

# Practical Differential Privacy

Georgios Kellaris

Ph.D. Dissertation at HKUST



Privacy Tools  
for Sharing Research Data

A National Science Foundation  
Secure and Trustworthy Cyberspace Project



with additional support from the Sloan Foundation and Google, Inc.

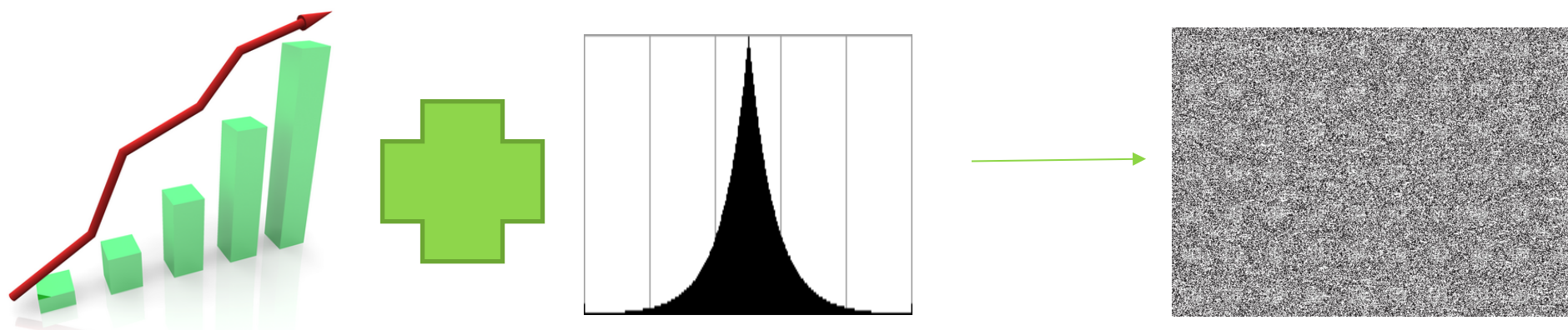
## $\epsilon$ -Differential Privacy [DMNS06]

- Ensure user privacy when publishing statistics, by hiding the presence of any user in the data
- How?
  - Add noise to the statistics before publishing
  - Noise is drawn from the Laplace Distribution
  - The noise scale must be proportional to the impact any user has on the statistics (sensitivity)



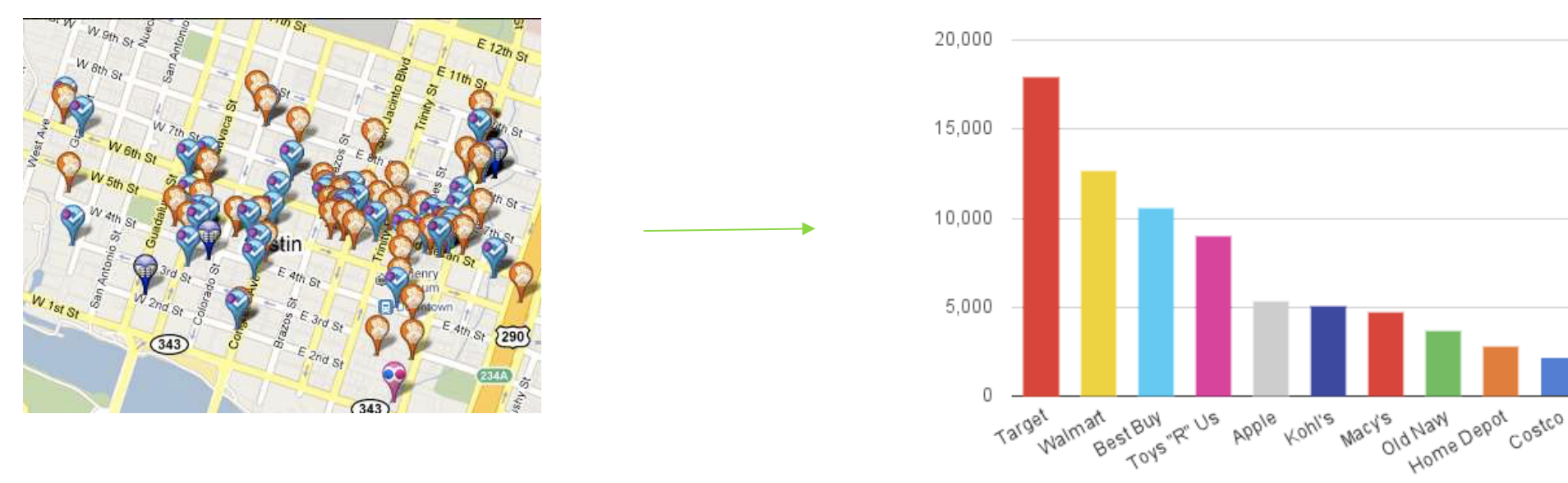
## Utility Issues

- The required noise may destroy the accuracy of the output
- For critical applications this might be unacceptable
  - E.g. for detecting disease outbreaks

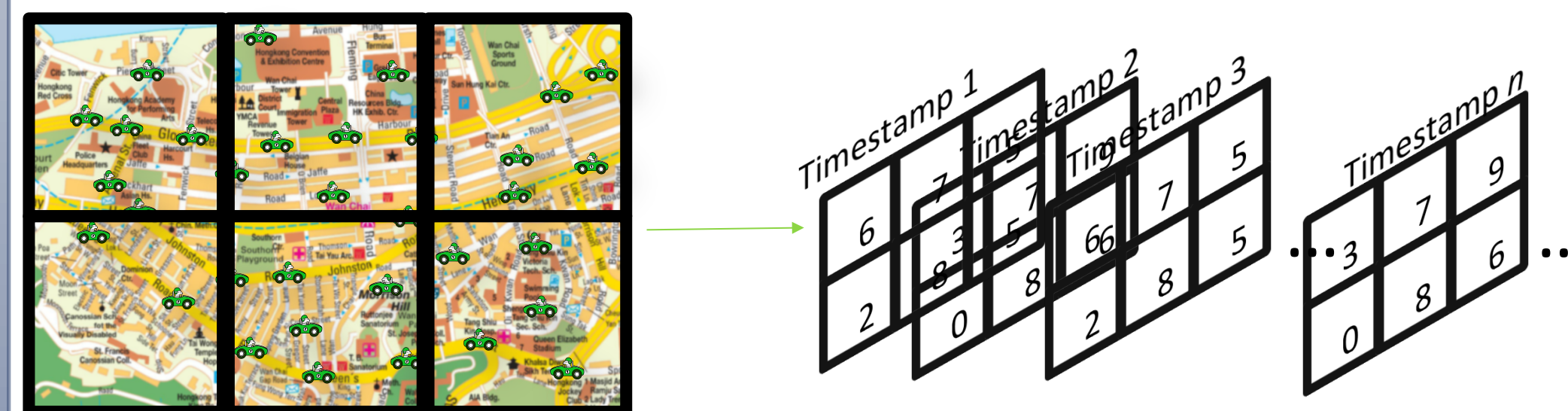


## Targeted Settings and Challenges

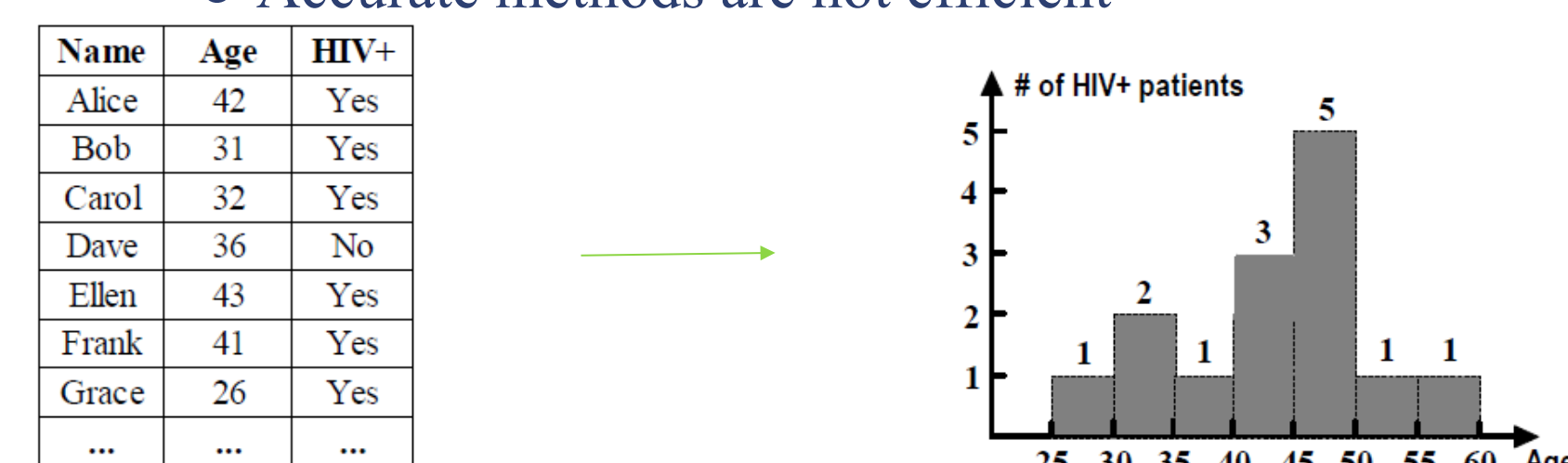
- Historical Data
  - High Sensitivity



- Infinite Streams
  - Event-level privacy [DNPR10] not enough

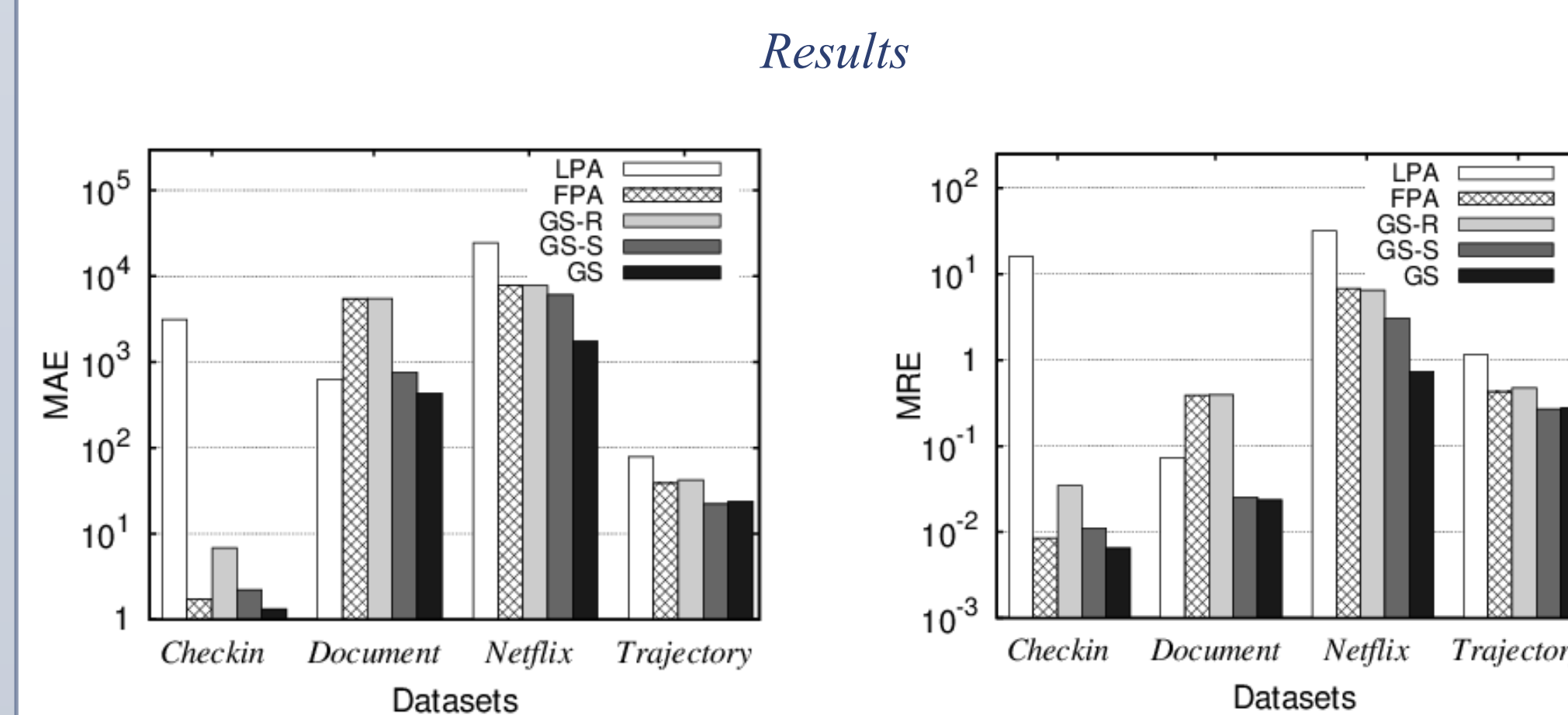
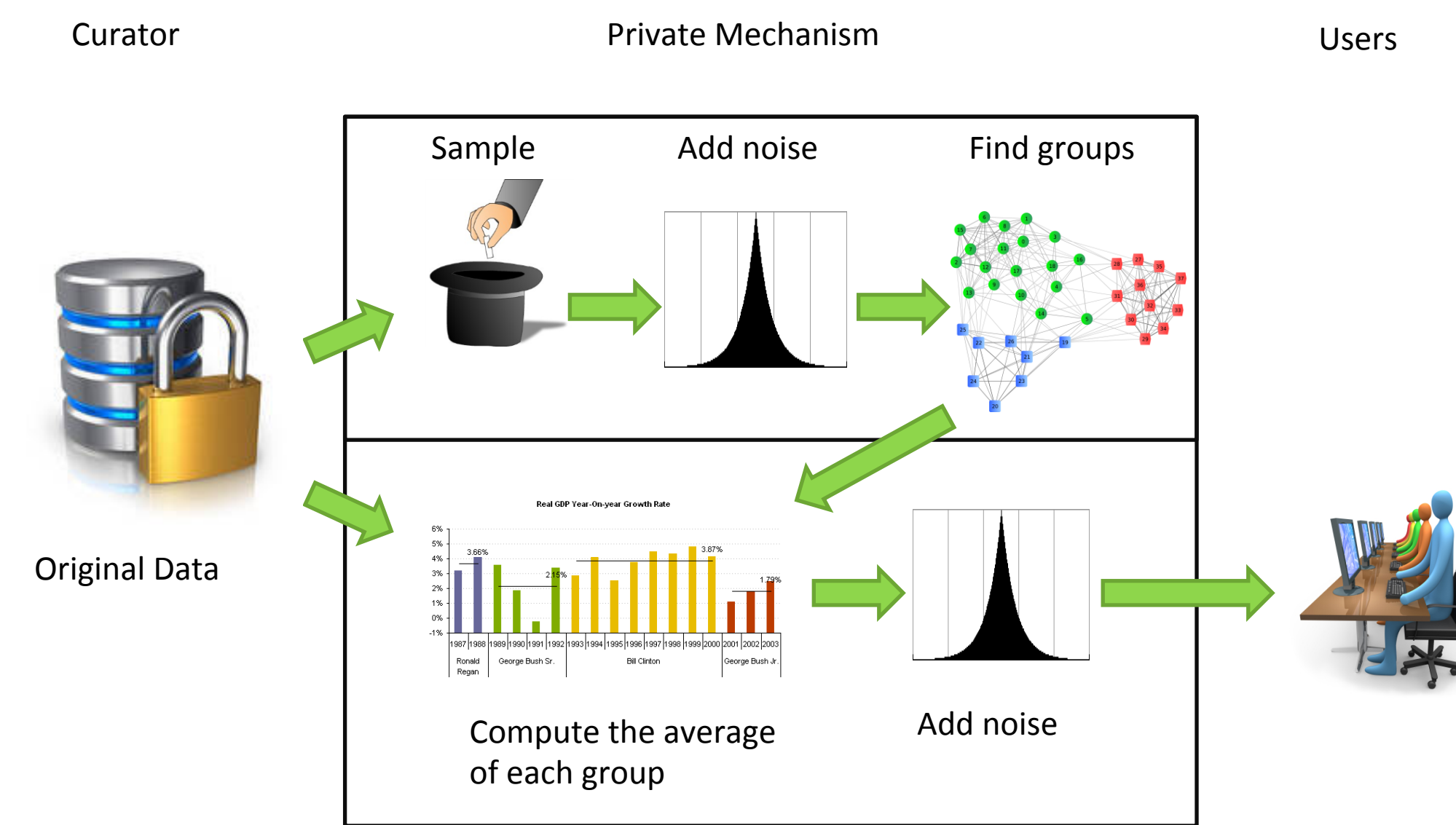


- Histograms for Range Queries
  - Noise error accumulates when computing ranges
  - Accurate methods are not efficient

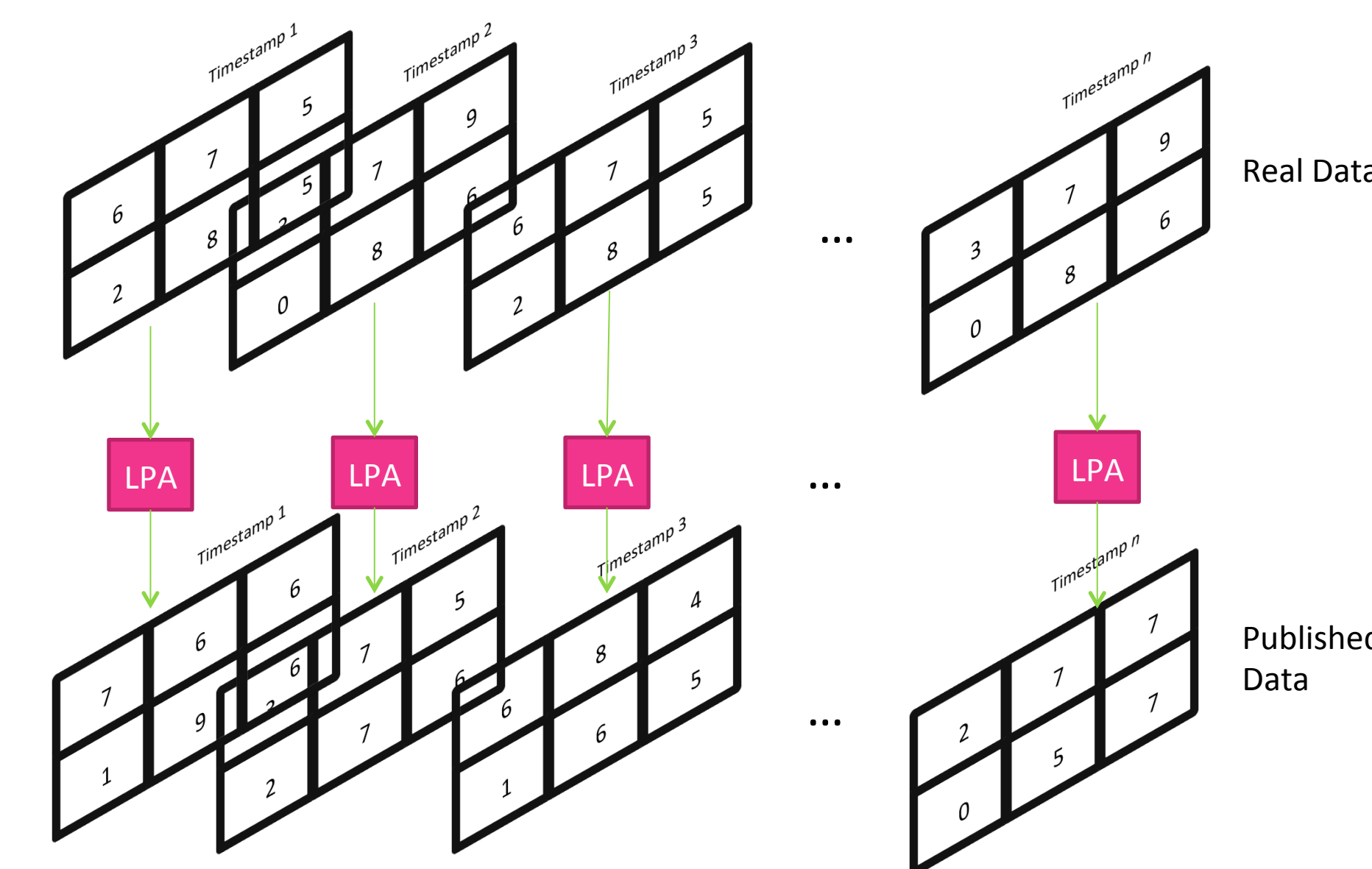


## Solutions

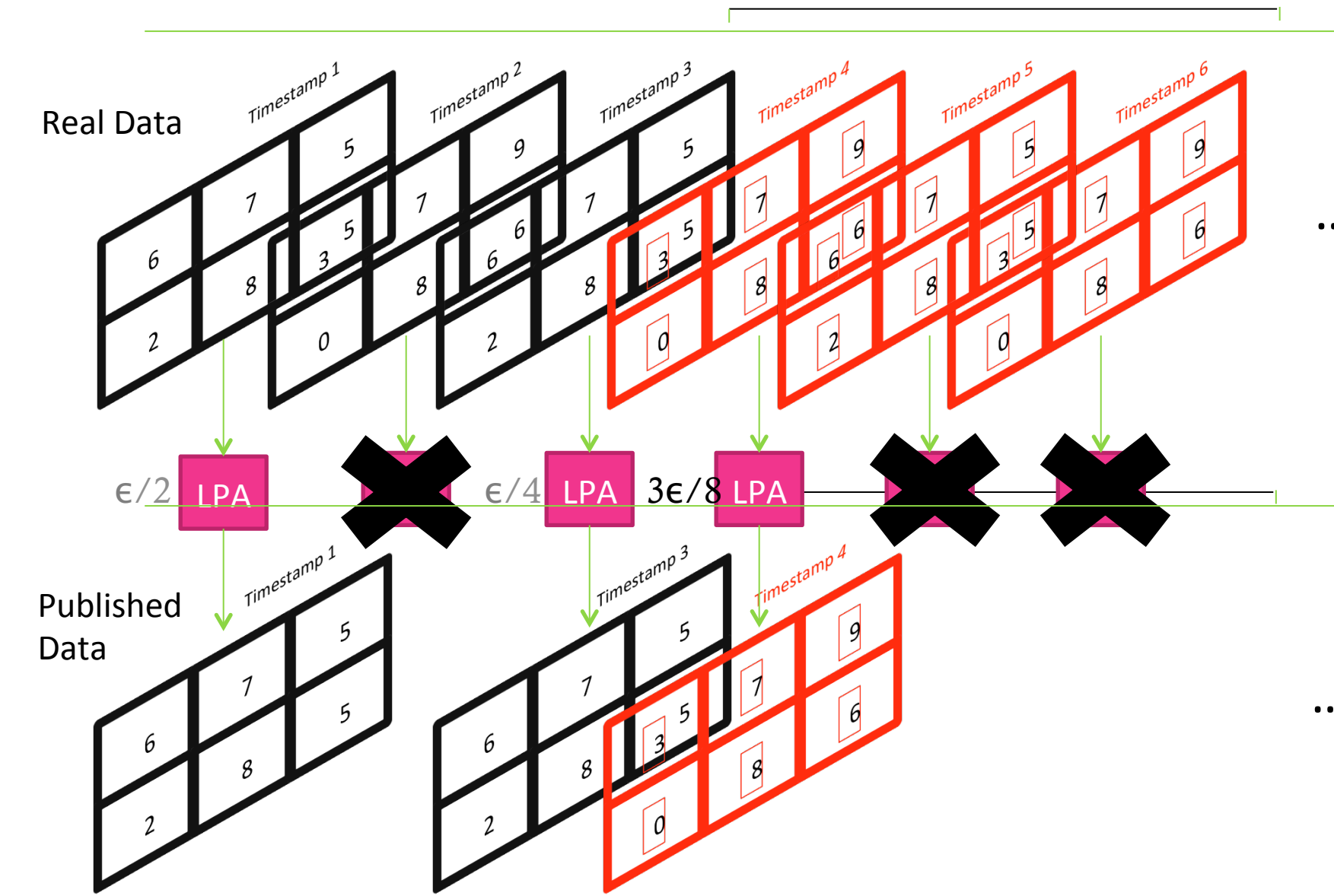
### Historical Data [KP13]



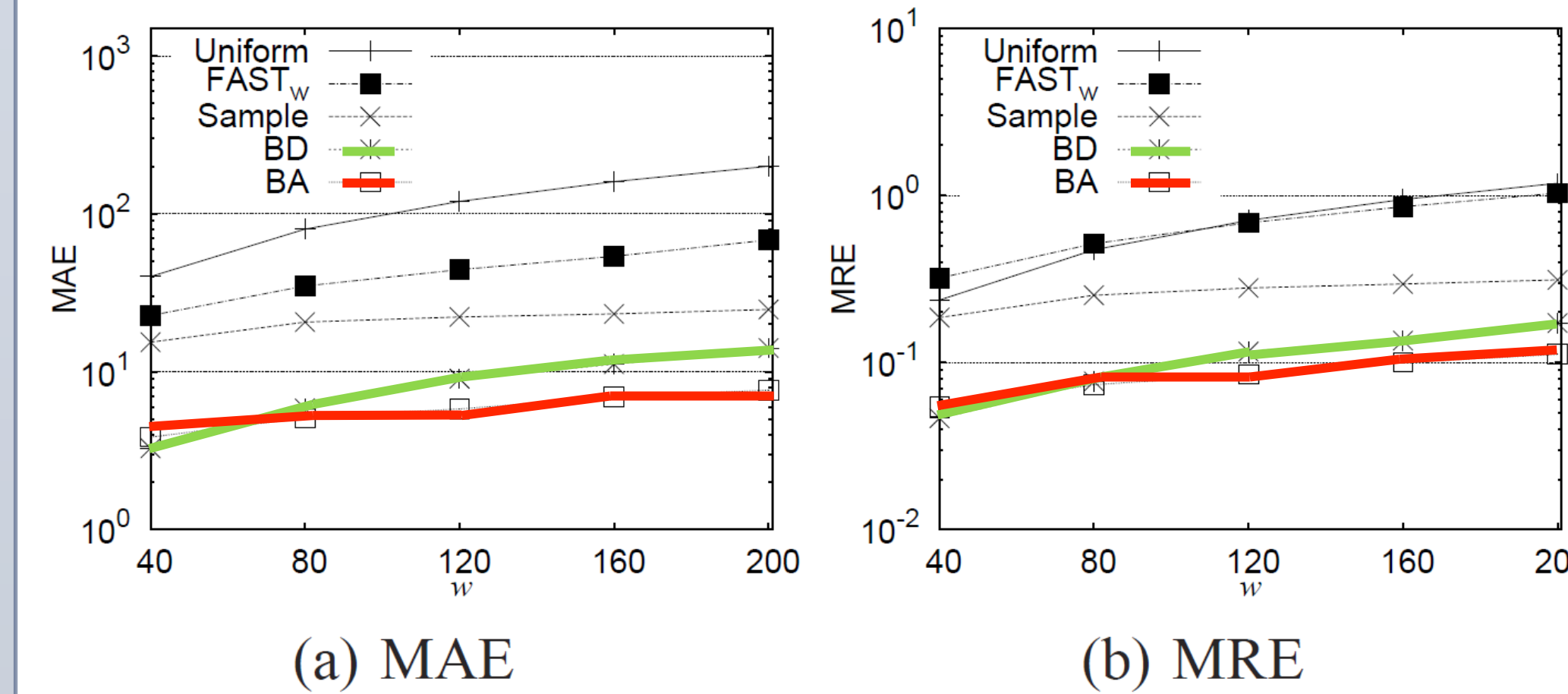
### Infinite Streams [KXP14]



- Event level
  - $\epsilon$ -differential privacy at any timestamp
  - it does not protect user movement
  - noise proportional to 1
- $w$ -Event level
  - $\epsilon$ -differential privacy at any  $w$  consecutive timestamps
  - it protects user movement that lasts at most  $w$  timestamps
  - noise proportional to  $w \ll T$
- User level
  - $\epsilon$ -differential privacy at all timestamps
  - it protects user movement
  - noise proportional to  $T$

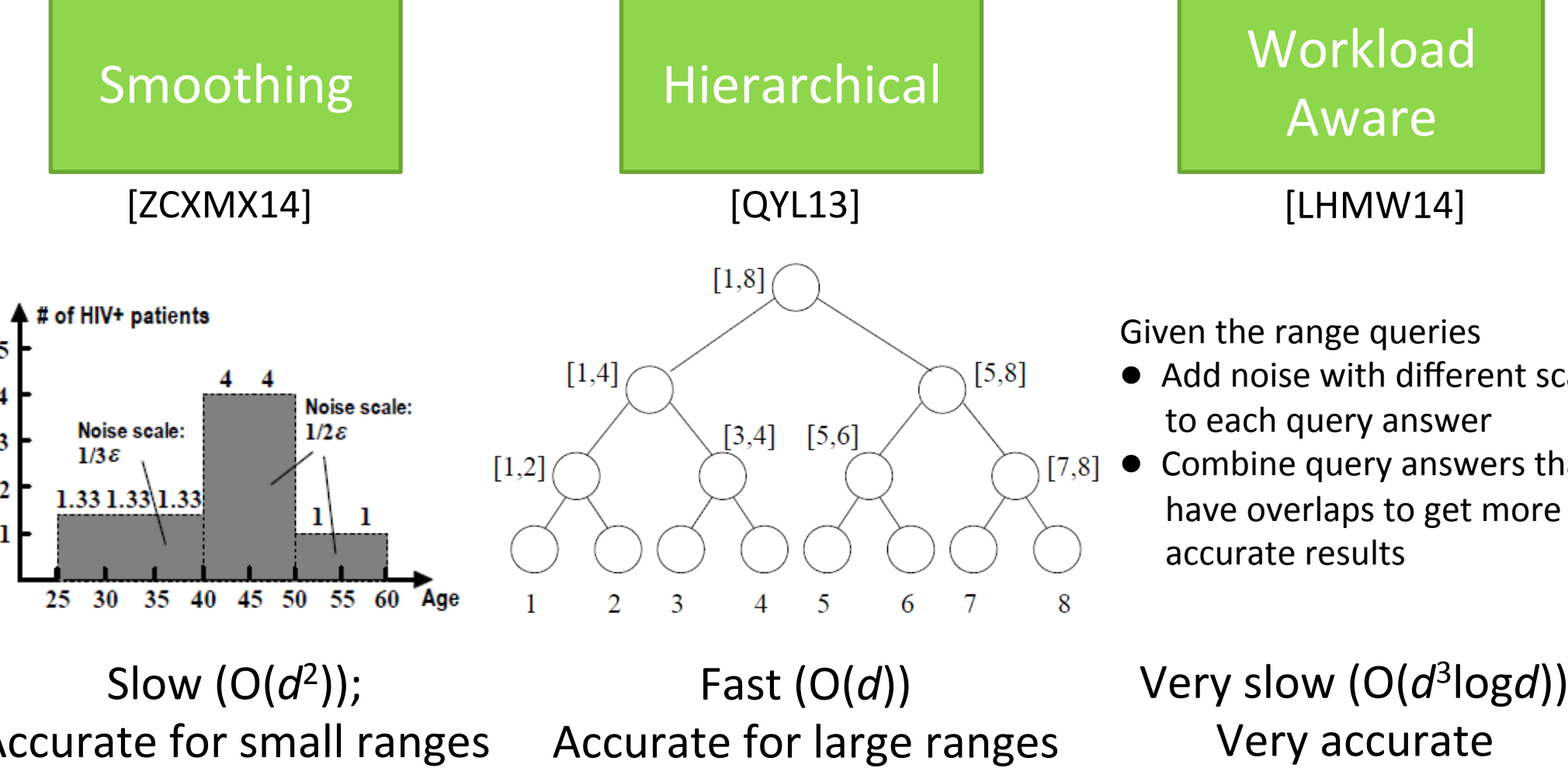


## Results

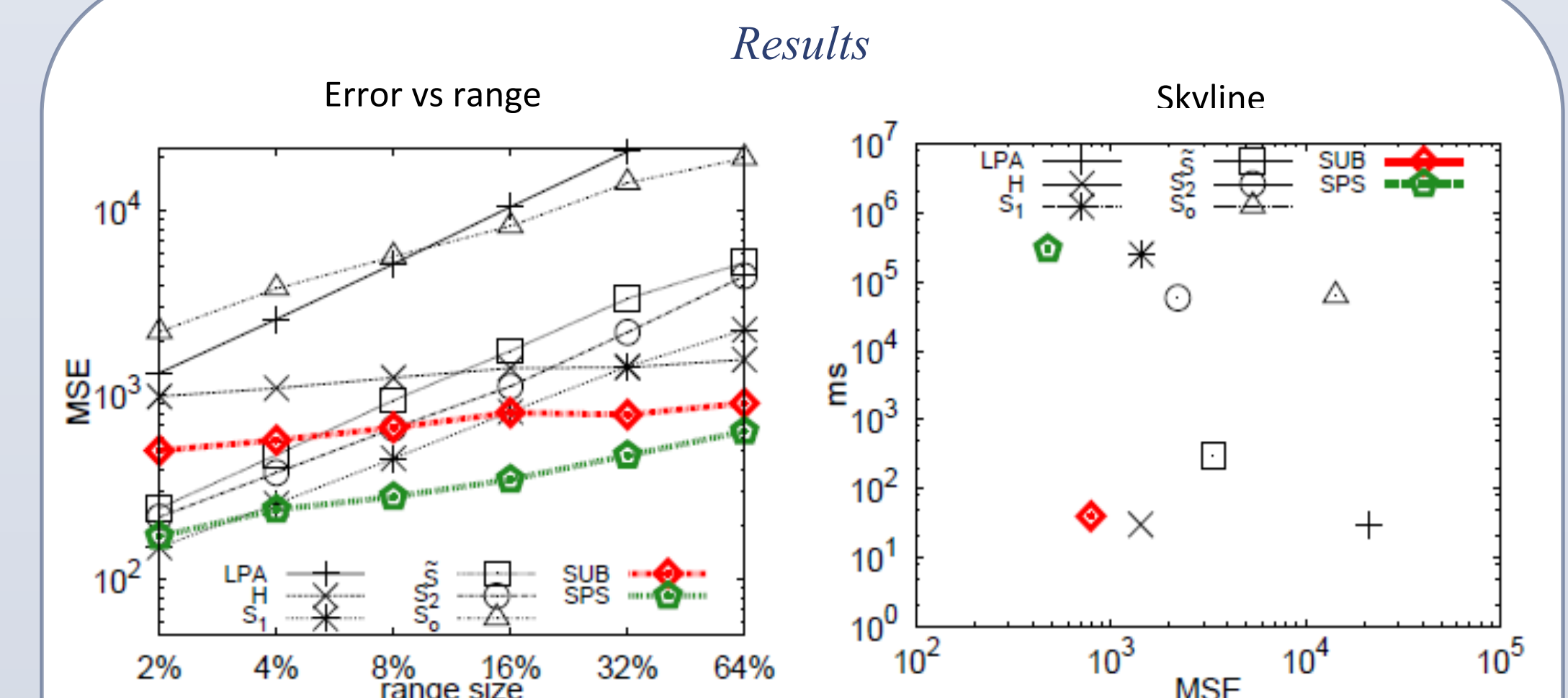
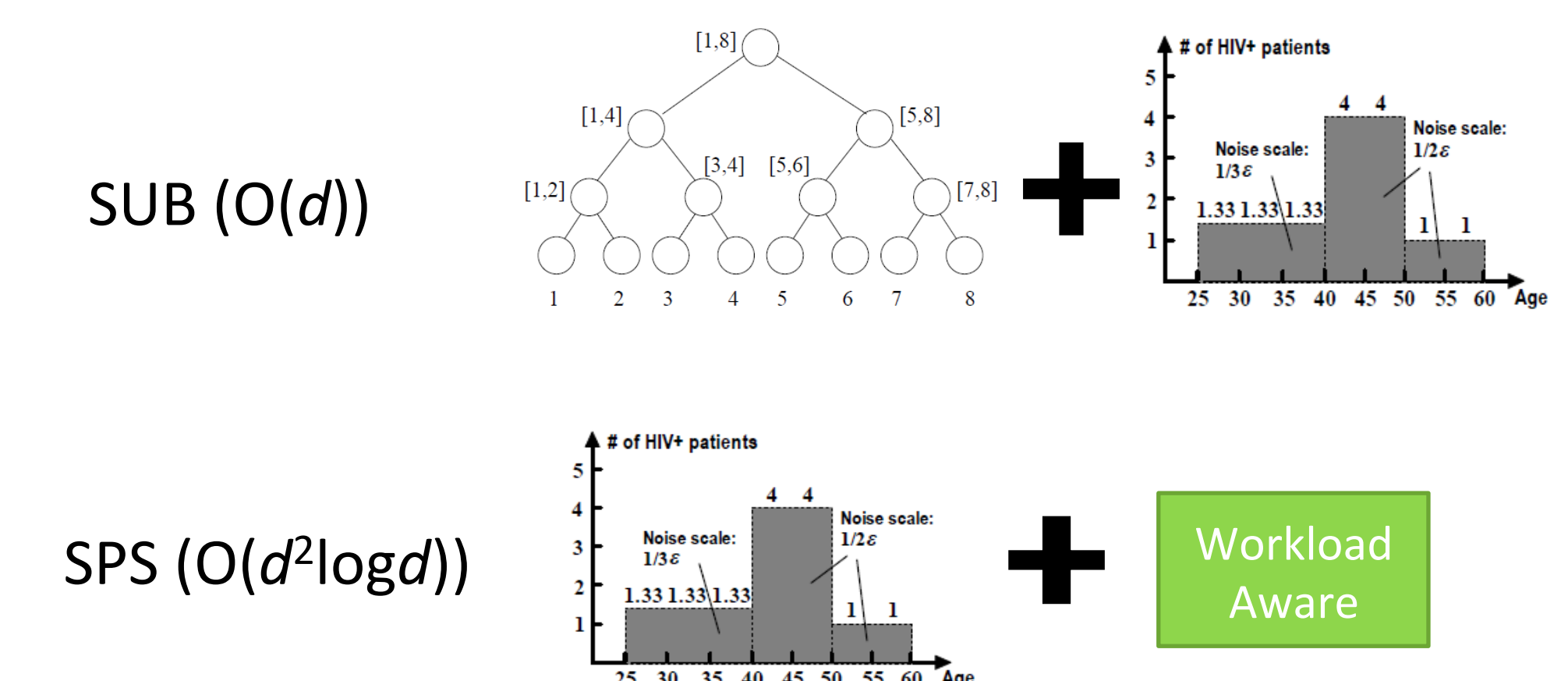


## Histograms for range queries

- Every existing method can be reproduced from three modules by parameterizing them

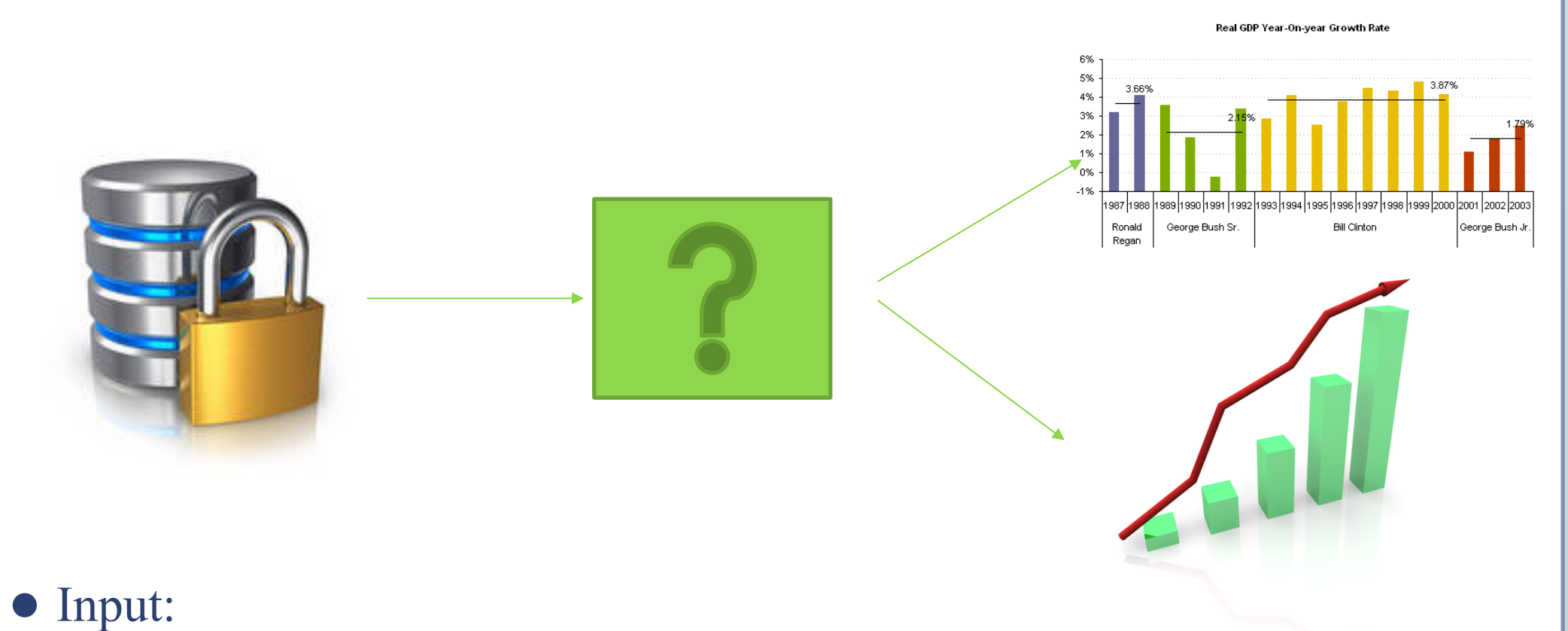


- Combine modules in order to devise new efficient and accurate methods



## Goal

- Design a system that privately decides which method should be used in order to increase the accuracy of the result



- Input:
  - Original Data
  - Statistics
- Choose the best method/combination privately
- Compute the result
- Output noisy statistics

## Challenges

- What are the data characteristics that render a method better than another in terms of accuracy?
- Can we determine the best method by utilizing non-private information on the input data?
- If not, can we utilize only statistics that have very low sensitivity?
- How can we choose the best method on-the-fly? (Streaming setting)

## References

[DMNS06] Calibrating noise to sensitivity in private data analysis. Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. TCC 2006.  
 [DNPR10] Differential privacy under continual observation. Cynthia Dwork, Moni Naor, Toniann Pitassi, and Guy N. Rothblum. STOC 2010.  
 [KP13] Practical Differential Privacy via Grouping and Smoothing. Georgios Kellaris, and Stavros Papadopoulos. VLDB 2013.  
 [KXP14] Differentially private event sequences over infinite streams. Georgios Kellaris, Stavros Papadopoulos, Xiaokui Xiao, and Dimitris Papadias. VLDB 2015.  
 [LHMW14] A Data- and Workload-Aware Algorithm for Range Queries Under Differential Privacy. Chao Li, Michael Hay, Jerome Miklau, and Yue Wang. VLDB 2014.  
 [QYL13] Understanding Hierarchical Methods for Differentially Private Histograms. Wahbeh Qardaji, Weining Yang, and Ninghui Li. VLDB 2013.  
 [ZCXM14] Towards accurate histogram publication under differential privacy. SDM 2014.