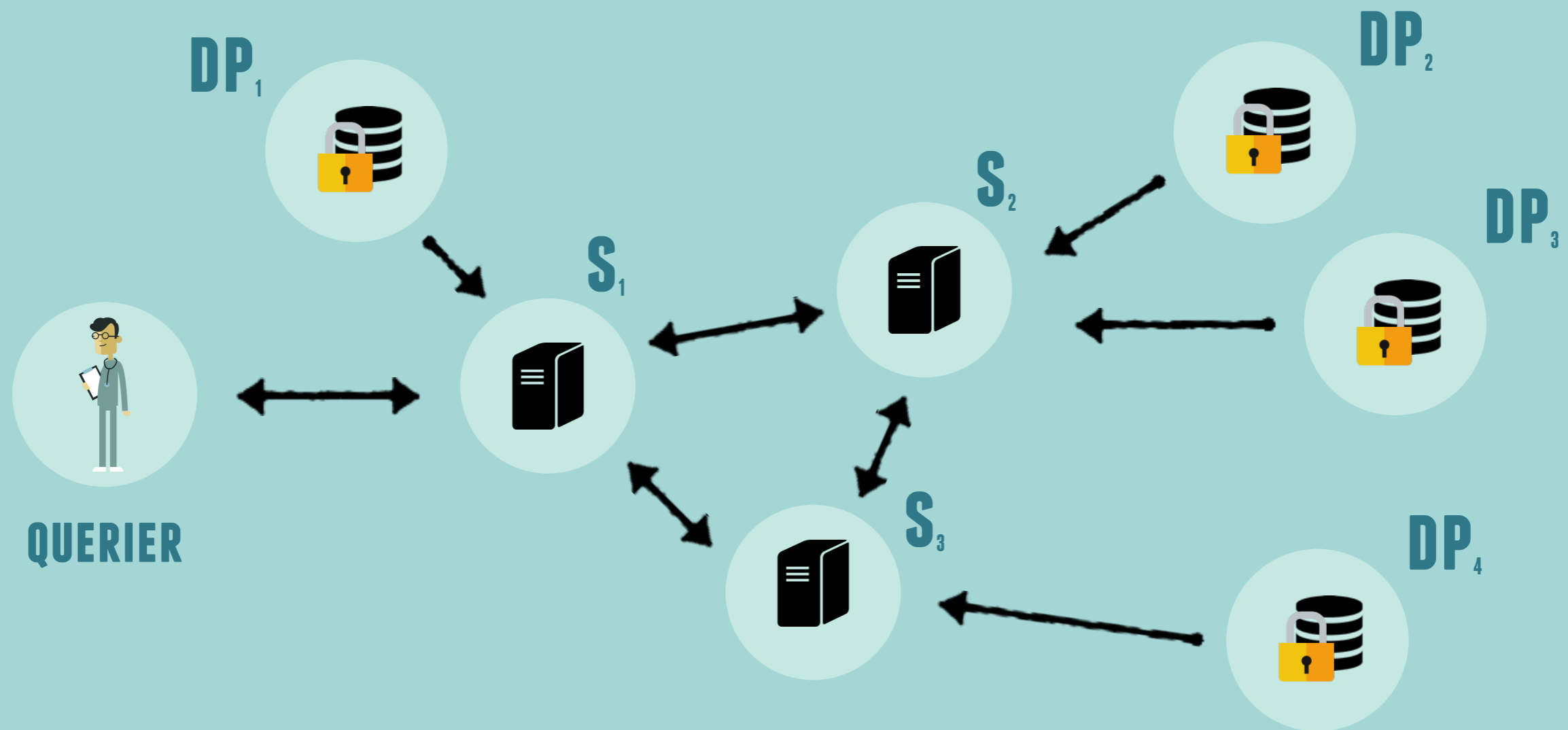
Additively-homomorphic ElGamal crypto scheme



Collective Authority



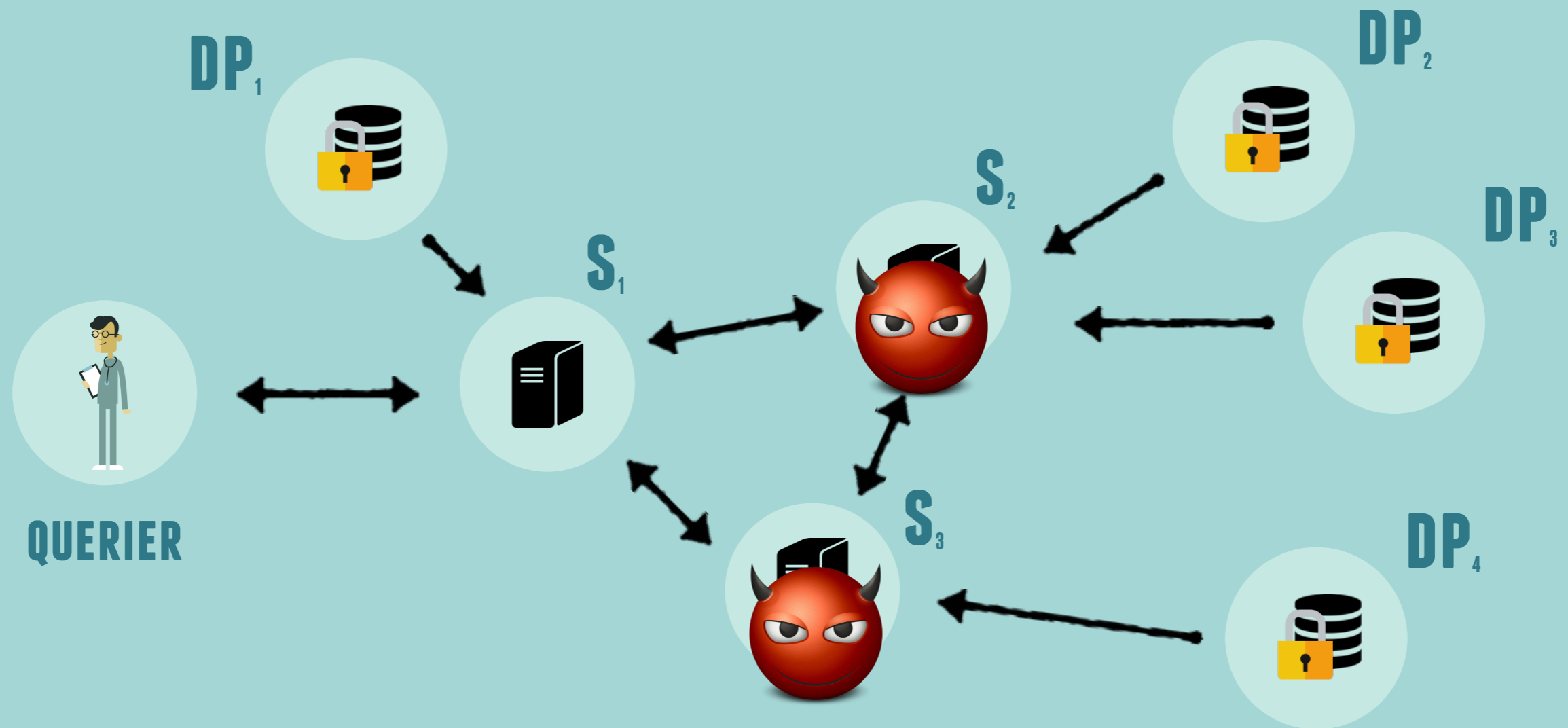
SYSTEM MODEL



- Collective authority of m servers S
- n Data Providers DPs
- Clients Q querying the system

DP = DATA PROVIDER
S = SERVER

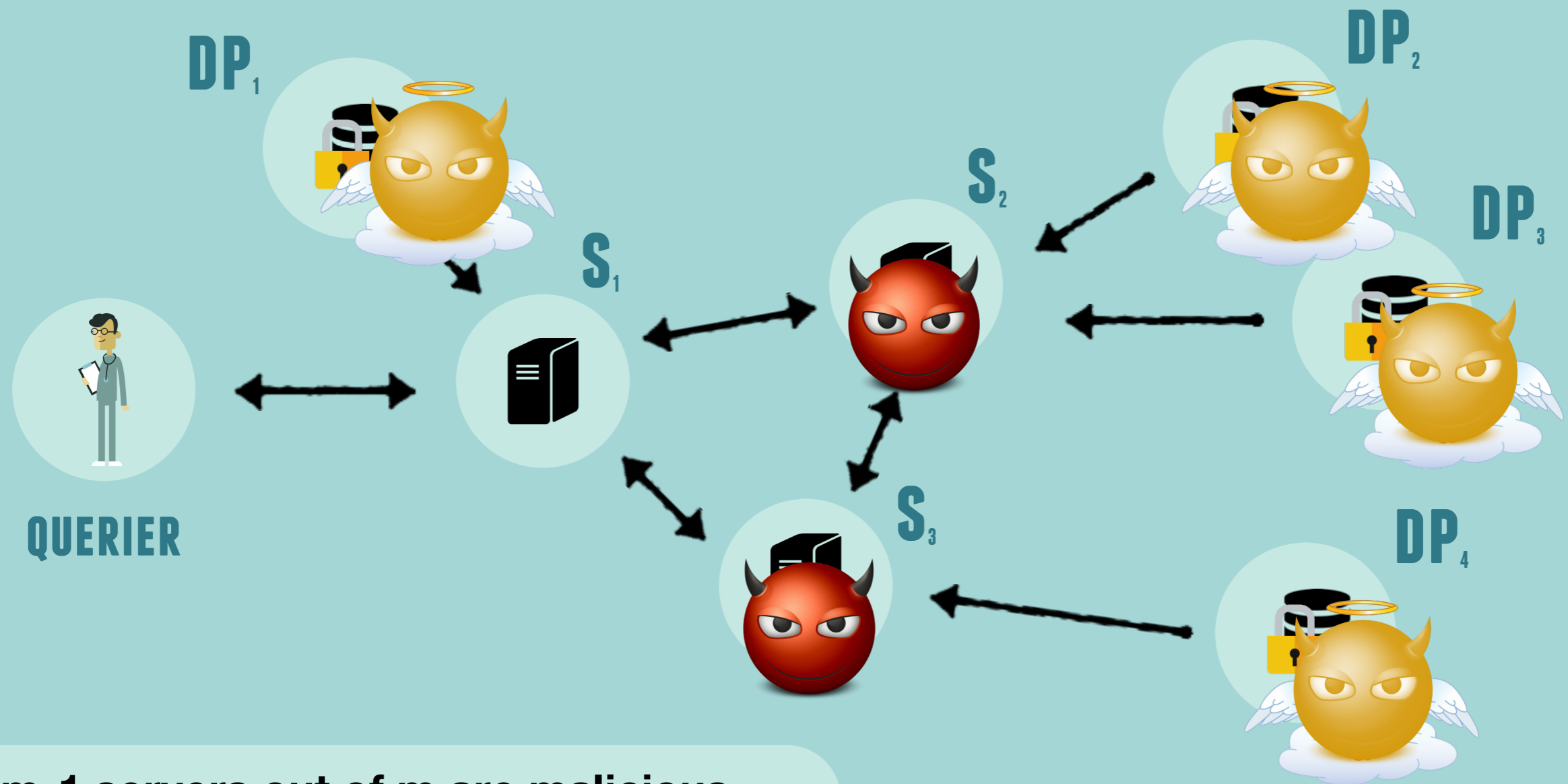
THREAT MODEL



- **m-1 servers out of m are malicious (Anytrust Model)**

DP = DATA PROVIDER
S = SERVER

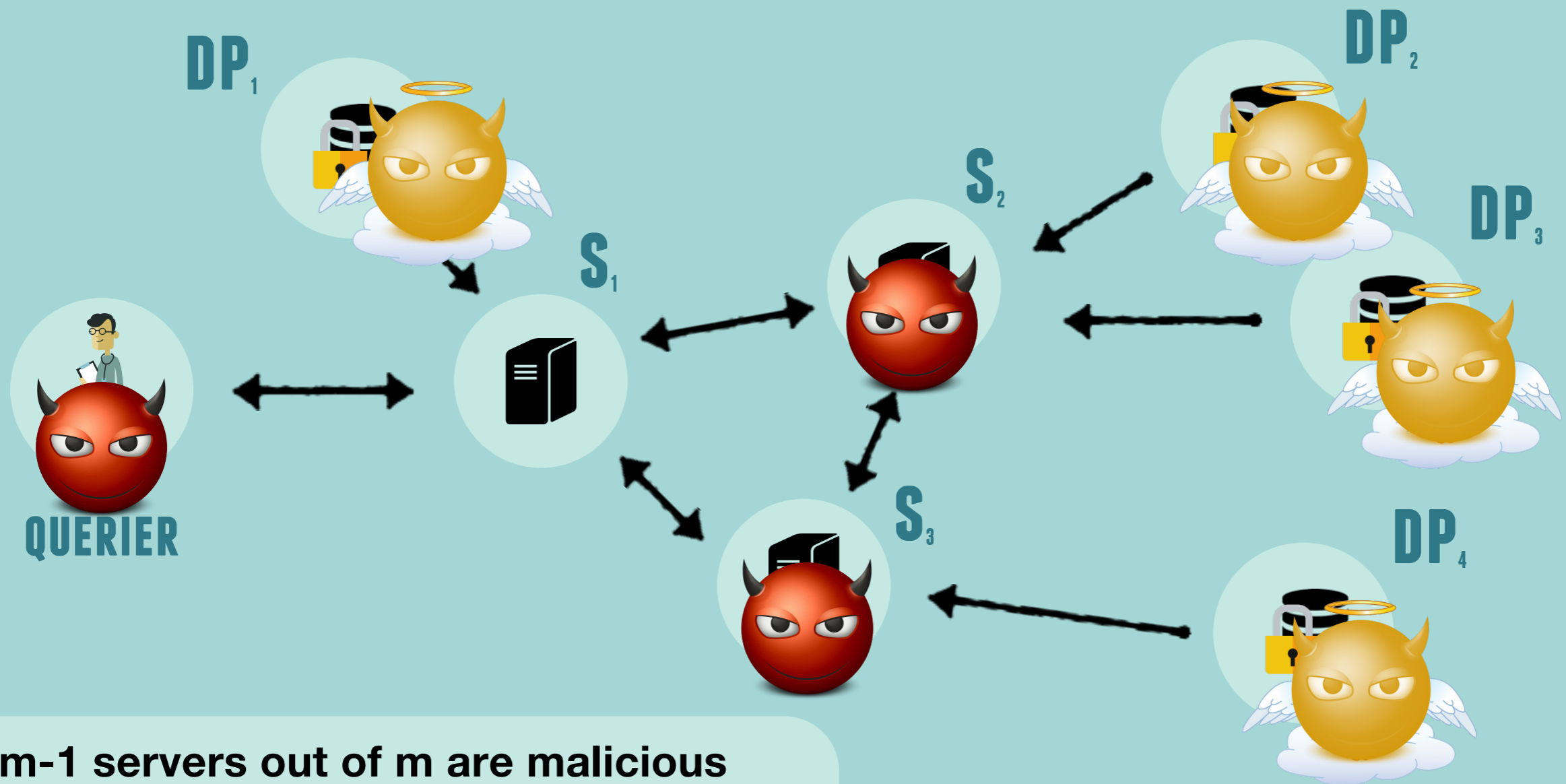
THREAT MODEL



- **m-1 servers out of m are malicious (Anytrust Model)**
- **Data Providers are honest-but-curious**

DP = DATA PROVIDER
S = SERVER

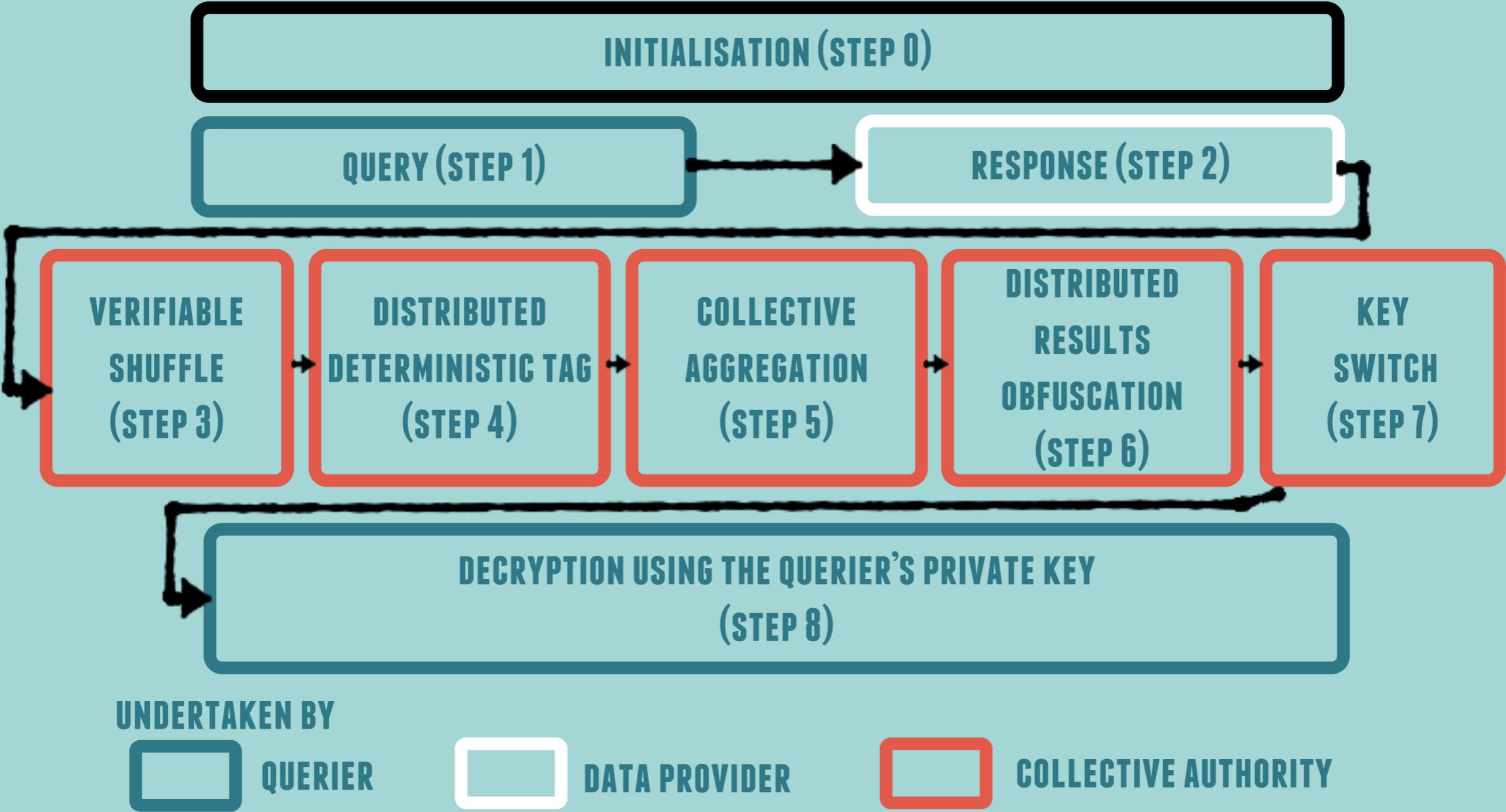
THREAT MODEL



- $m-1$ servers out of m are malicious (Anytrust Model)
- Data Providers are honest-but-curious
- Queriers are malicious

DP = DATA PROVIDER
S = SERVER

QUERY PROCESSING WORKFLOW

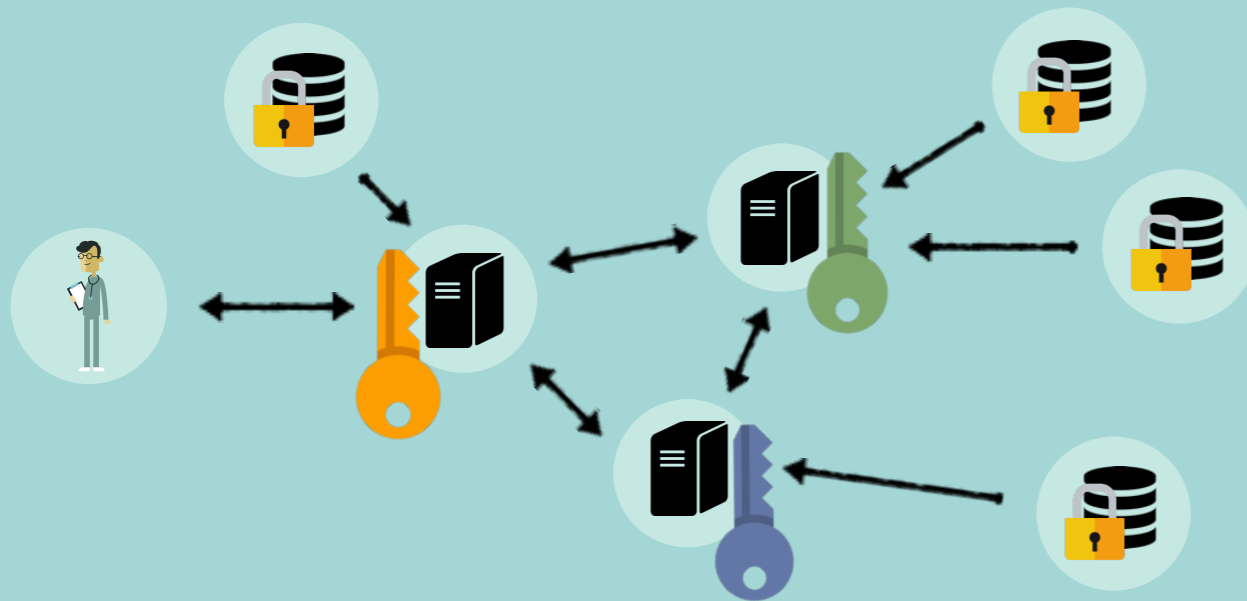


WORKFLOW - INITIALISATION (STEP 0)

INITIALISATION (STEP 0)



Each server constructs his public-private ElGamal Key pair.







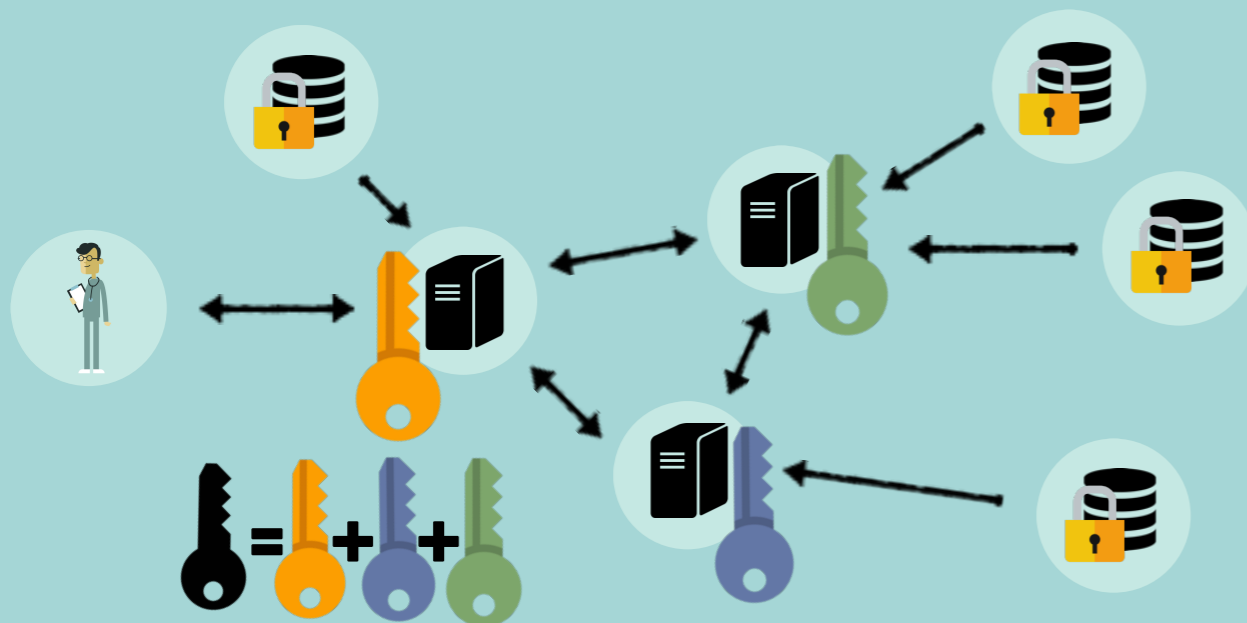
WORKFLOW - INITIALISATION (STEP 0)

INITIALISATION (STEP 0)



Each server constructs his public-private ElGamal Key pair.

Collective Key:  =  +  + 



WORKFLOW - INITIALISATION (STEP 0)

INITIALISATION (STEP 0)

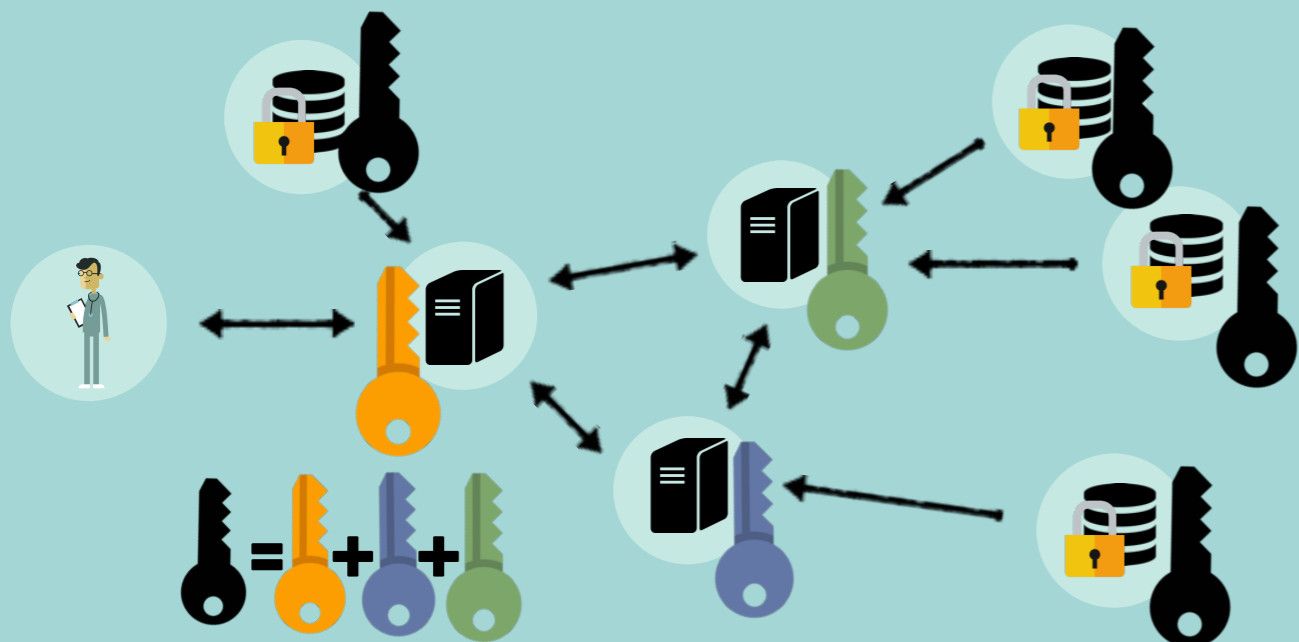


Each server constructs his public-private ElGamal Key pair.

Collective Key:  =  +  + 



Data Providers use the Collective Key to encrypt their data



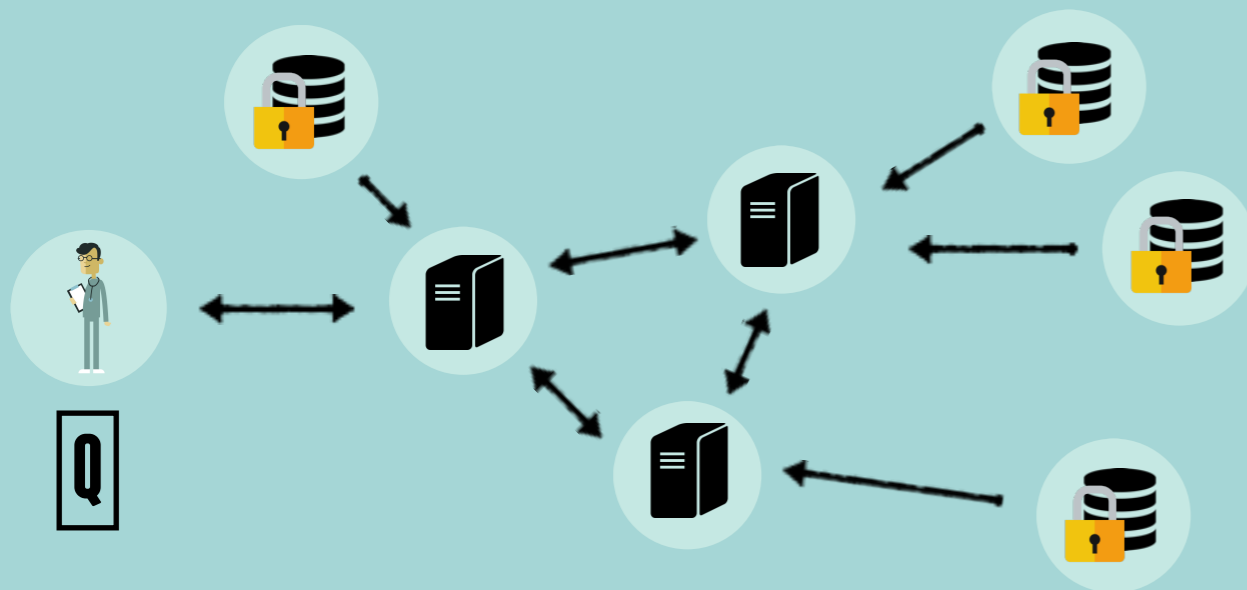
WORKFLOW - QUERY (STEP 1)

INITIALISATION (STEP 0)

QUERY (STEP 1)



```
SELECT SUM (CHOLESTEROL_RATE), COUNT(*)  
FROM DP1,...,DP20  
WHERE AGE IN [40:50] AND ETHNICITY = CAUCASIAN  
GROUP BY GENDER
```



WORKFLOW - QUERY (STEP 1)

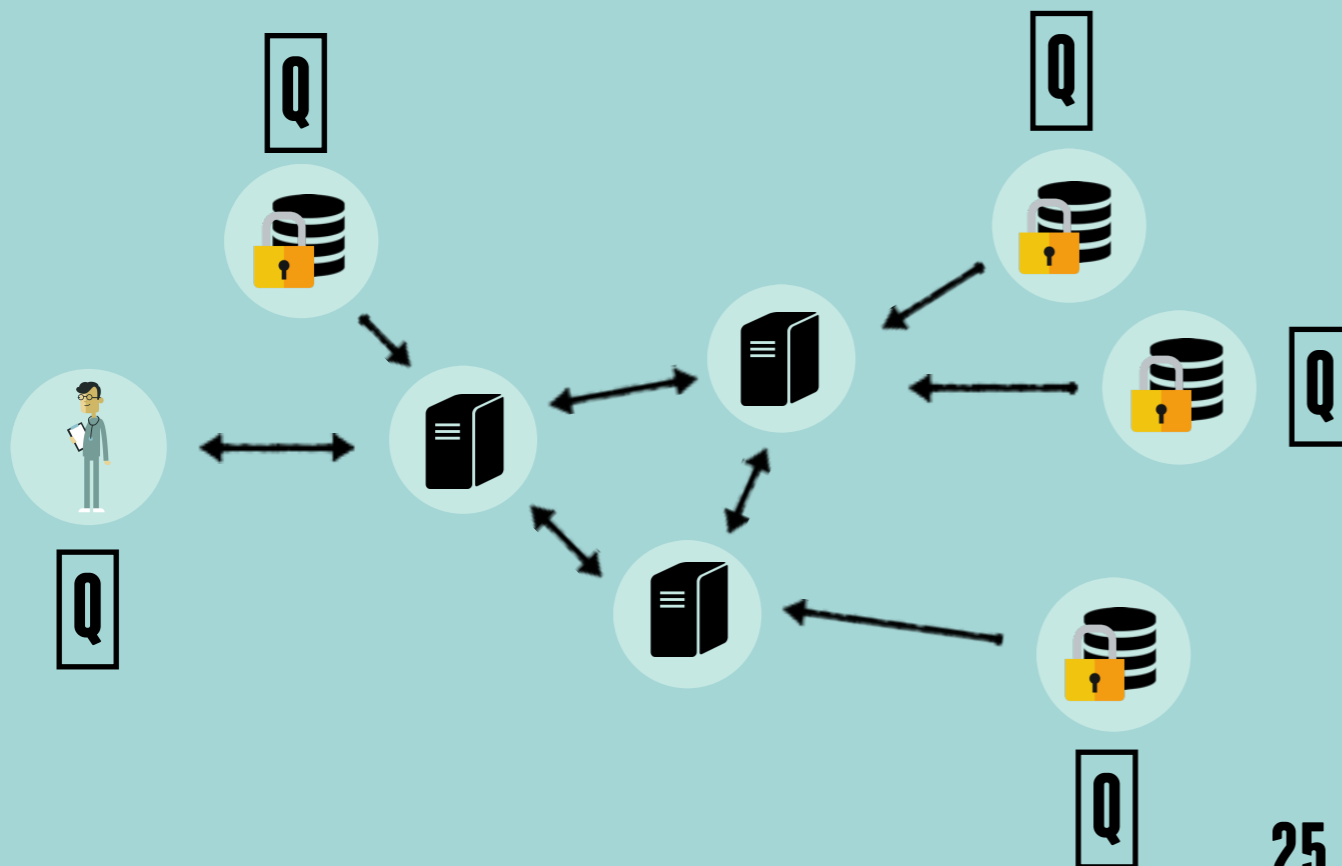
INITIALISATION (STEP 0)

QUERY (STEP 1)

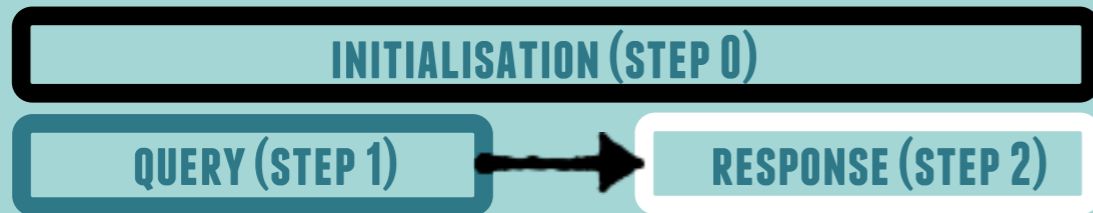


```
SELECT SUM (CHOLESTEROL_RATE), COUNT(*)  
FROM DP1,...,DP20  
WHERE AGE IN [40:50] AND ETHNICITY = CAUCASIAN  
GROUP BY GENDER
```

Query broadcasted to Data Providers



WORKFLOW - RESPONSE (STEP 2)



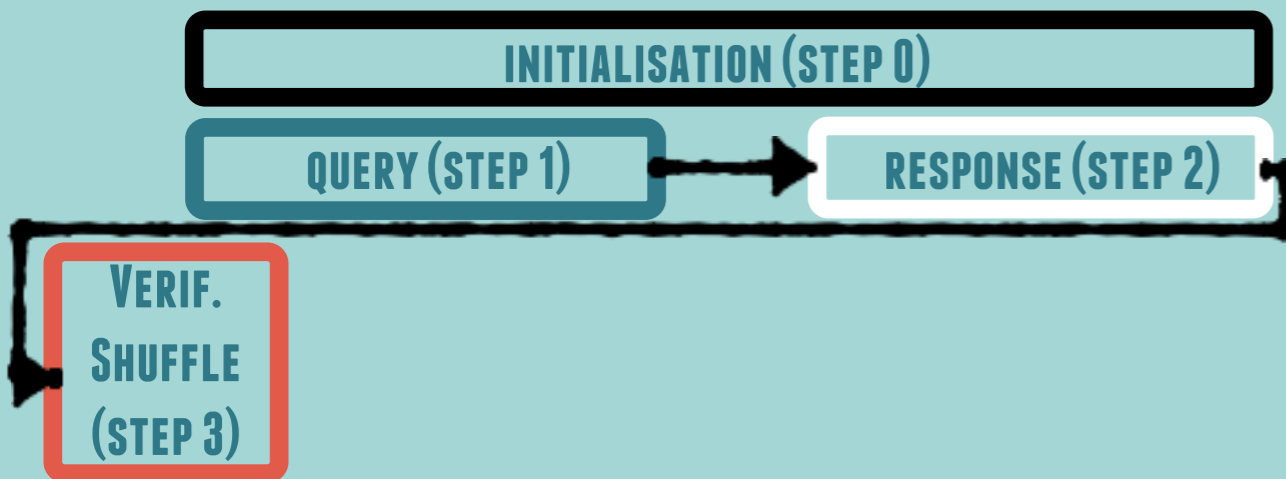
ID	Gender	Age	Ethnicity	flu	Cholesterol_rate	cancer
1	$E_{(1)}$	$E_{(40)}$	$E_{(1)}$	$E_{(1)}$	$E_{(23)}$	$E_{(0)}$
2	$E_{(2)}$	$E_{(40)}$	$E_{(2)}$	$E_{(0)}$	$E_{(34)}$	$E_{(0)}$
...



[[group. attr.], [where. attr.], [aggr. Attr.]]
 [$E_{(1)}$], [$E_{(40)}, E_{(1)}$], [$E_{(23)}, E_{(0)}$]
 ...



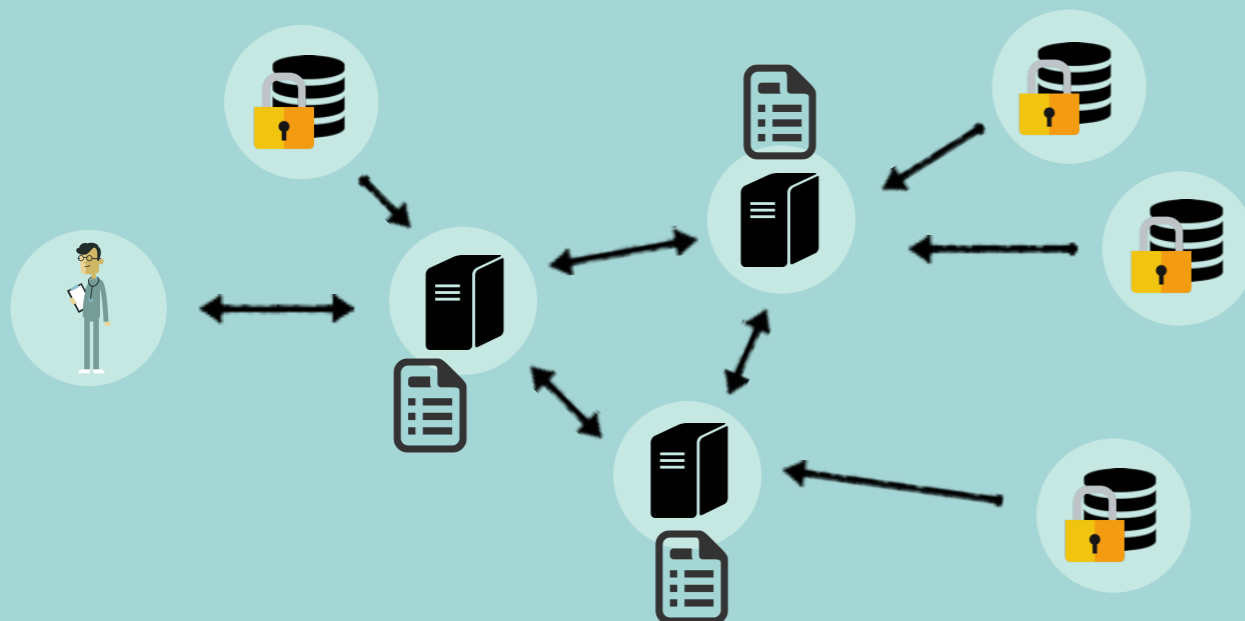
WORKFLOW - VERIF. SHUFFLE (STEP 3)



Each server starts a **verifiable shuffle protocol**:

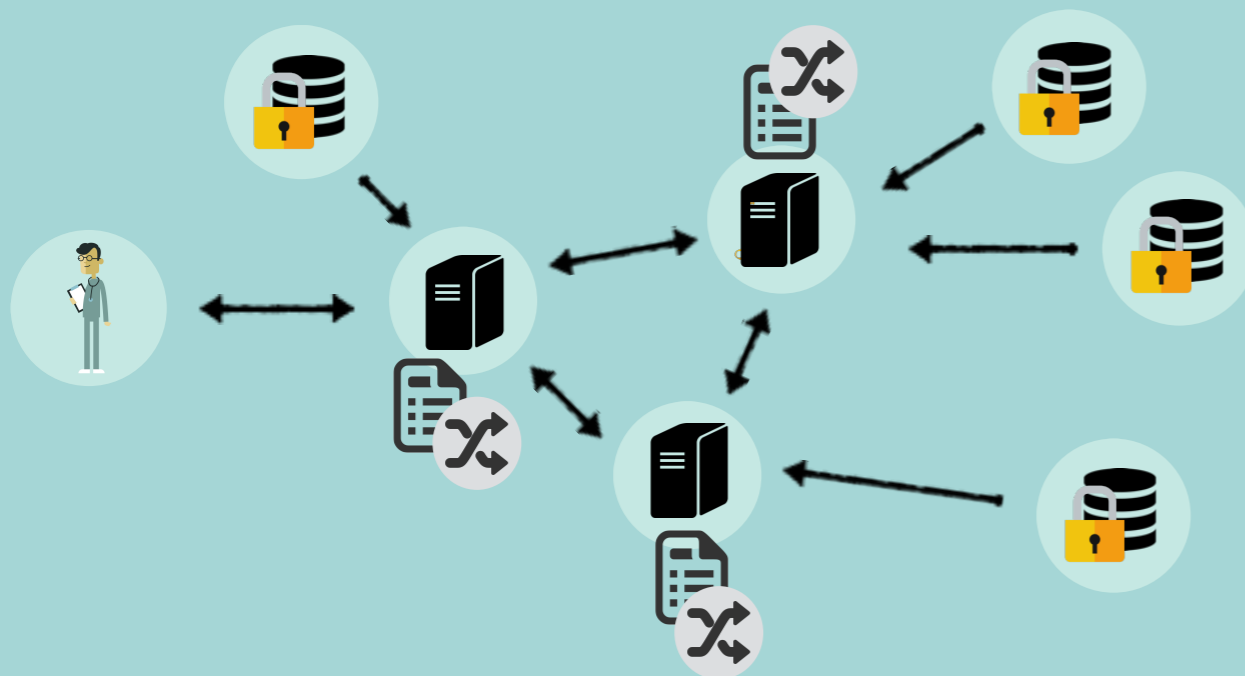
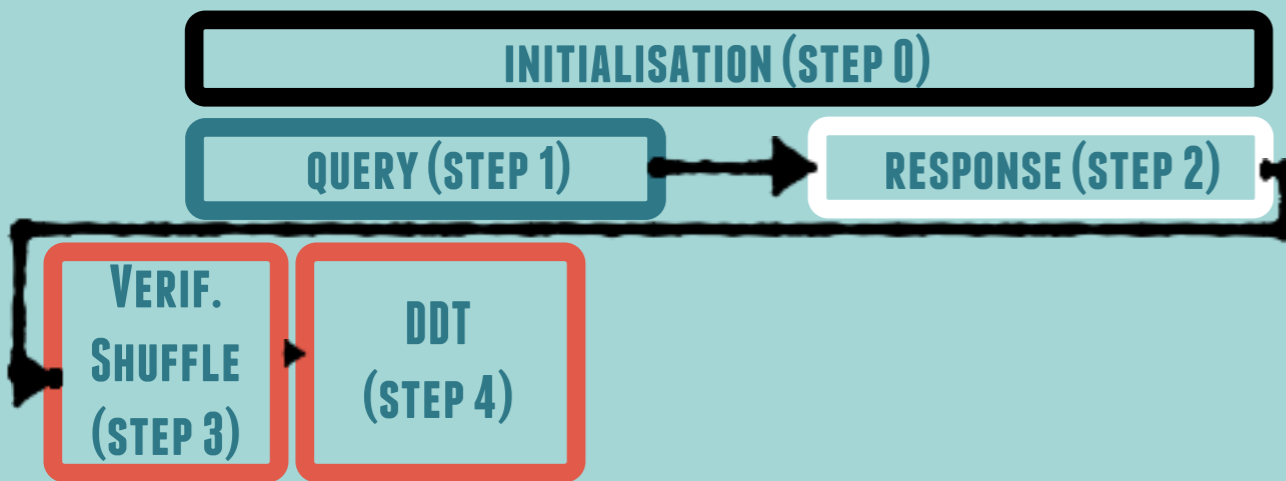
In this protocol each server sequentially:

- **Shuffle** the list of responses
- **Rerandomize** (re-encryption) all the ciphertexts



Using Neff Shuffle and the corresponding zero-knowledge proof [1]

WORKFLOW - DDT (STEP 4)



Each server starts a **distributed deterministic tagging protocol**:

Query:

WHERE $age = E_k(40)$ AND $ethnicity = E_k(2)$



WHERE $age = DT(40)$ AND $ethnicity = DT(2)$

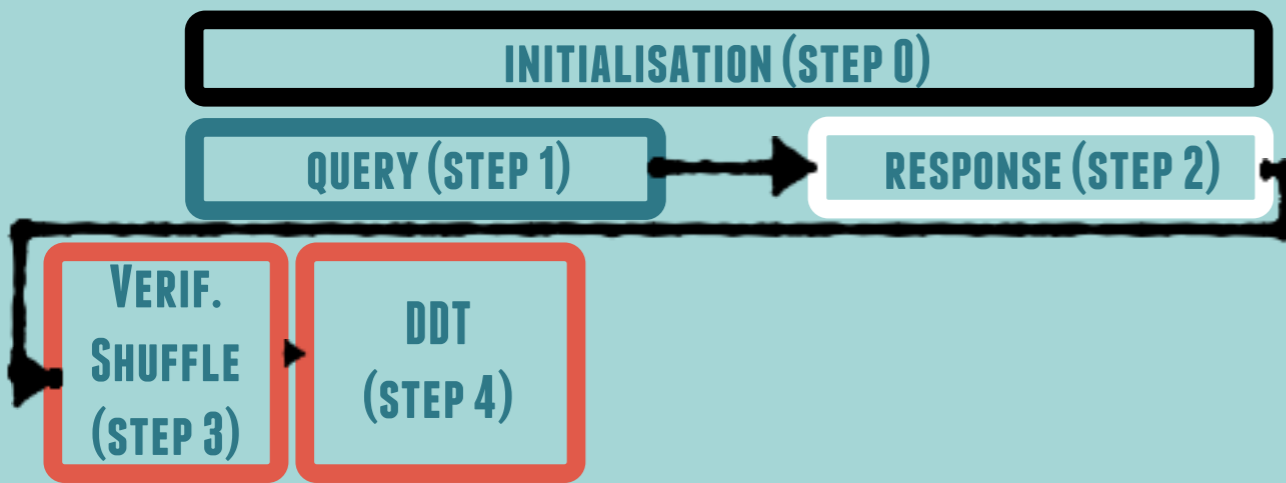
Data:

$[[E_k(1)], [E_k(40), E_k(2)], [E_k(23), E_k(1)]]$



$[[DT(1)], [DT(40), DT(2)], [E_k(23), E_k(1)]]$

WORKFLOW - DDT (STEP 4)



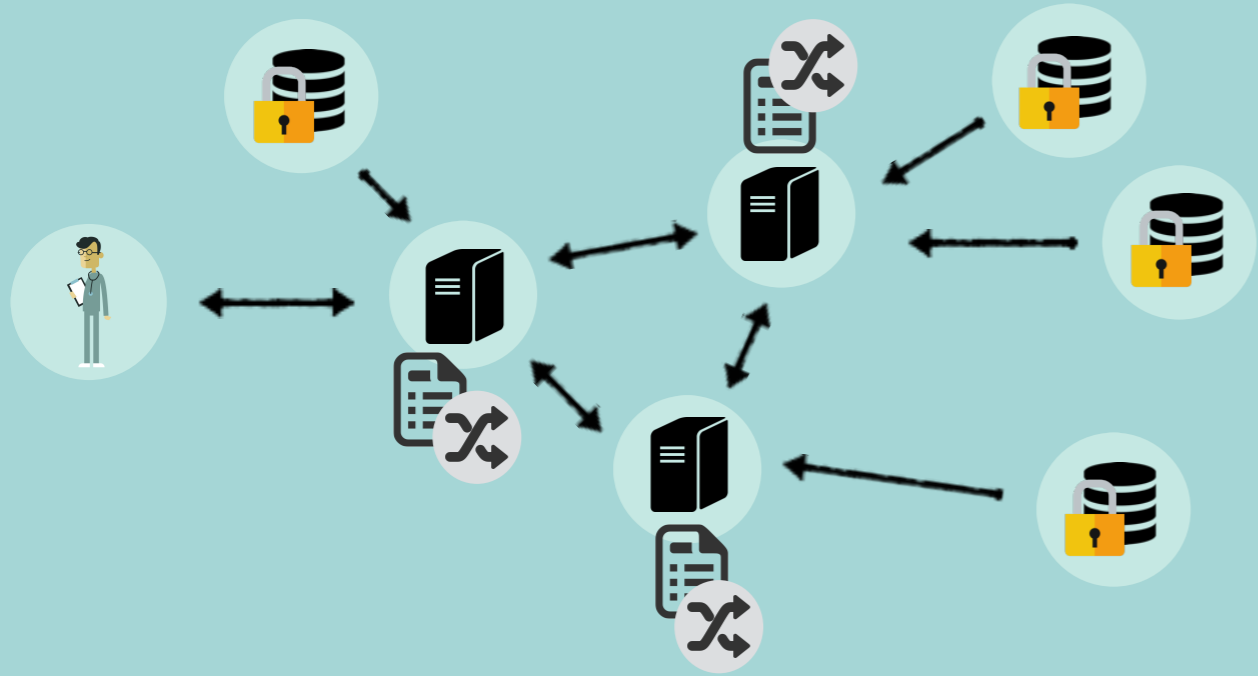
Each server starts a **distributed deterministic tagging protocol**:

In this protocol each server sequentially:

- **partially decrypt** the ciphertexts
- **Blinds** the message by multiplying the ciphertexts with a random **ephemeral secret key**

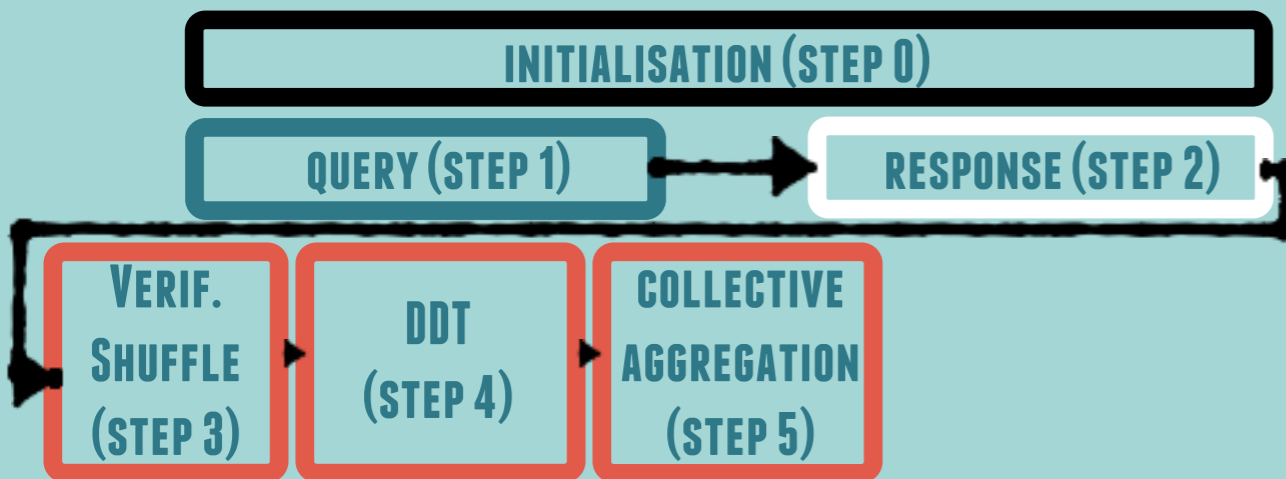
→ deterministic tag depending on the value of the encrypted message

All operations are done with zero-knowledge proofs from Camenisch et al.

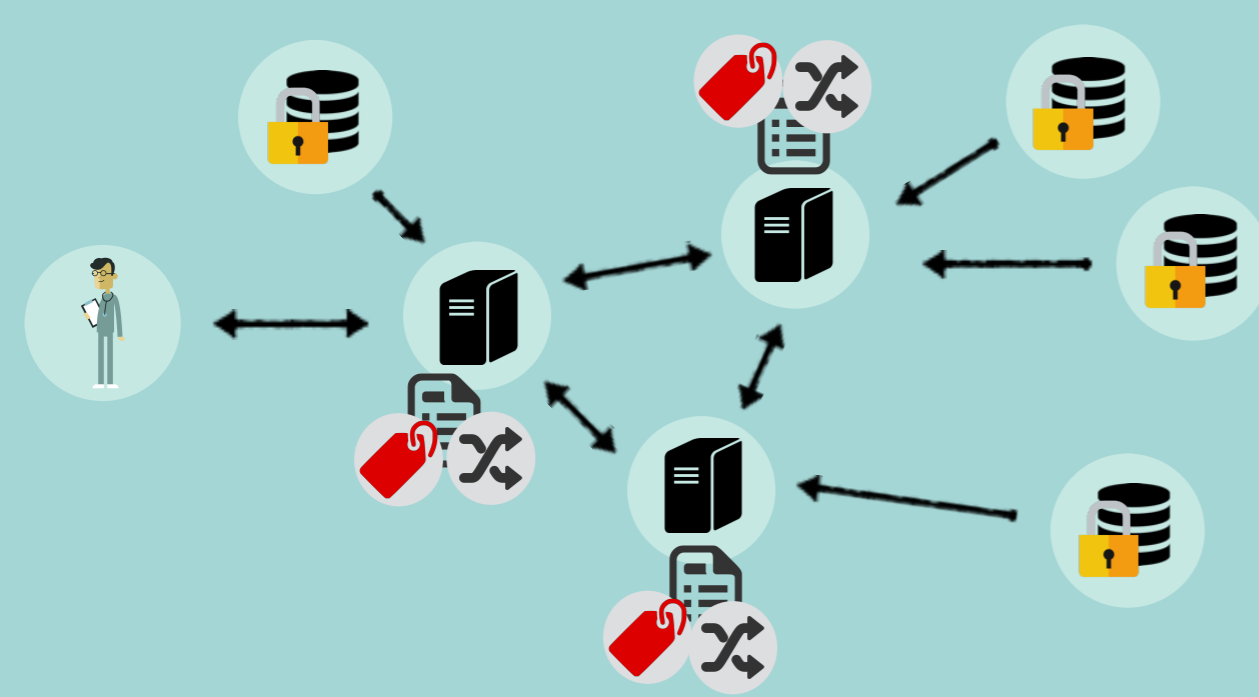


[1] Jan Camenish and Markus Stadler. Proof systems for general statements about discrete logarithms. (1997)

WORKFLOW - COLLECTIVE AGGR. (STEP 5)

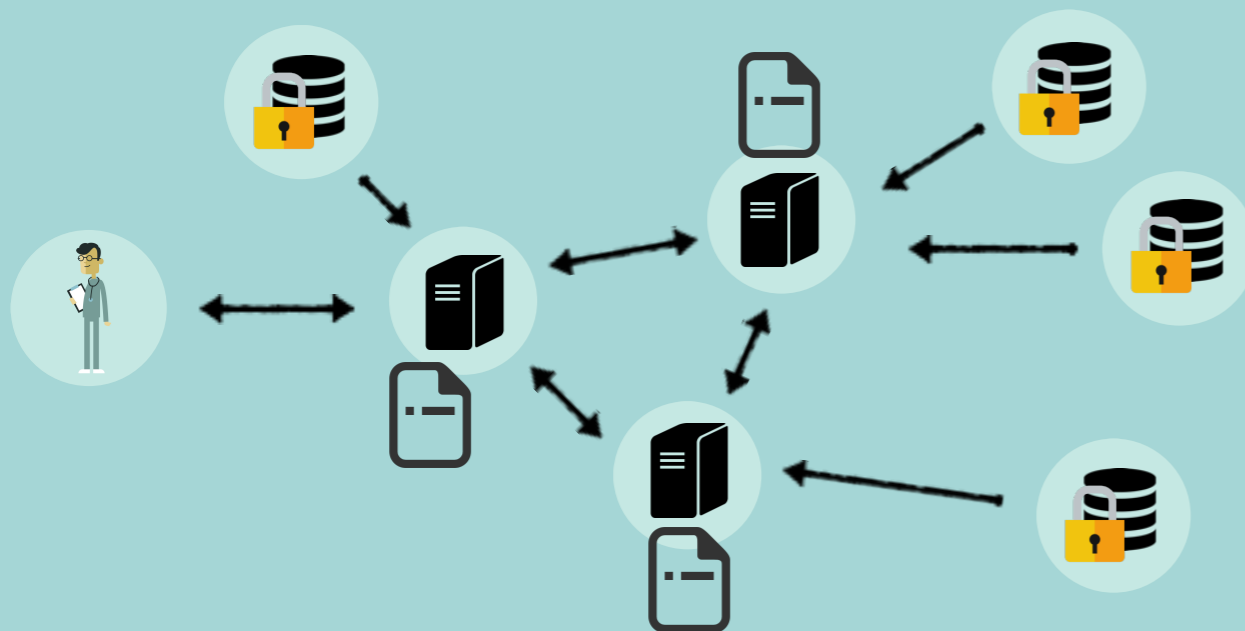
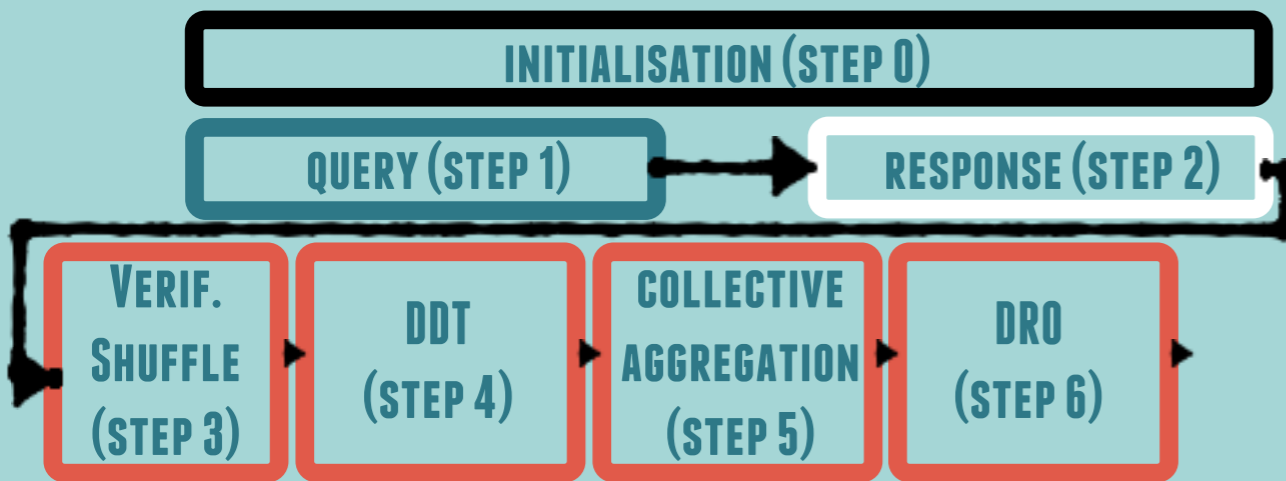


Servers **collectively aggregate** the responses by group.

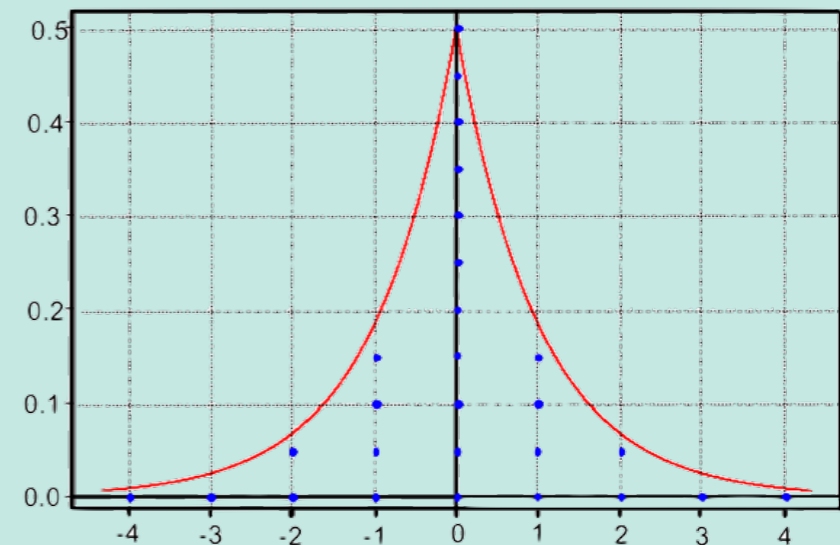


Proofs consist in publishing the ciphertexts and the result

WORKFLOW - DRO (STEP 6)

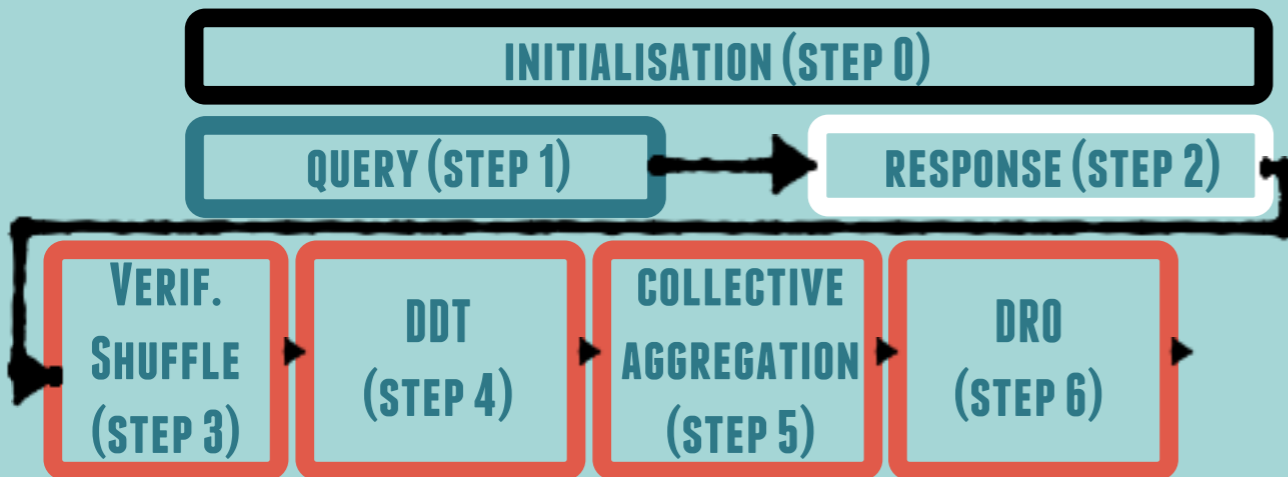


Distributed Results Obfuscation:
Setup:
Servers agree on (ϵ, δ) -differential privacy parameters and produce:



→ $[0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, \dots]$ = list of noise values satisfying (ϵ, δ) -differential privacy.

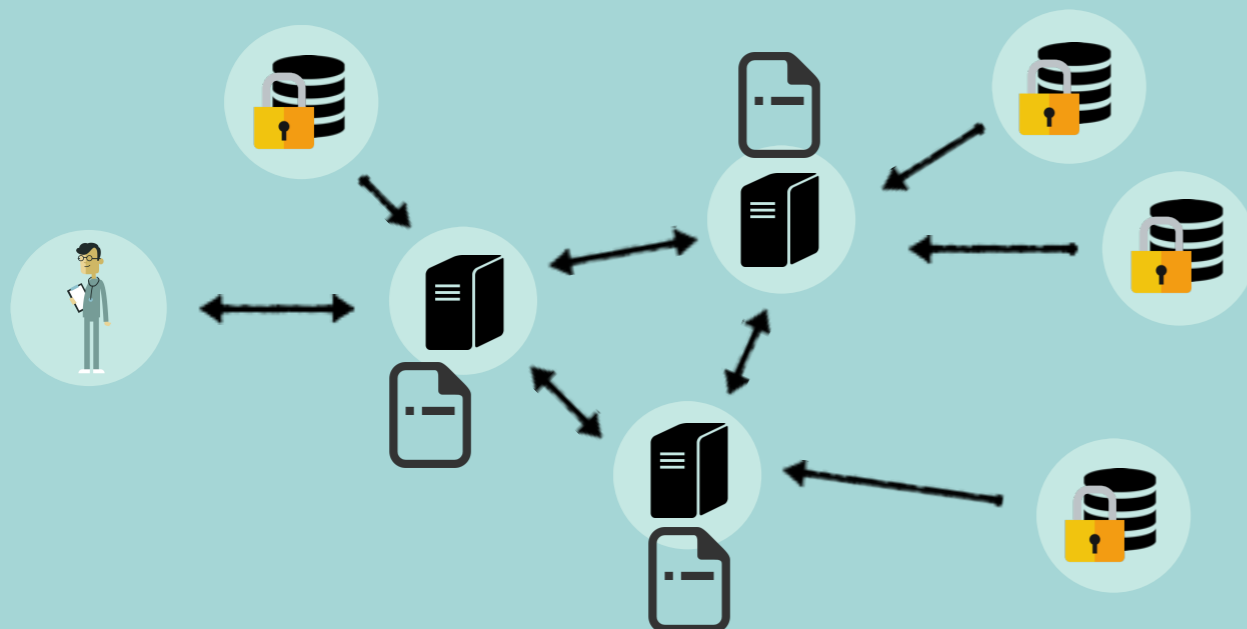
WORKFLOW - DRO (STEP 6)



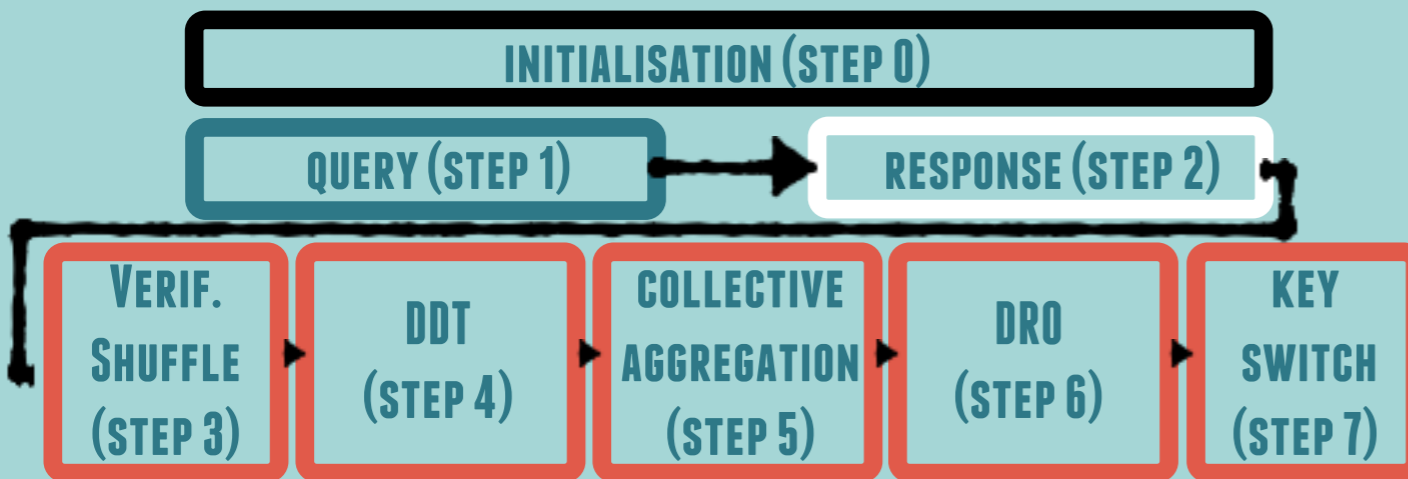
Distributed Results Obfuscation: Runtime:

- A server starts a **collective shuffling** of the list of noise values
- **adds the first noise value** in the list to the query result.

→ Oblivious noise addition (shuffling **encrypts and shuffles** the list of noise values).



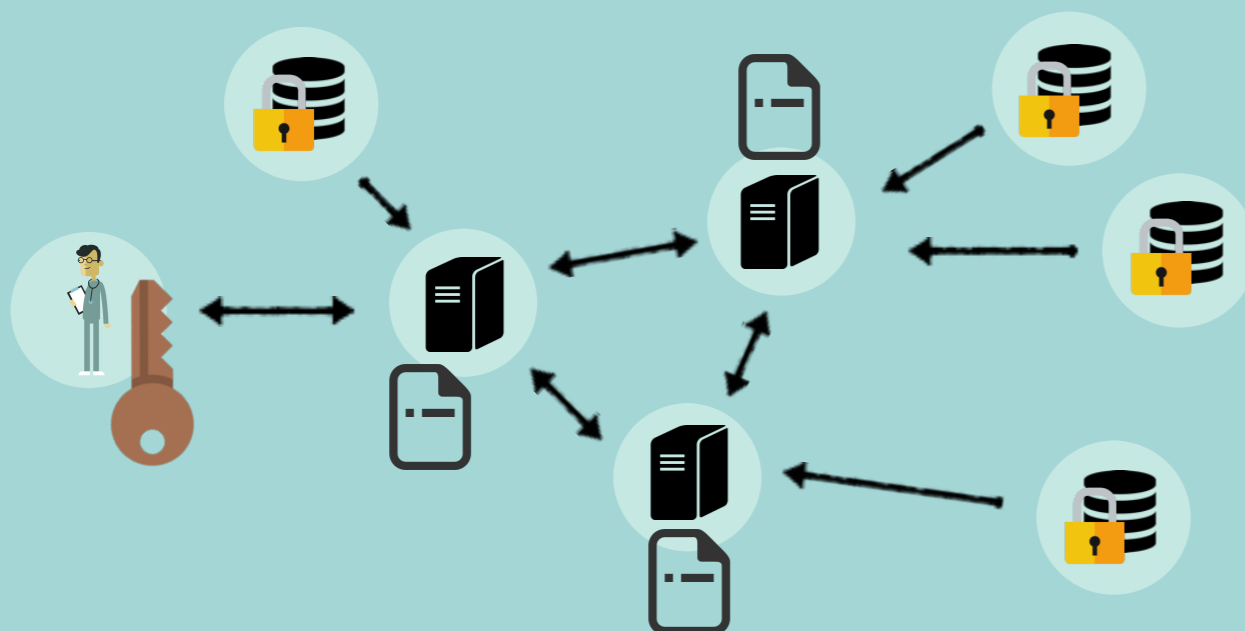
WORKFLOW - KEY SWITCH (STEP 7)



In the **key switch protocol** each server:

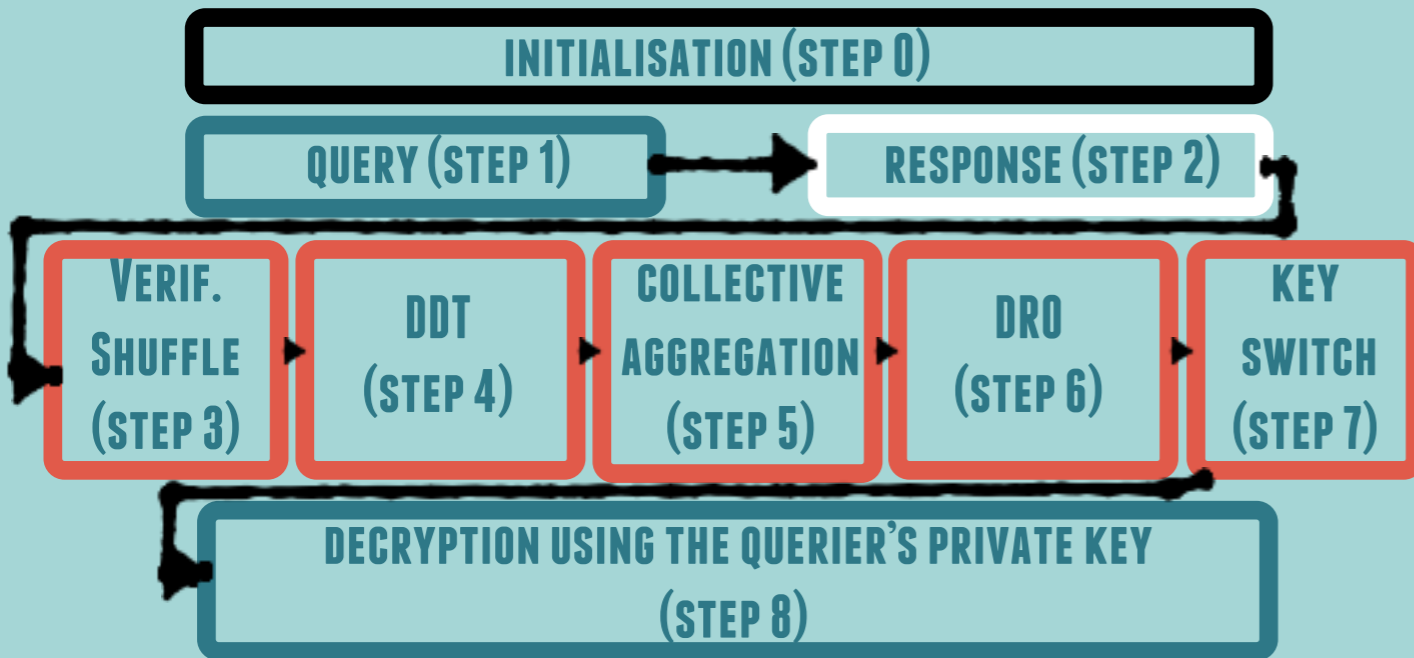
- **partially decrypt**
- **encrypt** with a new key all the ciphertexts.

→ **Encryption is switched** from the Collective Key to the querier's public key.

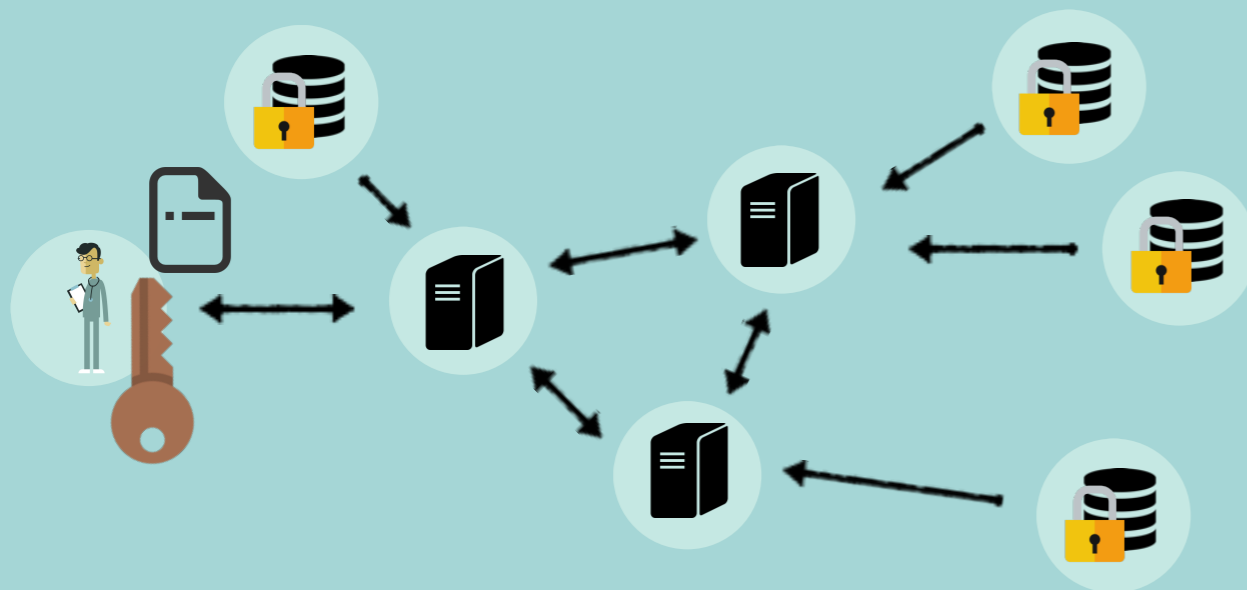


All operations are done with zero-knowledge proofs from Camenish et al.

WORKFLOW - DECRYPTION (STEP 8)



Querier **decrypts** the result with his secret key



PERFORMANCE EVALUATION

Servers configuration

- Memory: 256GB RAM
- Processor: Intel Xeon E5-2680 v3 (Haswell)
- **Cores: 24 (with 48 threads)**
- Frequency: 2.5GHz
- Bandwidth capacity: 1Gbps

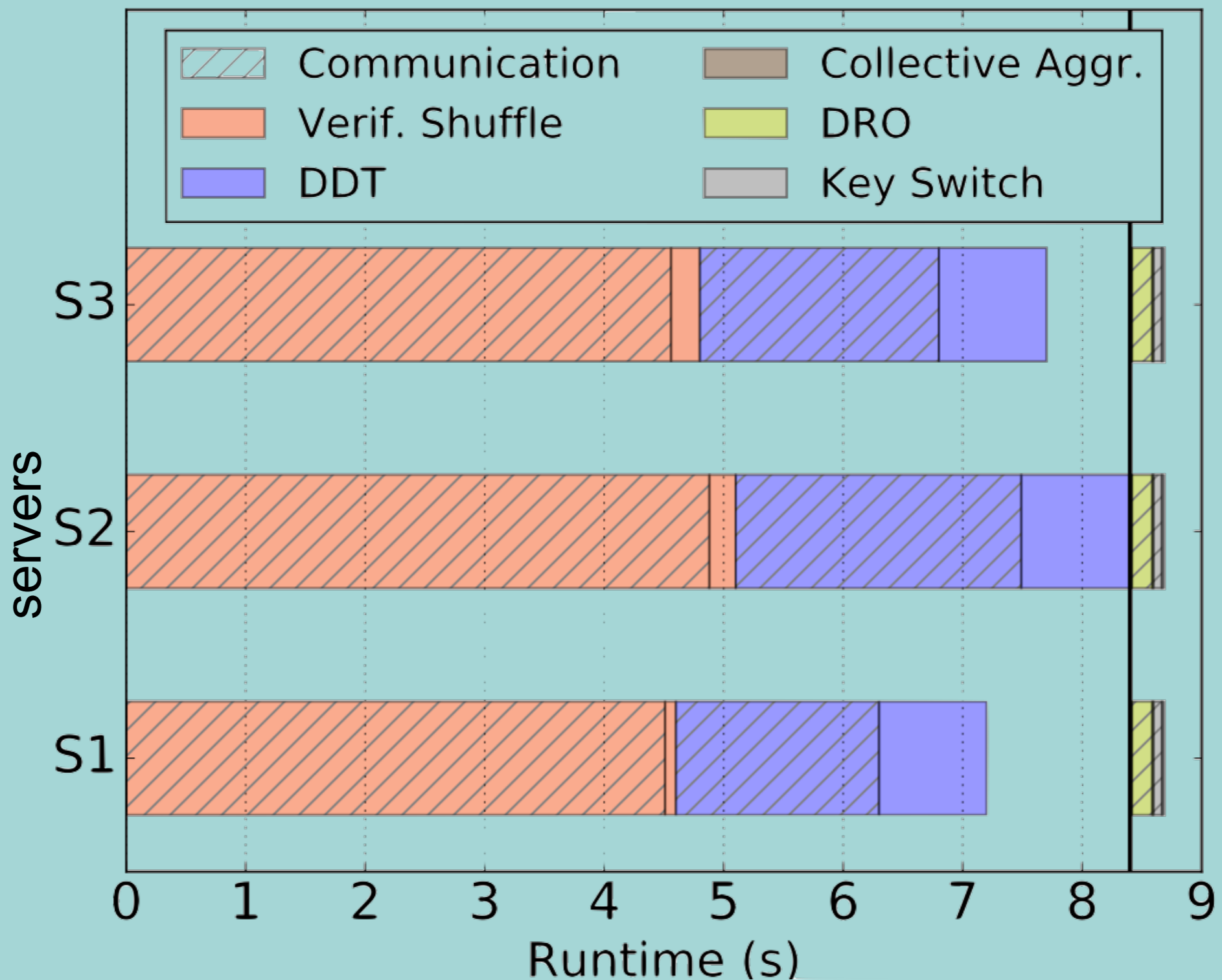
Network and Crypto

- Realistic virtual network emulation tool with 10ms delays btw. servers
- DeDiS' Onet library
- DeDiS' implementation of Ed25519 Elliptic Curve (**128-bit security**)

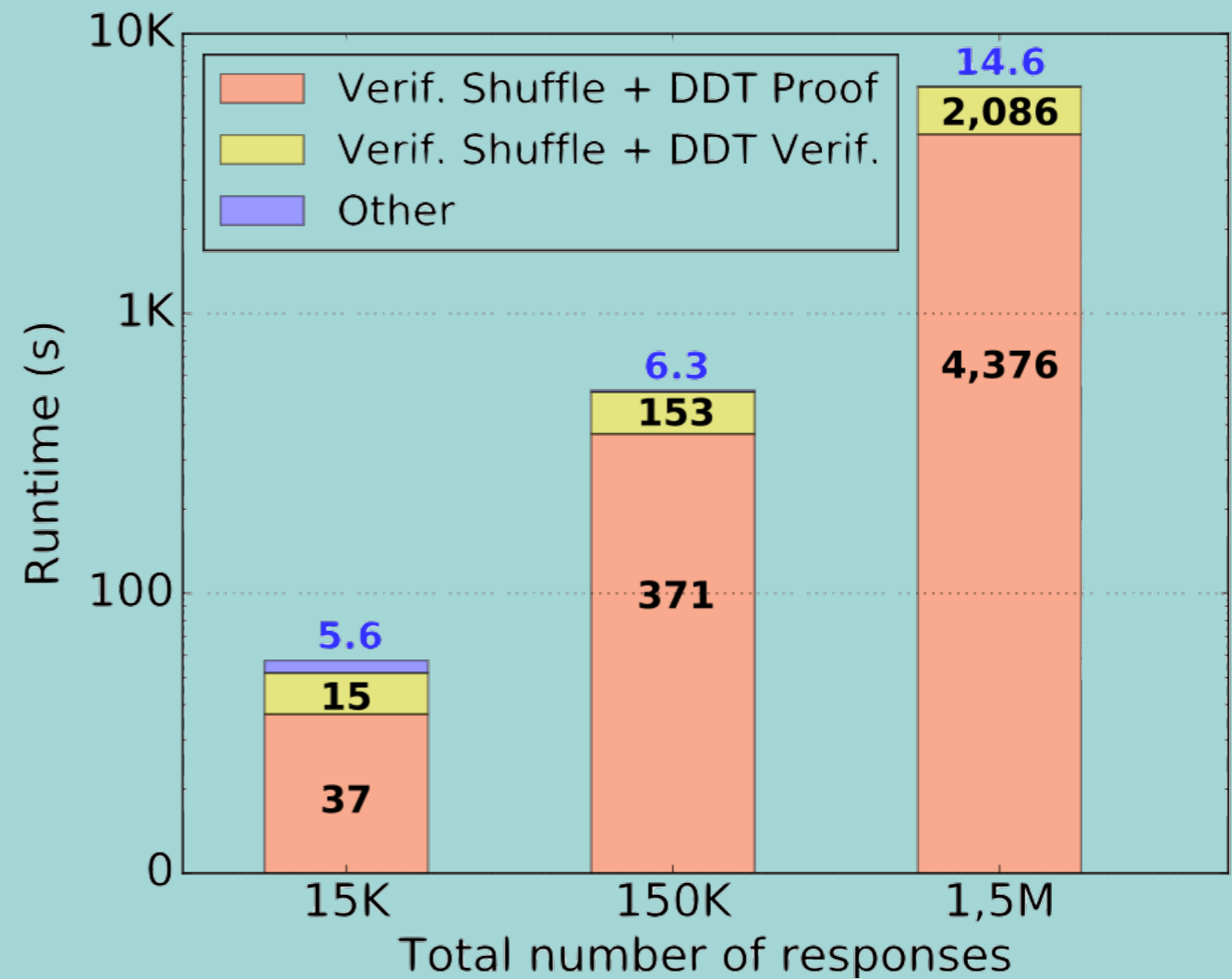
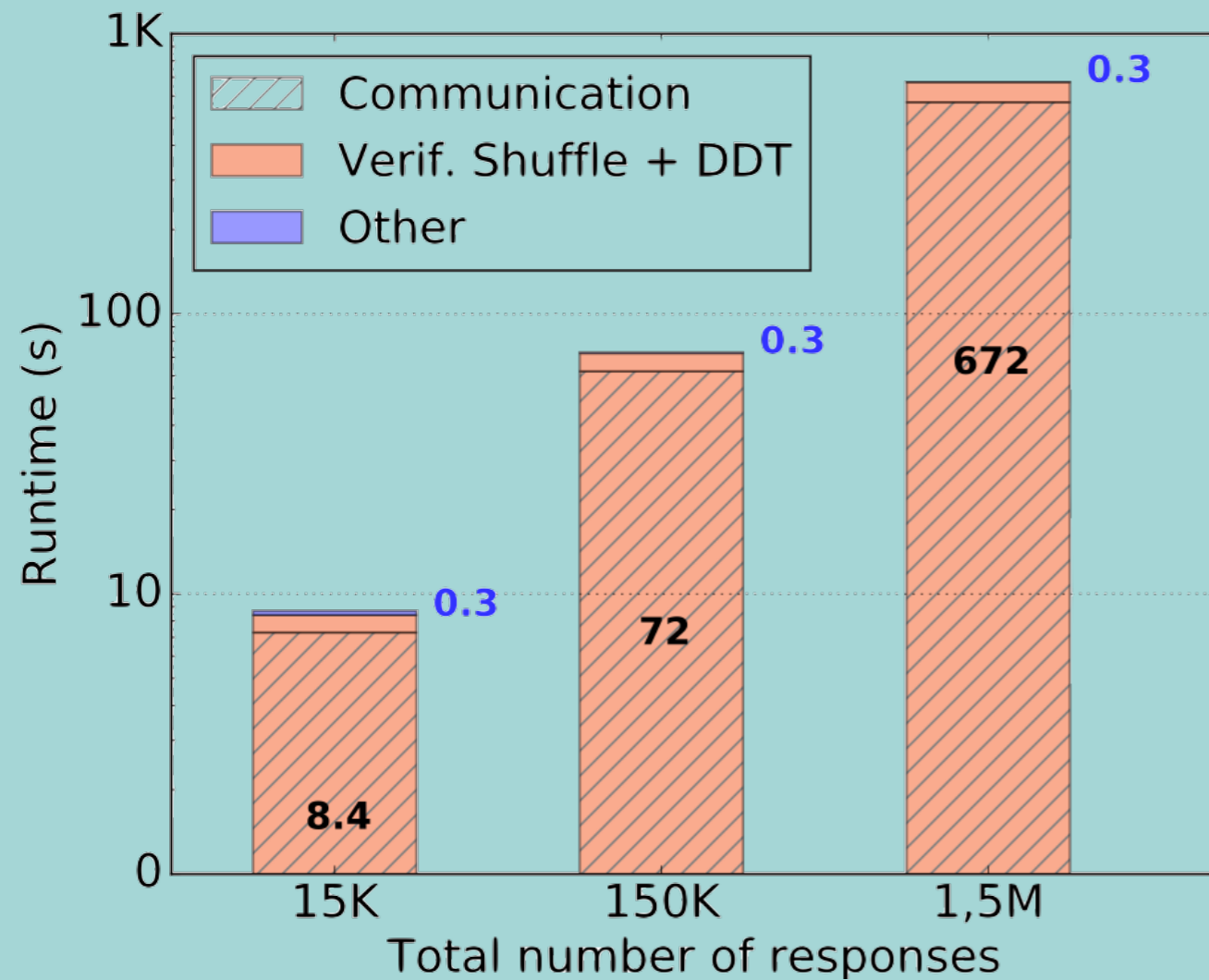
Default parameters

- **3 servers**
- **15,000 responses** in total (equally distributed in servers)
- **1 GROUP BY** attribute with 10 possible values , **1 WHERE** and **10 aggregating attributes**
- 1000 noise values

SERVICES COLLABORATION



RUNTIME VS. NBR. OF RESPONSES

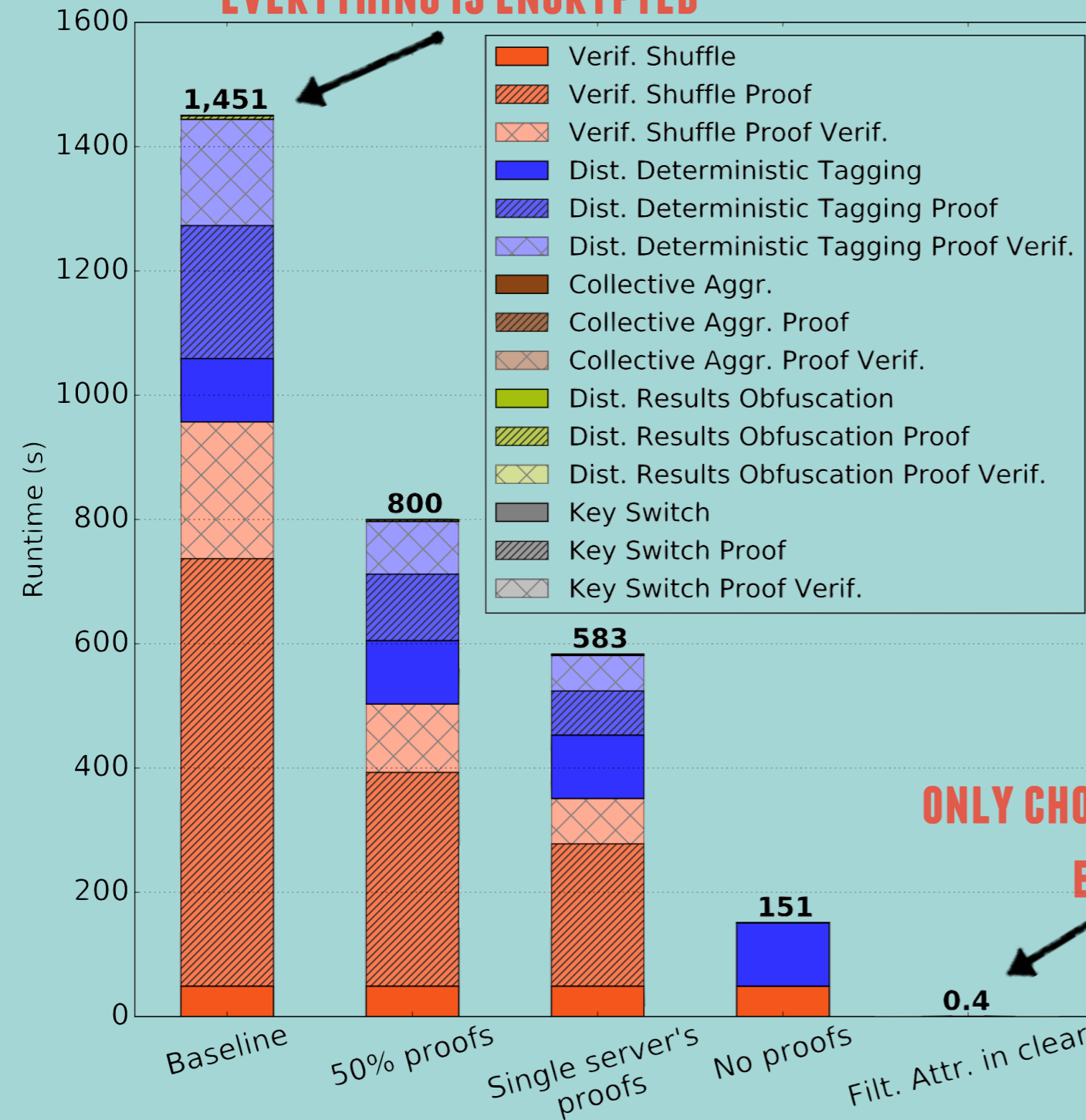


PERFORMANCE/SECURITY TRADEOFFS

EVERYTHING IS ENCRYPTED

```
SELECT SUM (CHOLESTEROL_RATE), COUNT(*)
FROM DP1,...,DP20
WHERE AGE IN [40:50] AND ETHNICITY = CAUCASIAN
GROUP BY GENDER
```

3 servers
 400K responses with
 1 GROUP BY attribute
 2 WHERE attributes
 2 aggregating attributes



ONLY CHOLESTEROL_RATE IS ENCRYPTED

CONCLUSION

A Decentralized System for Privacy-Conscious Data Sharing

- SQL statistical queries based on Boolean conditions
 - Strongest-link security
 - Data confidentiality
 - Distributed differential privacy
 - Distributed deterministic tagging of probabilistic ciphertexts
 - Collective encryption key switching
- Runtime linear with the amount of data to process

github.com/lca1/unlynx

david.froelicher@epfl.ch

