

BLIP: Non-interactive differentially-private similarity computation on Bloom filters

Mohammad Alaggan¹, Sébastien Gambs^{1,2}, Anne-Marie Kermarrec²

1 Université de Rennes 1 – IRISA, Rennes, France

2 INRIA Rennes Bretagne-Atlantique, Rennes, France

4 October 2012

Context

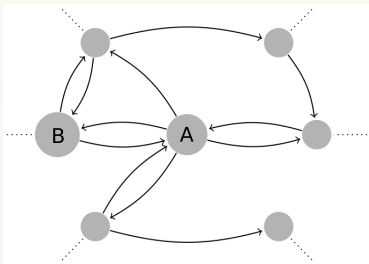
- ▶ A **collaborative social system**
- ▶ Example: each user has a **profile**, which is a list of **items** he has tagged/liked



Application (example)

Decentralized **personalized** search BFGKL^{a10}

^aBertier, Frey, Guerraoui, Kermarrec and Leroy.



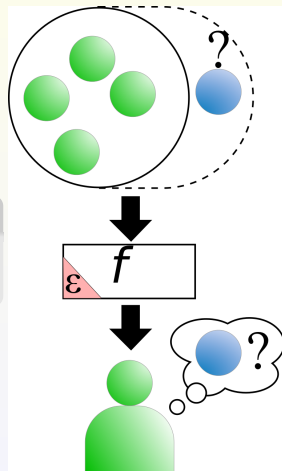
Challenge

Using the personal information while preserving **privacy**

Differential Privacy

Differential privacy (Dwork 2006) (informal)

Inclusion of a particular item should not affect the probability of a particular output except by very small amount



Differential Privacy

Definition (ϵ -Differential privacy (Dwork 2006))

A randomized algorithm \mathcal{A} satisfies ϵ -differential privacy if for all neighboring databases \mathbf{d}, \mathbf{d}' , and outputs t :

$$\Pr[\mathcal{A}(\mathbf{d}) = t] \leq \exp(\epsilon) \Pr[\mathcal{A}(\mathbf{d}') = t] ,$$

where \mathbf{d} and \mathbf{d}' differ on *at most* one item and $\epsilon > 0$ is the privacy parameter

Implemented *interactively* to our system in AGK^{a11}

^aAlaggar, Sébastien and Kermarrec.

Motivation for non-interactiveness

Composability (McSherry 2009)

If function f is ϵ -differentially private then $f(x; r_1)$ and $f(x; r_2)$ *together* are (2ϵ) -differentially private (less private)

Informally: having computed the same differentially private function twice, the randomness could be *averaged out*

Motivation for non-interactiveness

Privacy budget (DMNS^{a06})

If ϵ_0 privacy is desired, and the function f to compute is ϵ -differentially private, do *not* compute and release f more than ϵ_0/ϵ times

A problem in very large scale systems

^aDwork, McSherry, Nissim and Smith.

Randomized response differential privacy (BNO^a08)

Given a profile of ℓ bits ($\{0, 1\}^\ell$), flip each bit with probability $1/(2 + \epsilon)$.

^aBeimel, Nissim and Omri.

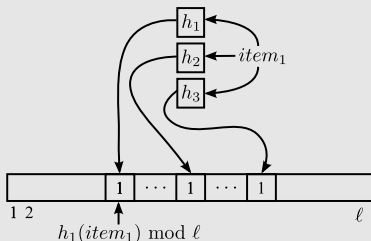
Questions

- ▶ Can this technique be applied to profiles of unbounded size on infinite universe ($\{0, 1\}^*$)? (infinitely many zeroes to flip!)
- ▶ Can the probability of flipping be reduced (less noise, more utility) while preserving differential privacy?

Our contribution

- ▶ Yes to both questions

Bloom filter^a



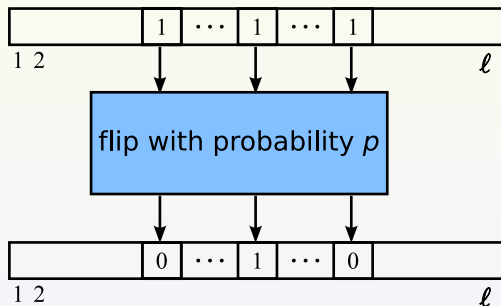
- ▶ A compact representation of profiles using k different hash functions
- ▶ Bloom filters have a **false-positive** probability (could be tweaked)
- ▶ Could be exhaustively **queried** to *approximately* reconstruct the profile

^aBloom 1970

BLIP mechanism

Bloom-then-flip mechanism

Flip each bit of the Bloom filter with probability $p < 0.5$



Challenge: How does flipping individual *bits* affect the privacy of the *items* encoded in the Bloom filter in several bits?

BLIP

Theorem (Optimal flipping probability)

Flipping each bit with probability $p = 1/(1 + e^{\epsilon/k})$, where k is the number of hash functions, satisfies differential privacy for items.

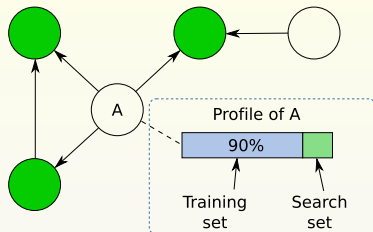
Moreover, this probability is optimal (no lower flipping probability can satisfy differential privacy.)

Theorem (Error bound)

Given two flipped Bloom filters, the additive error of their scalar product (the similarity metric) is $\Theta(\sqrt{\ell})$ with constant probability which is asymptotically optimal. (Lower bound by MMPRTV^a 11.)

^aMcGregor, Mironov, Pitassi, Reingold, Talwar and Vadhan.

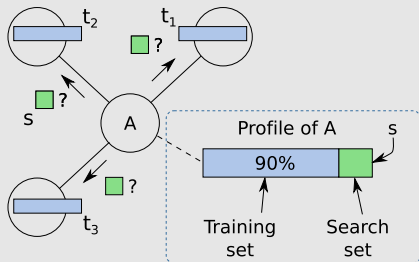
Experimental setup



Utility

Recall: Percent of **search items** found in the collective sets of his neighbors

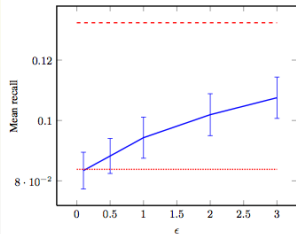
$$r = \frac{|s \cap (\cup_j t_j)|}{|s|}$$



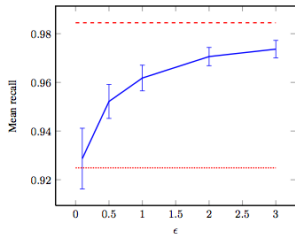
Datasets

- ▶ **Delicious**
 - ▶ Tagged bookmarks, 500 users, 52000 items, 135 average profile size
- ▶ **Digg**
 - ▶ Shared URLs, 500 users, 1237 items, 317 average profile size
- ▶ **Survey**
 - ▶ News items liked, 120 users, 196 items, 68 average profile size

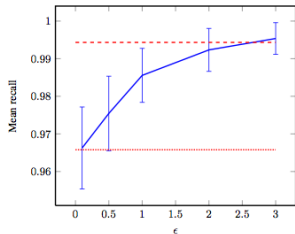
Utility measurement



(a) Delicious (100 runs)



(b) Digg (100 runs)



(c) Survey (500 runs)

Choosing the privacy parameter ϵ

- ▶ AACP^a11 works out the amount of information leaked about each item
- ▶ LC^b11 requires knowledge of the queries to be computed and is not applicable to non-interactive differential privacy

^aAlvim, Andrés, Chatzikokolakis and Palamidessi.

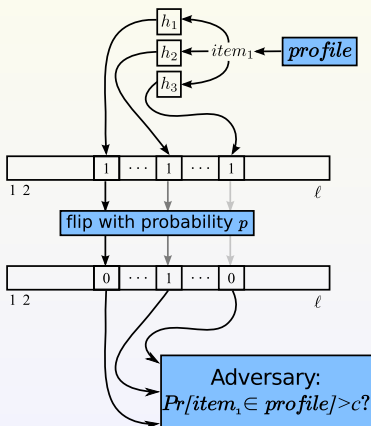
^bLee and Clifton.

Our contribution

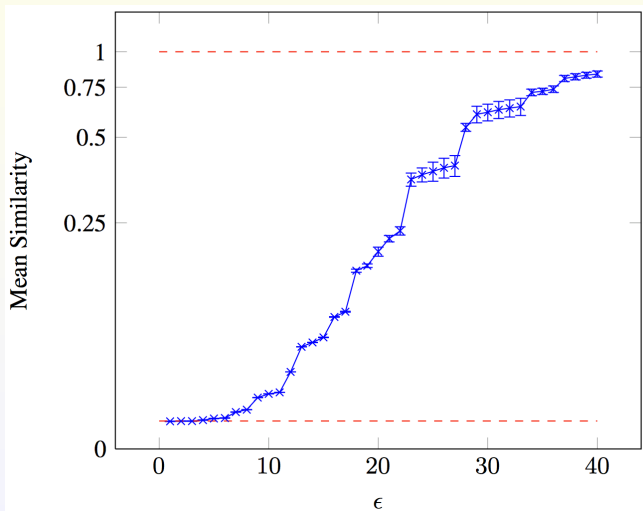
- ▶ A technique to upper-bound ϵ in our setting

Privacy measurement: Profile reconstruction attack

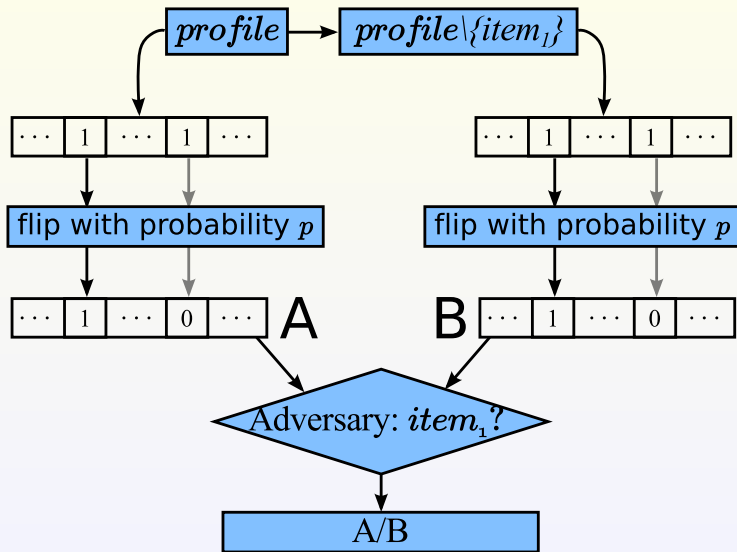
- ▶ Adversary should not be able to reconstruct the private profile
- ▶ Gives an upper bound for ϵ



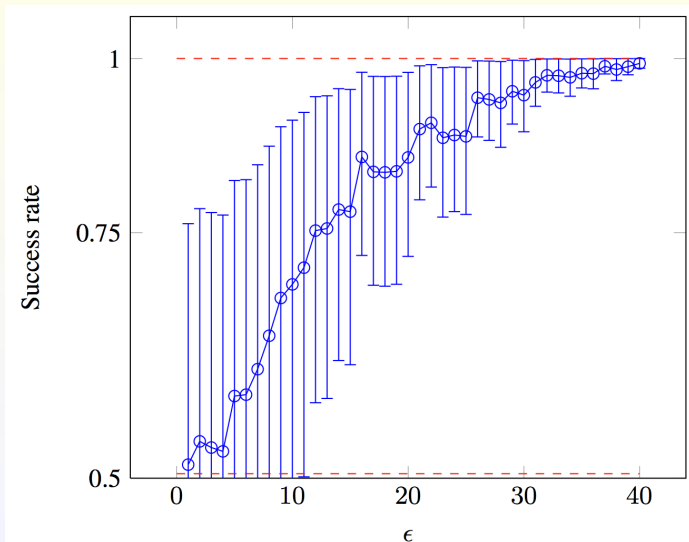
Profile reconstruction attack



Profile distinguishing game



Profile distinguishing game



Conclusion

Conclusion

- ▶ Non-interactiveness for very large scale systems
- ▶ Utility vs. privacy trade-off

Potential applications

- ▶ Non-interactive “matchmaking” protocols

Perspectives

- ▶ Handling dynamic updates
- ▶ Addressing the effect of correlation between items (raised by KM^{a11})
- ▶ Design of more sophisticated attacks

^aKifer and Machanavajjhala.

Thanks.

Questions ?