# PROCHLO: Strong Privacy for Analytics in the Crowd

Andrea Bittau[1], Úlfar Erlingsson[1], Petros Maniatis[1], Ilya Mironov[1], Ananth Raghunathan[1], David Lie[2], Mitch Rudominer[3], Ushasree Kode[3], Julien Tinnes[3], Bernhard Seefeld[3]

[1]Google Brian          [2]Google Brian and U. Toronto          [3]Google

Presenter: Jinyang Li (jinyang7)

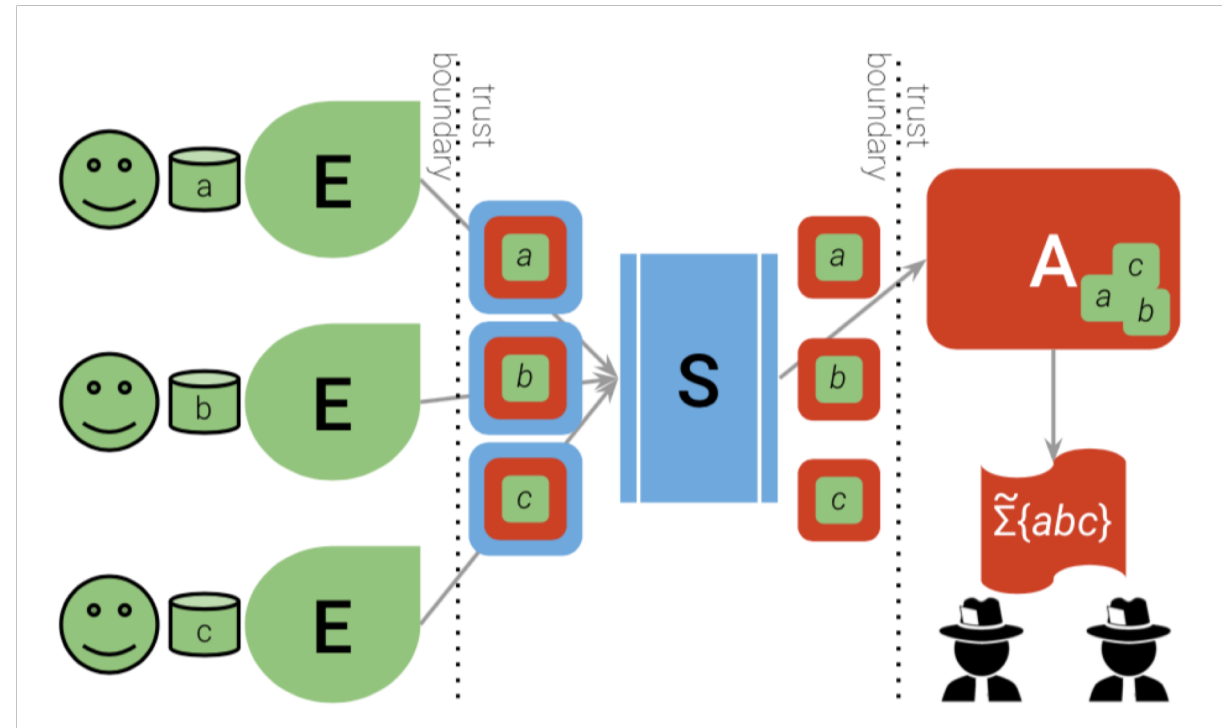# Systems Analytics and Privacy

Monitor API usage on software platforms

How to do the analytics?

How to handle the <u>private</u> data carefully?
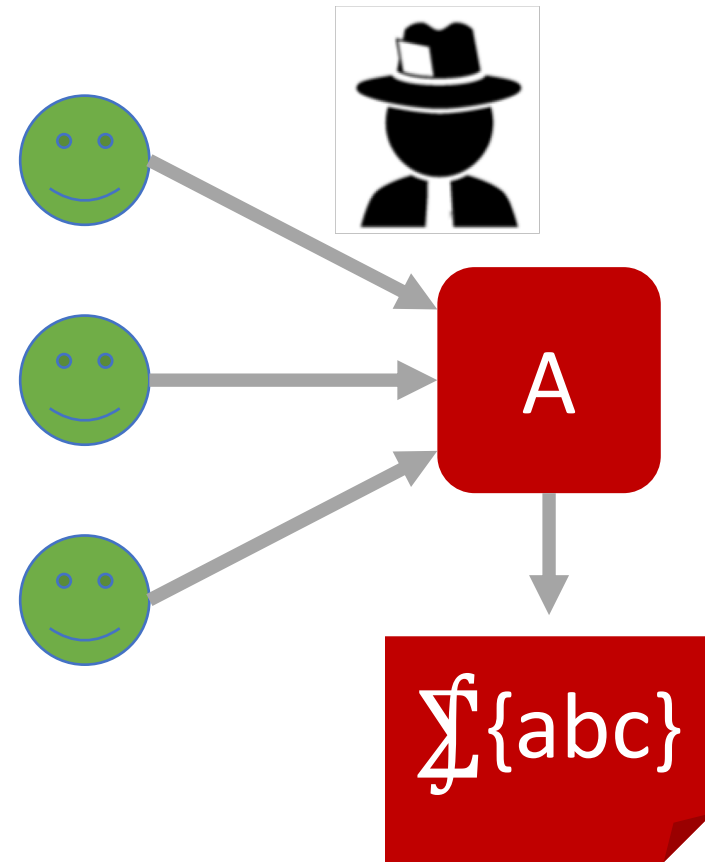
# ESA Architecture and Prochlo Realization

- Perform such monitoring with
  - High utility
  - Strong privacy guarantees
- Encode, Shuffle, Analyze
  - Framework for monitoring
  - Privacy protection
  - Fit to software engineering
- Prochlo
  - A hardened ESA realization
  - SGX
  - Oblivious shuffling
  - Threshold crypto & blinding

# Naïve API Monitoring

- What could go wrong?

- Uncommon API in an unpopular App

- At-least-K uses of an API

- Hard to get right!
  - Certain groups favor certain app features
  - IP address may reveal location
  - Etc.
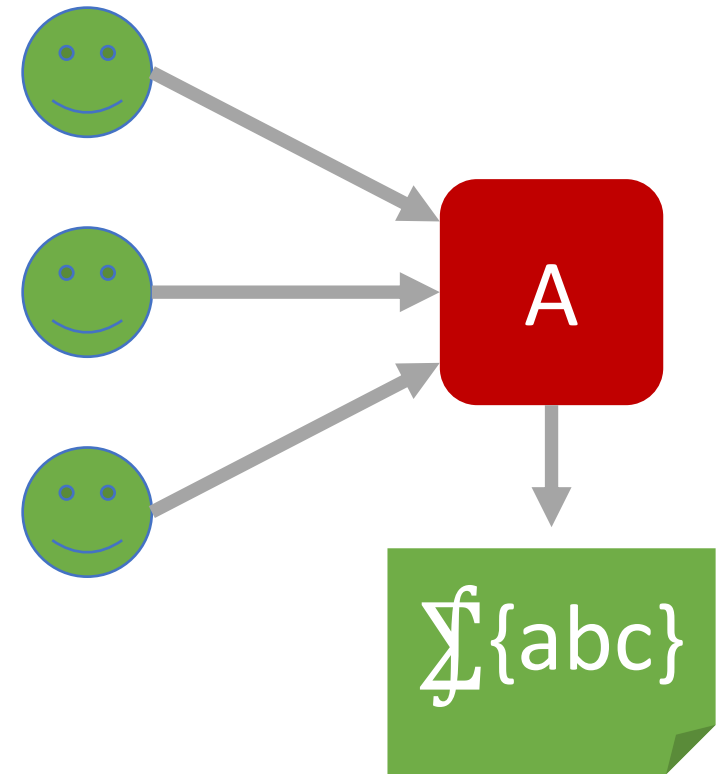
# Differential-Privacy Data Analysis

- DP gives (ε, δ) upper-bound on privacy loss
- Holds for all questions & for all attackers priors

$$Pr[M(D) \in S] \leq e^{\epsilon} Pr[M(D') \in S] + \delta$$
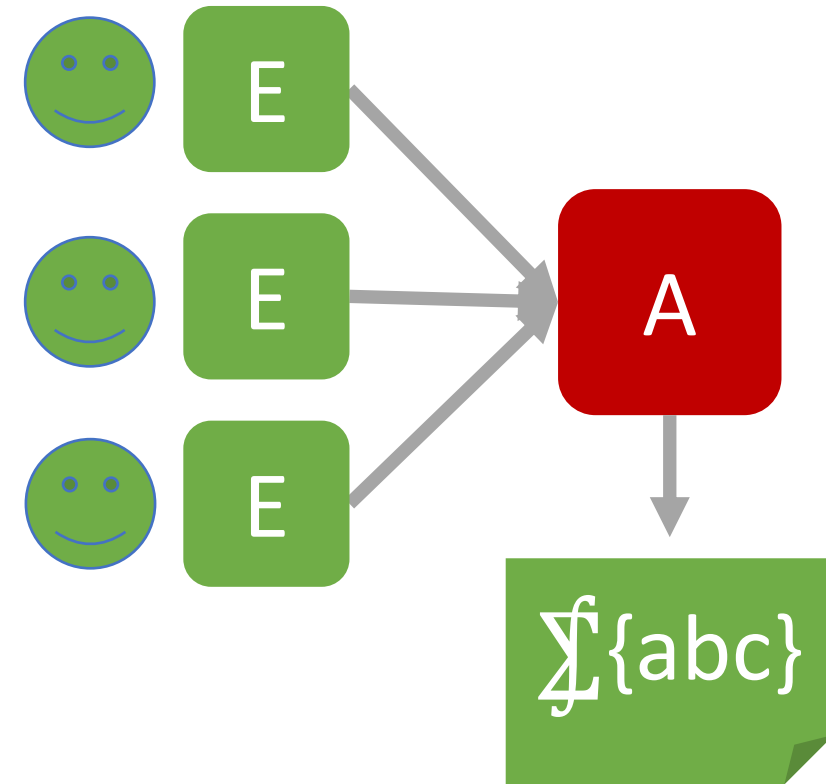
Multiplicative upper bound

Very small failure rate

- Fundamental flaw with using database in DP
- Bad fit for software engineering
  - New algorithms and systems
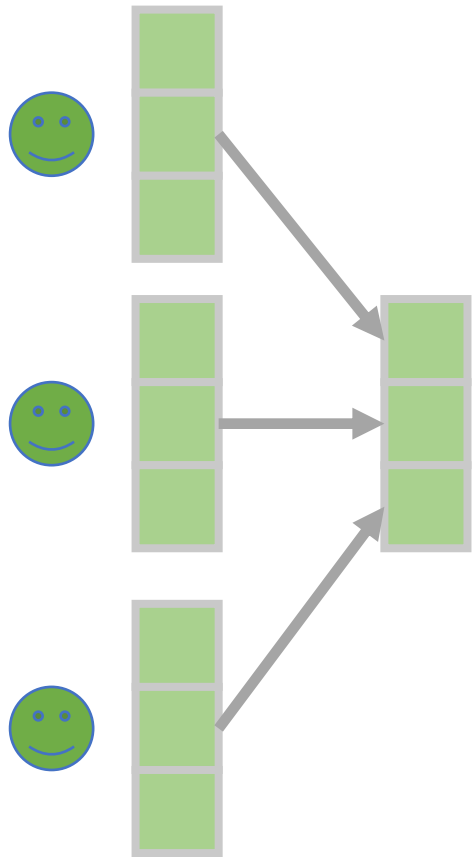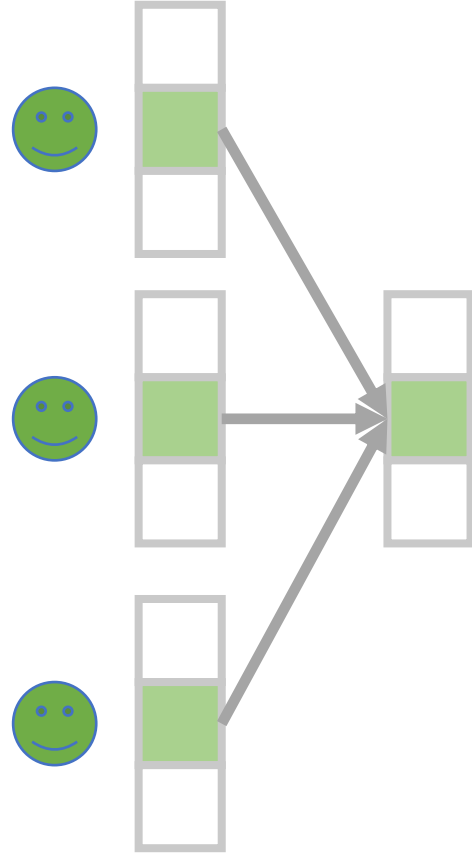  - Protect the databases forever

# Randomization & Local DP

- Randomized response
- No central (hackable) DB of real, private user data
- Google's RAPPOR system
  - Software monitoring system for Chrome
  - Since 2014
  - Largest deployed differential private mechanism solution
  - Dozens of purposes, billions of randomized daily reports
- Limitations:
  - Only good for very popular things and very large datasets
  - Too statistical
  - Too much noise (grows as sqrt(#reports))
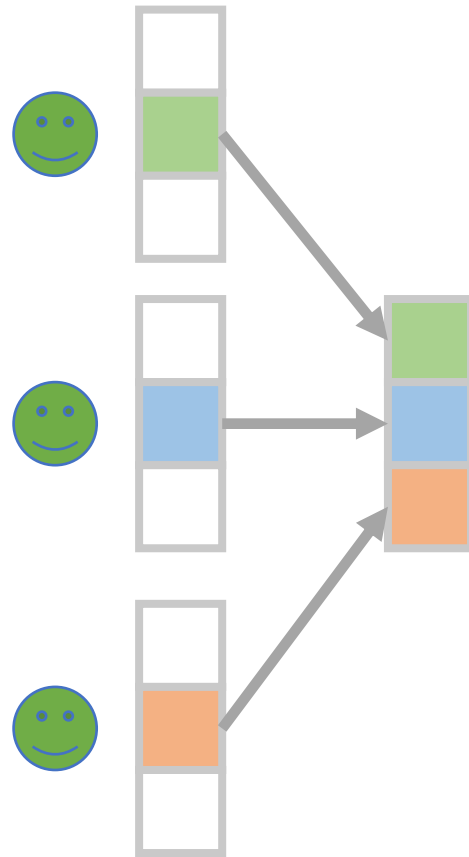
# Encoding Fragments and Crowds
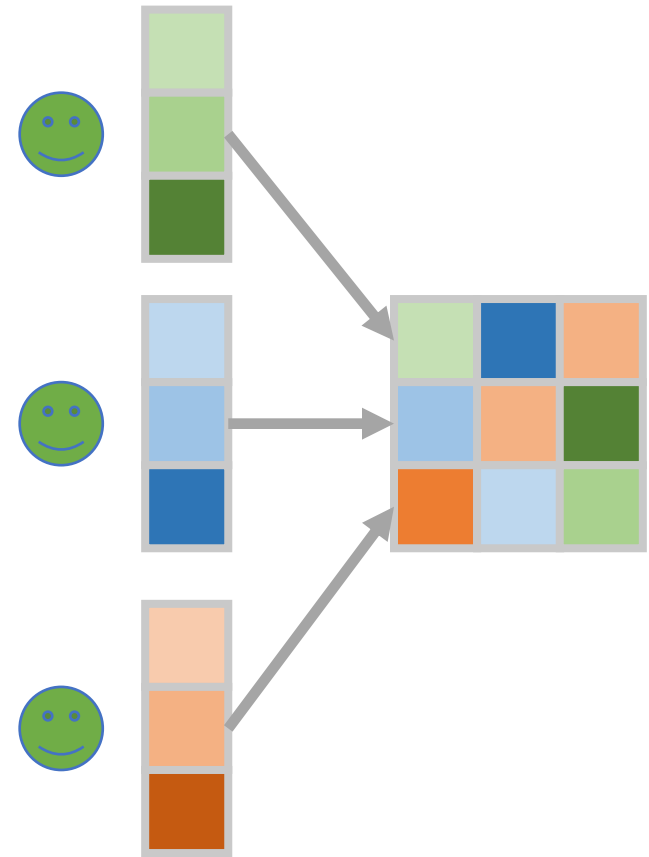


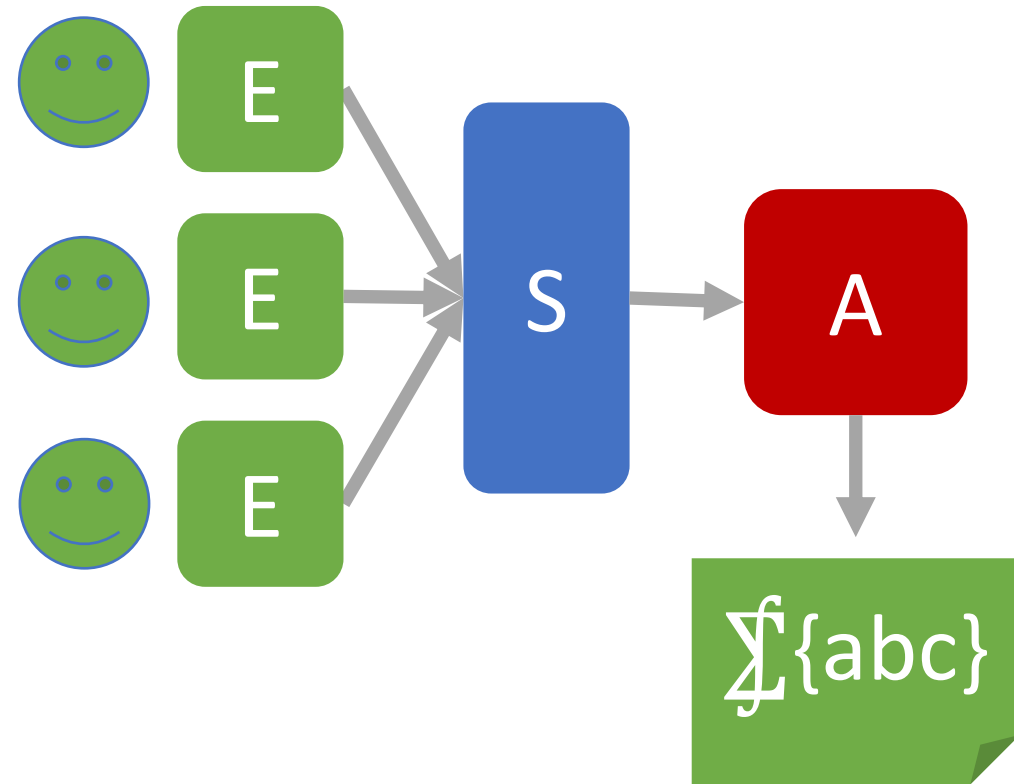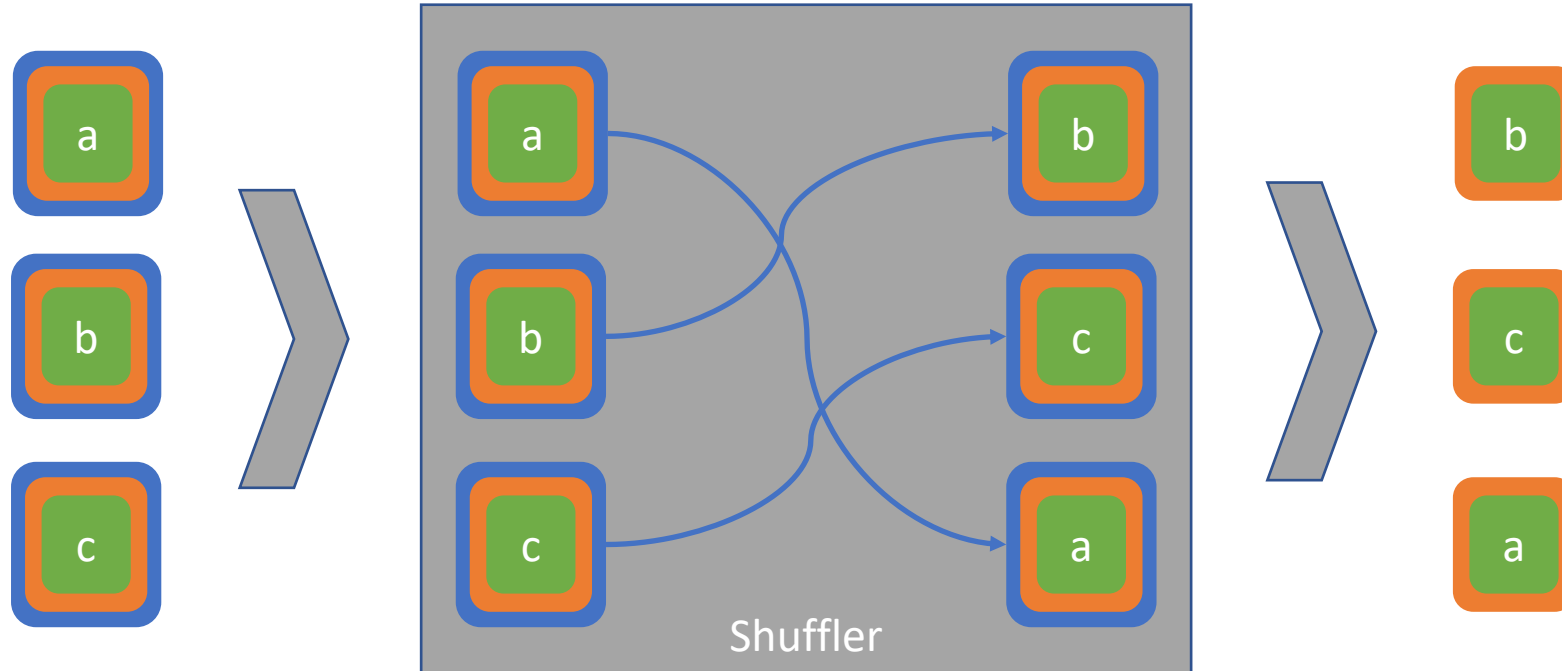Secret Share          Filter          Sample          Fragment

# Anonymity, Batching and Shuffling

- Shuffling provides anonymity
  - Strips IP address & metadata
- Create big crowds
  - By delaying and batching
  - Per-day, in 100s of millions
- Randomly shuffle the reports
  - Break linkability between fragments
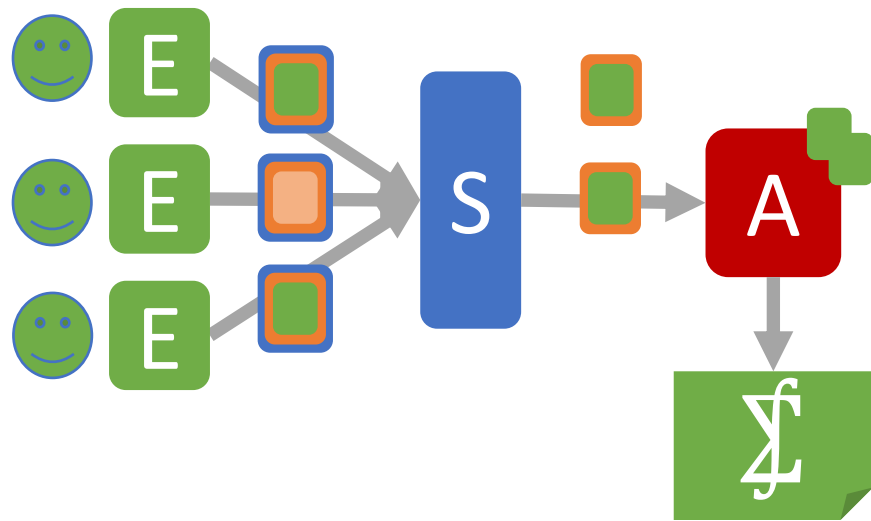  - Hide ordering and timing information
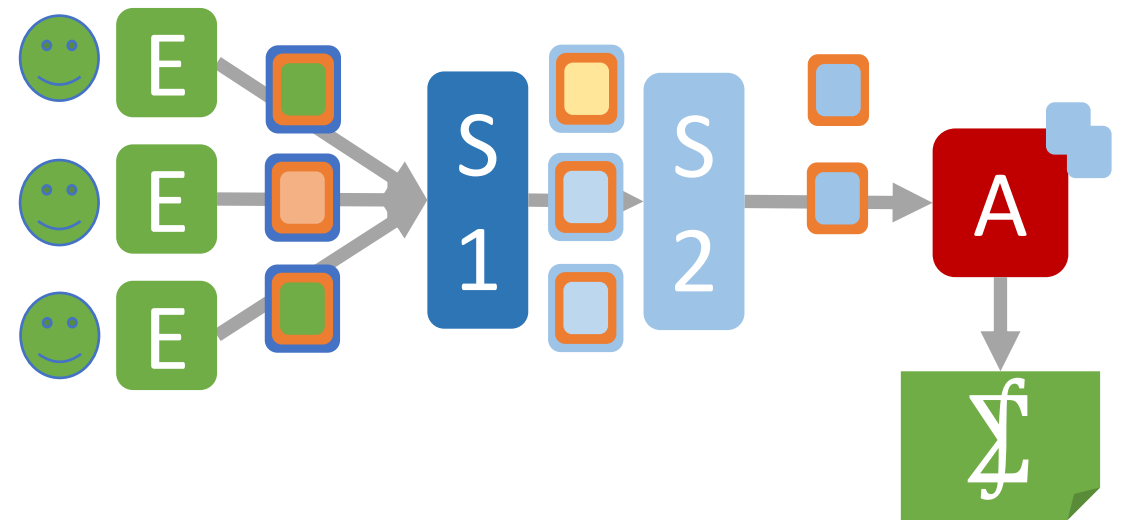
# Shuffling and Nested Encryption

# Randomized Thresholds, Blinding, and Crowd-based DP

## Randomized thresholding gives another form of DP
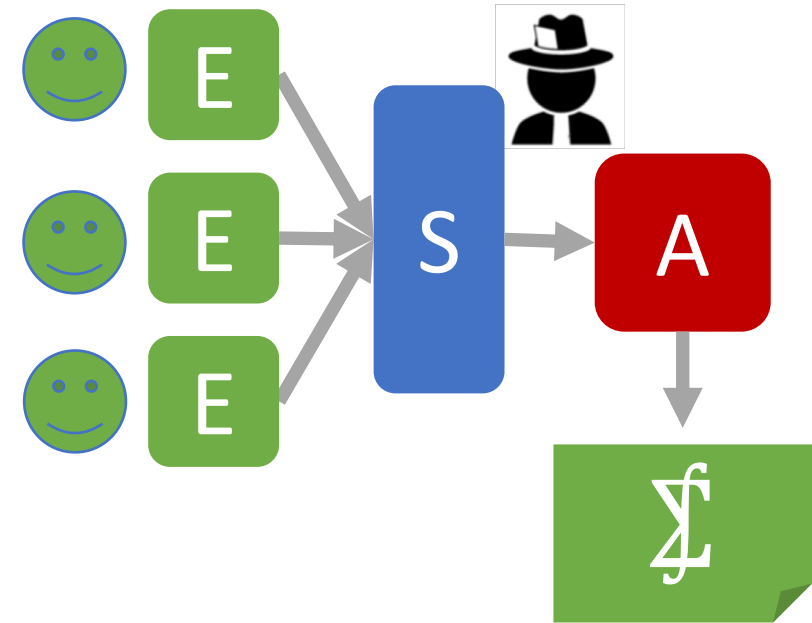
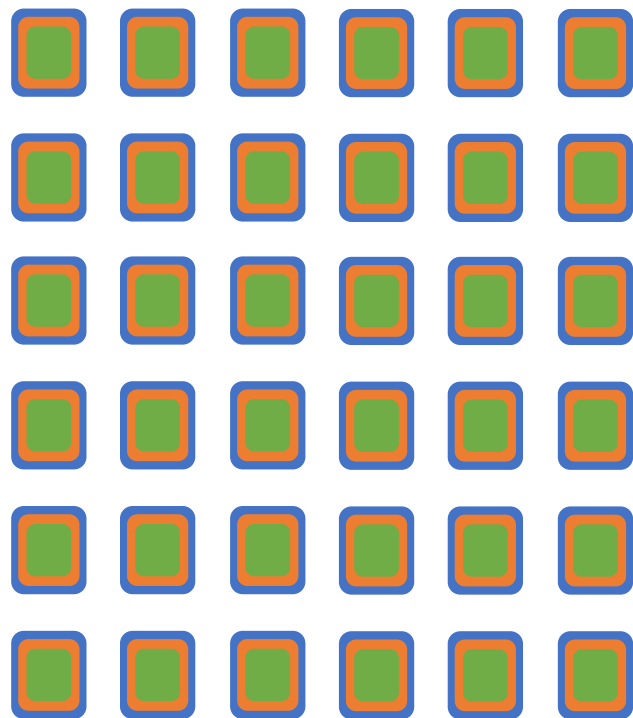

Blinding & Crowd Thresholding

+ Cryptographic blinding of crowds

# Risks in the ESA Shuffler

- Shuffling must be protected, isolated & opaque

- Insider risk, accidental server logs, etc.

- Malicious traffic analysis

- Prochlo
  - Hardened implementation of ESA shuffler
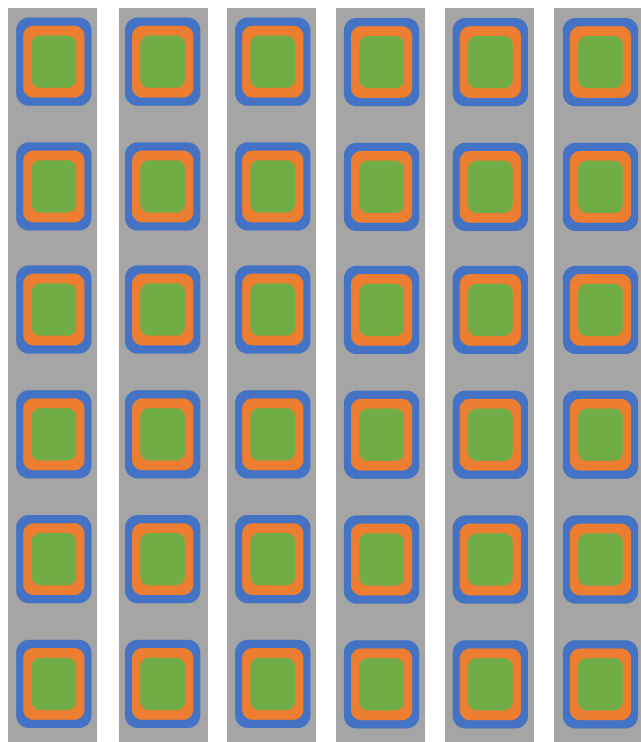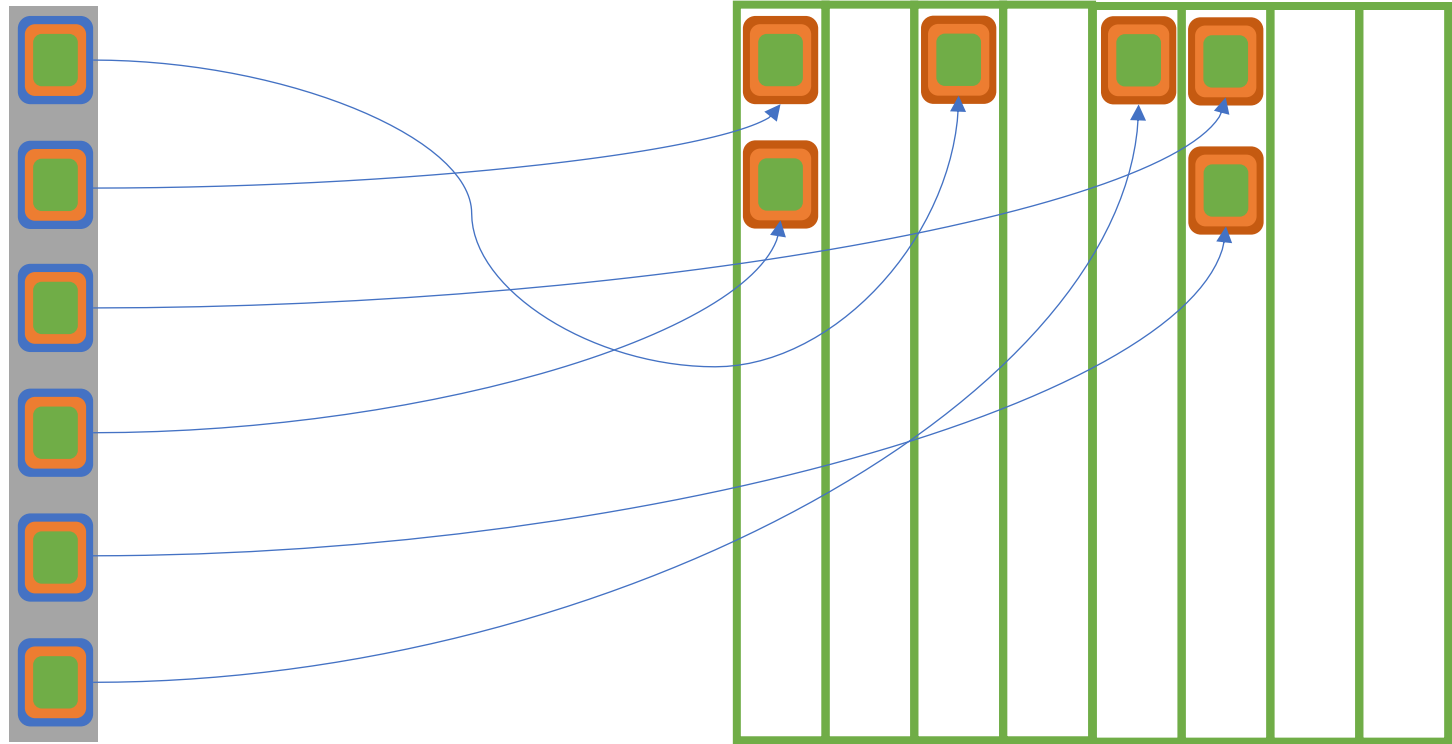  - SGX + oblivious shuffling
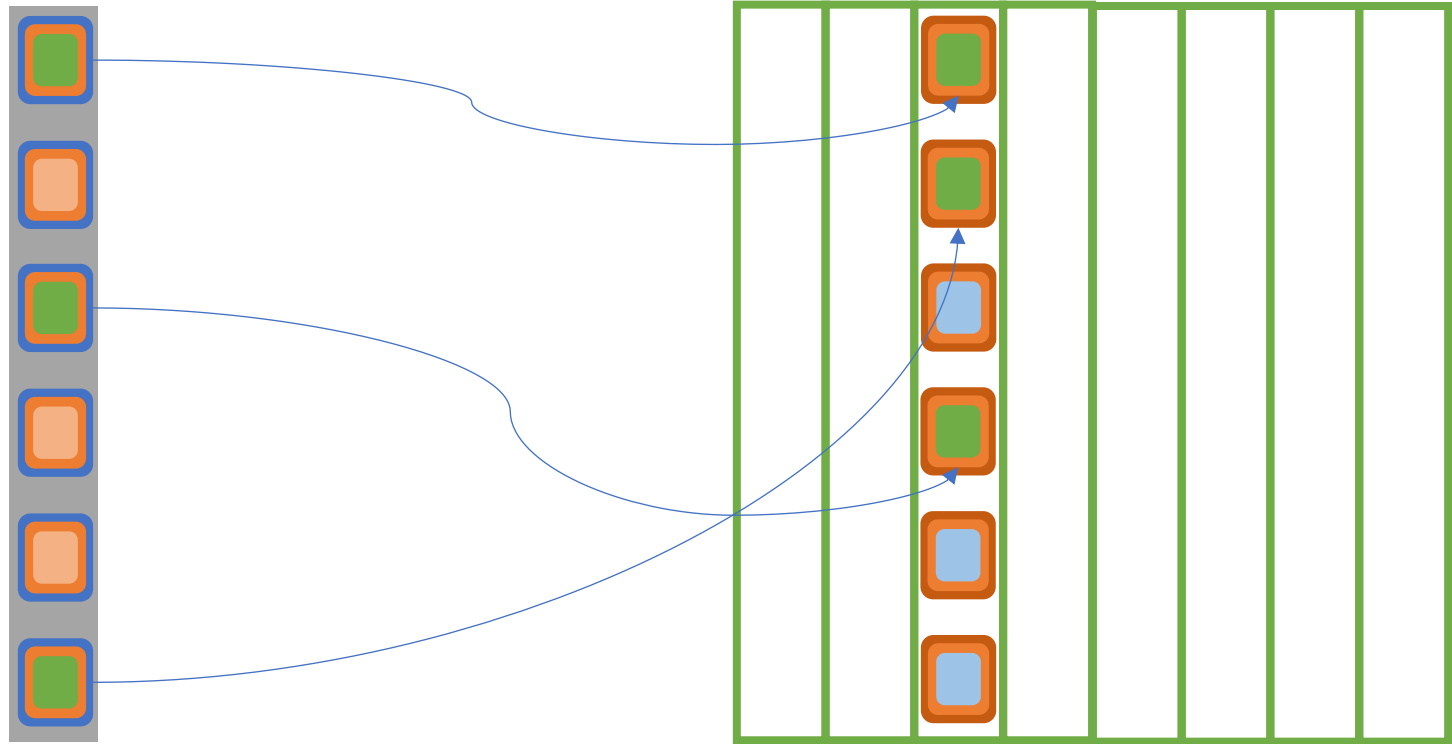
# Prochlo StashShuffle
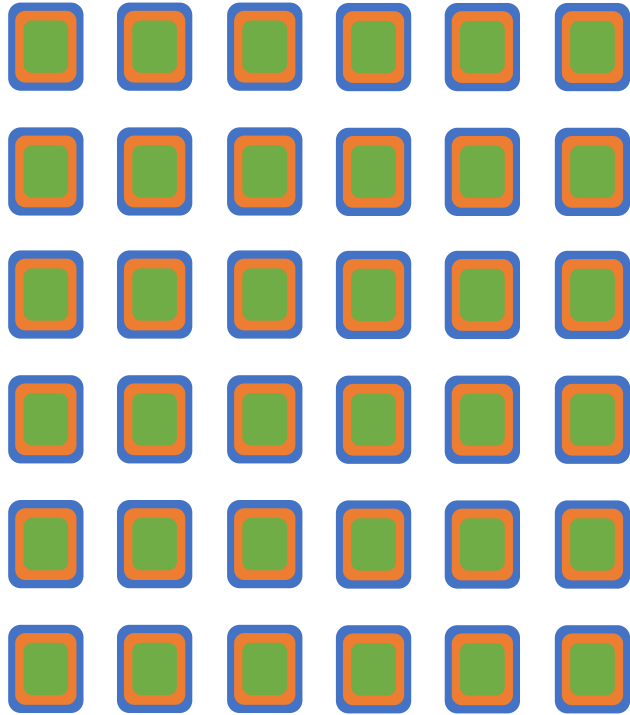


SGX

# StashShuffle Buckets

# StashShuffle Distribution



Intermediate Array

# StashShuffle Compression
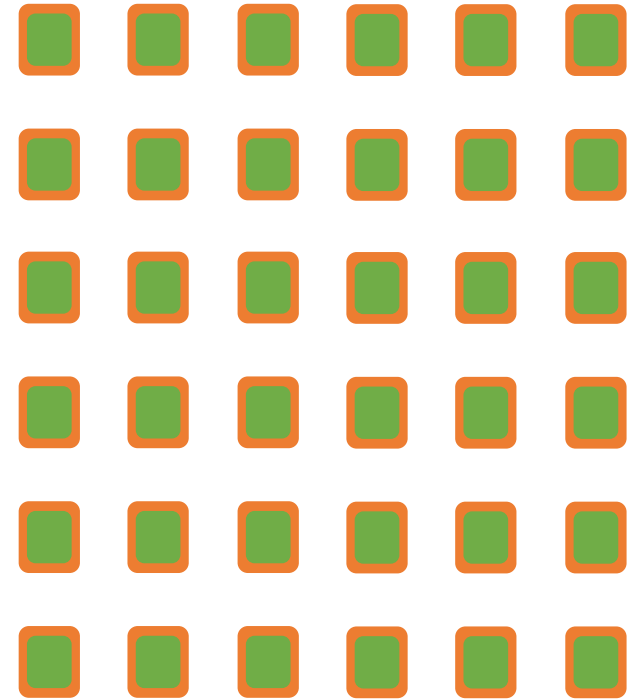


Output bucket

# Prochlo StashShuffle



SGX

# Shuffler Performance

| N | Permutation strength | Time to shuffle | SGX mem used | Overhead of passes |
|---|---|---|---|---|
| 10M | $2^{-80.1}$ | 738 s | 22 MB | 3.5x |
| 50M | $2^{-81.8}$ | 1 h | 52 MB | 3.4x |
| 100M | $2^{-81.9}$ | 2.1 h | 78 MB | 3.7x |
| 200M | $2^{-64.5}$ | 4.1 h | 69 MB | 3.3x |

# Utility Performance

# More experiments

- Perms: User Action Regarding Permissions
  - Multidimensional like API example
  - High utility with strong privacy $\varepsilon = 1.2$, $\delta = 10^{-7}$
- Suggest: Predicting the Next Content Viewed
  - High utility with intuitive privacy guarantee due to fragments
- Flix: Collaborative filtering
  - Utility equals state-of-the-art joint-distribution model
  - Strong privacy ($\varepsilon = 2.2$) + anonymity = no chance of re-identification

# Conclusion

- Making strong privacy suitable for use in standard software engineering

- Open source:
  - https://github.com/google/rappor
  - https://github.com/google/prochlo
  - https://fuchsia.googlesource.com/cobalt/

# Discussion

- "Just trust Intel" vs. "Just trust Google"
- Attack model: "Shuffler is honest-but-curious"
- Large latency