



Privacy Preserving Group Linkage

Fengjun Li¹, Yuxin Chen¹, Bo Luo¹, Dongwon Lee², and Peng Liu²

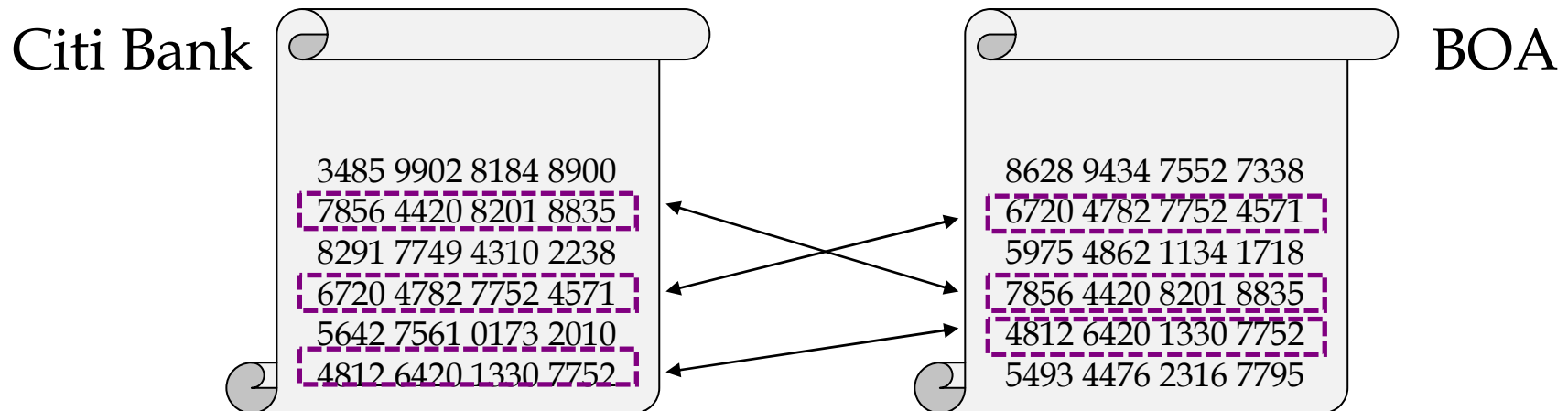
¹EECS Department, University of Kansas,

²College of IST, Penn State University



Record Linkage

- Record linkage is to identify related *records* associated with the same entity from multiple databases

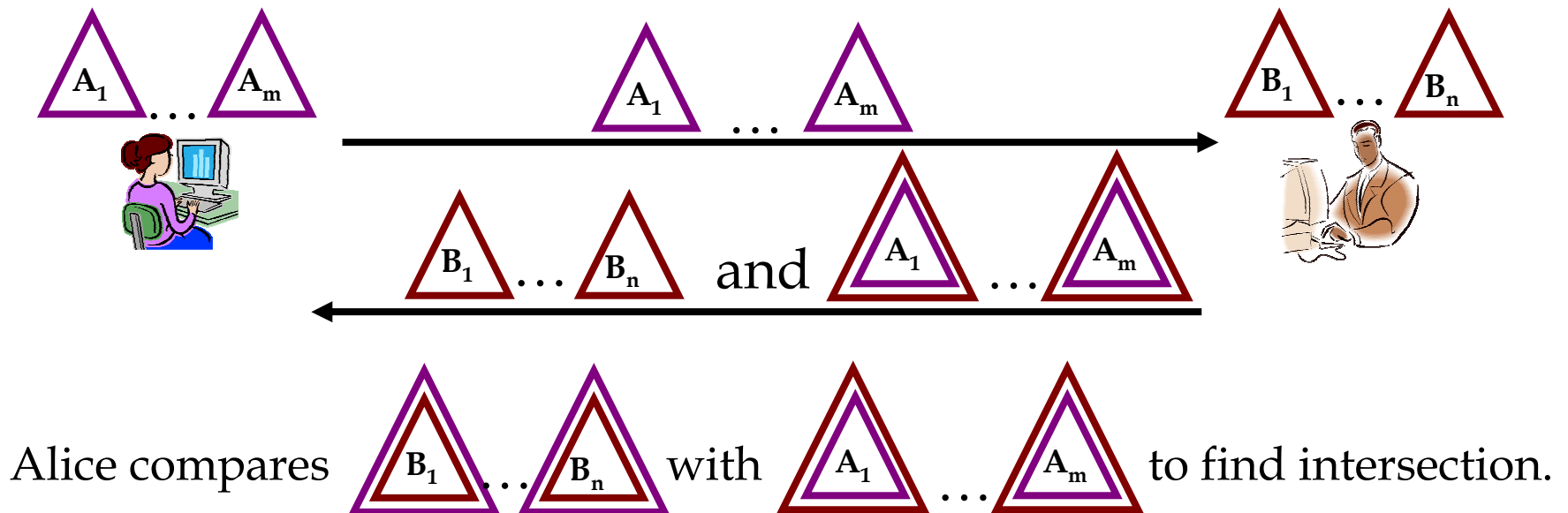


Privacy-Preserving Record Linkage

- Privacy becomes an issue when data is sensitive.
 - I will only share with you on the “linked records”
 - I will not give you the plain text of my primary keys.
- Secure multi-party set intersection problem
 - Solutions based on commutative encryption
 - Solutions based on homomorphic encryption

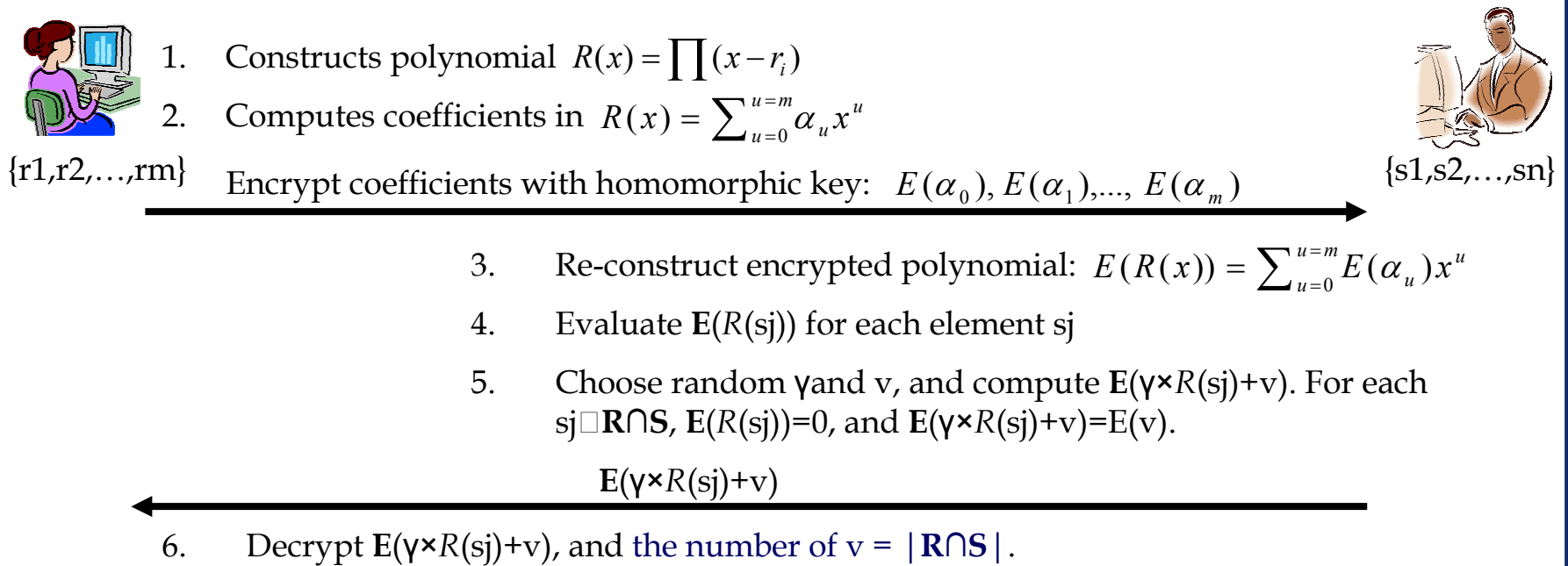
AES Protocol – Commutative Encryption Based

- **Commutative Encryption:** *using the same set of commutative keys, the encrypted content can be recovered in **any arbitrary order**.*
- **AES Protocol** [Agrawa et. al., SIGMOD 2003]:



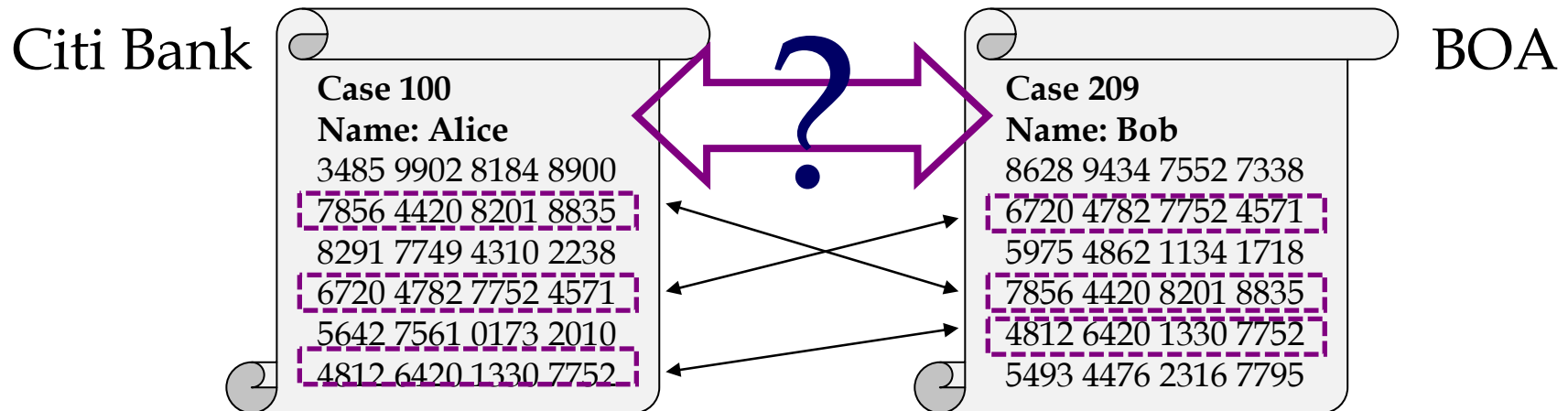
FNP Protocol – Homomorphic Encryption Based

- **Homomorphic encryption:** *allows certain algebraic operations in the plaintext to be performed on the ciphertext without decryption.*
- **FNP Protocol** [Freedman et. al., EUROCRYPT 2004]:



Group Linkage

- Extended from record linkage [On et. al., ICDE 2007]
 - Records -> groups of records

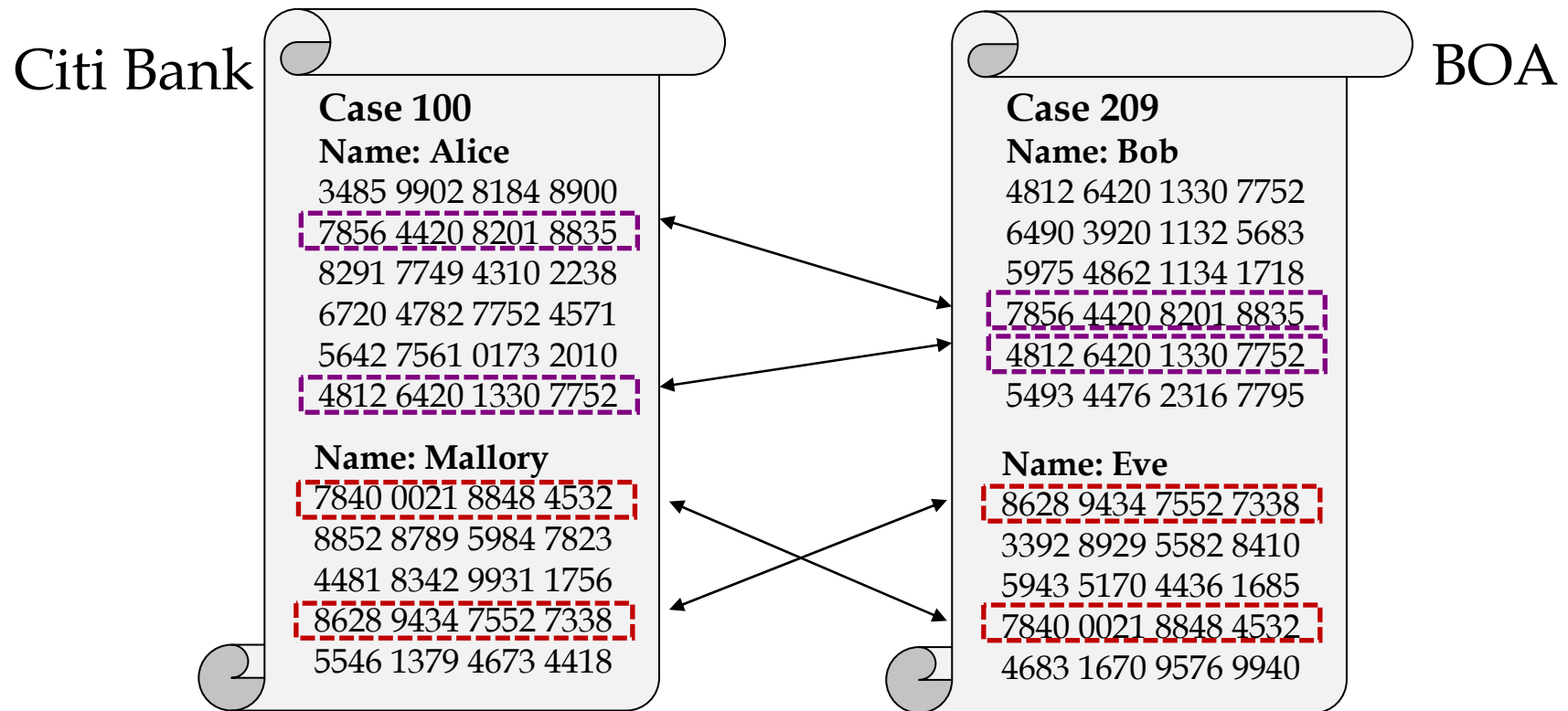


- Group linkage is to identify related *groups of records* associated with the same entity from multiple databases

Group Linkage

- For two sets of groups of records $\mathbf{R}=\{R_1, \dots, R_u\}$ and $\mathbf{S}=\{S_1, \dots, S_v\}$, GL calculates *group-level similarity* $\mathbf{SIM}(\mathbf{R},\mathbf{S})$, and determines if \mathbf{R} and \mathbf{S} are associated with the same entity
 - For $R=\{r_1,\dots,r_m\}$ and $S=\{s_1,\dots,s_n\}$, calculate *record-level similarity* $\mathbf{sim}(r,s)$
 - $\mathbf{SIM}(\mathbf{R},\mathbf{S})$ is a function of $\mathbf{sim}(r,s)$

Group Linkage: Exact Matching Example



Group Linkage: Approximate Matching Example

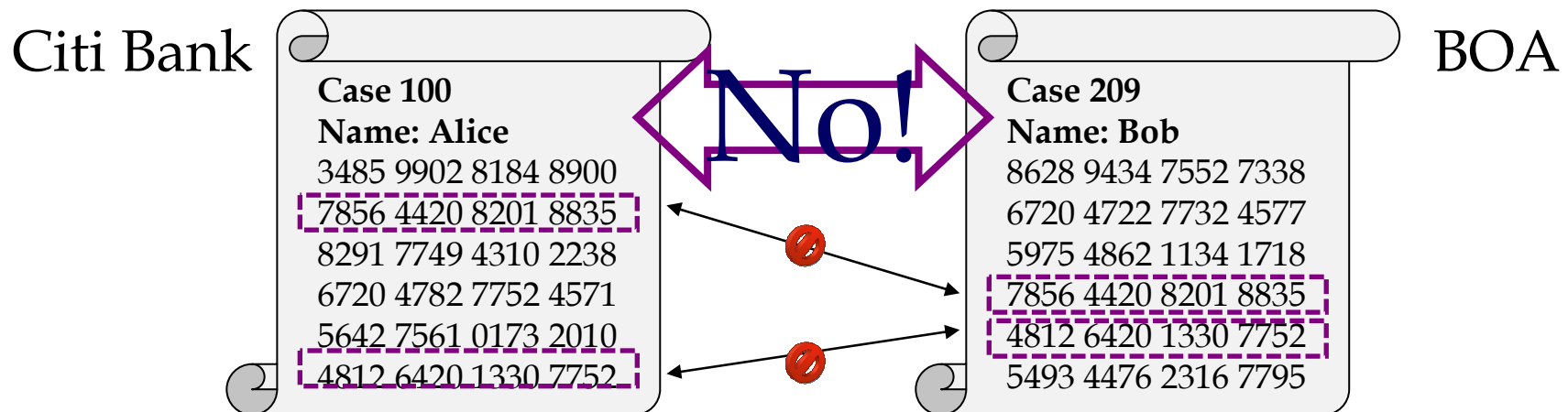
- modeling and representation of data, metadata, ontologies, and processes
- querying of scientific data



- modeling and representation of data and knowledge for scientific domains
- querying and analysis of scientific data.

Group Membership Inference Problem

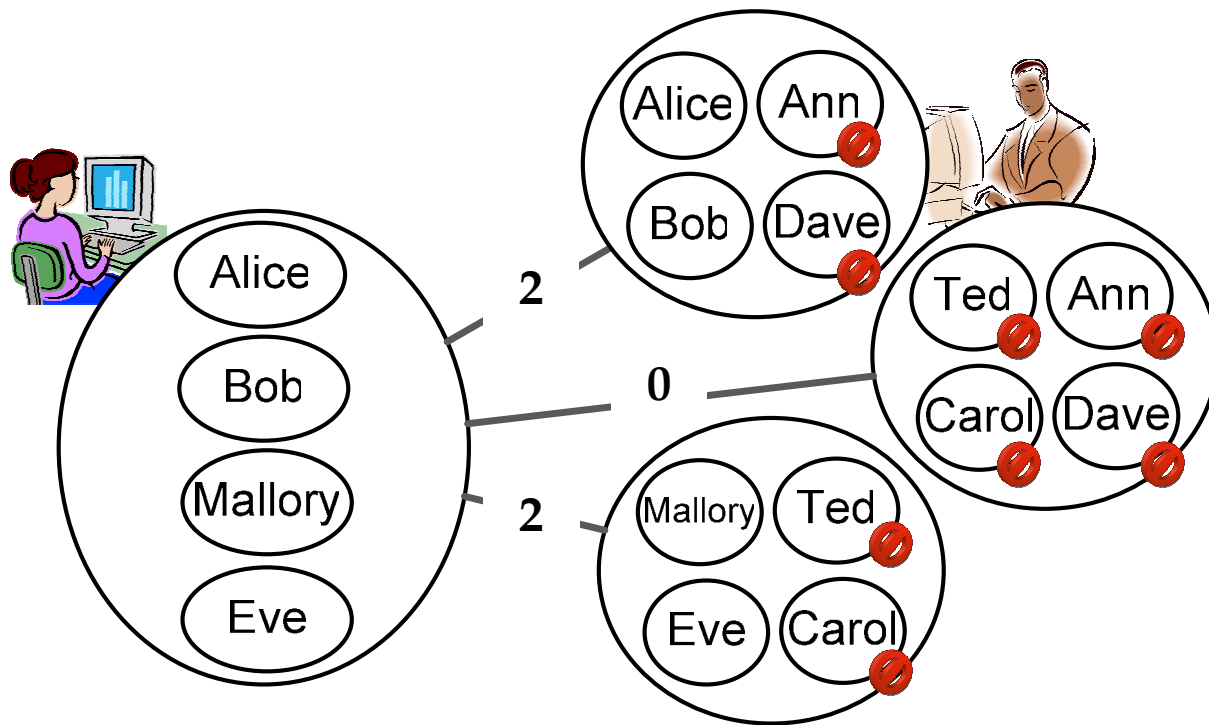
- Two parties share two groups after they confirm both groups are associated with the same entity.
- Privacy?
 - Cannot share “intersect” records when two groups are not linked.



Privacy-Preserving Group Linkage (PPGL)

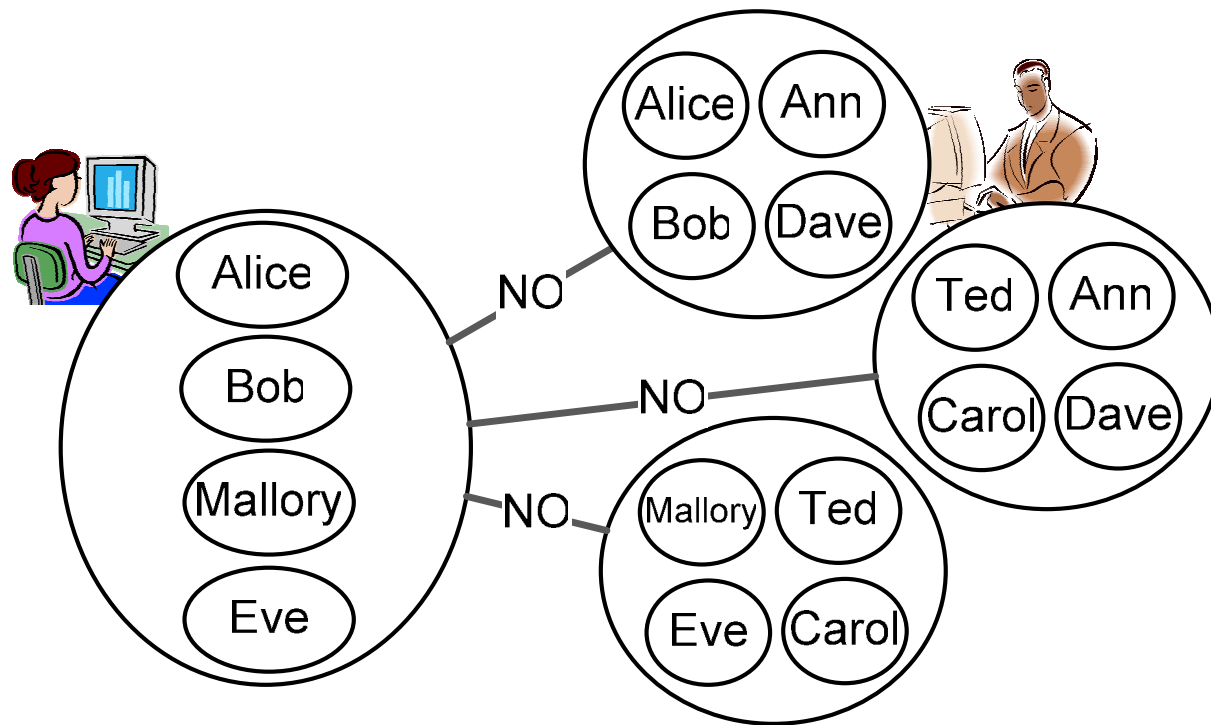
- PPRL protocols can be applied in PPGL
 - Secure set intersection size
 - The intersection size can be used to calculate group-level similarity
- However, directly applying PPRL protocol suffers from *group membership inference* problem

Group Membership Inference Problem



- Identities of overlapped group members can be inferred
- An attacker can manipulate the group members to infer more

Ideal PPGL Protocol



- Alice and Bob negotiate a similarity threshold
- For each group-wise comparison, Bob answers only “YES” or “NO”, instead of calculated similarity value

Threshold Privacy-Preserving Group Linkage

- **TPPGL Problem:** Alice and Bob preset a threshold θ , and follow the protocol to match two groups \mathbf{R} and \mathbf{S} . In the end, they learn only $|\mathbf{R}|$, $|\mathbf{S}|$, and a Boolean result B , where $B = \text{true}$ iff $\mathbf{SIM}(\mathbf{R}, \mathbf{S}) \geq \theta$.
- We propose three TPPGL protocols for both exact matching and approximate matching
 - K-combination approach for TPPGL-E
 - Homomorphic encryption approach for TPPGL-E
 - TPPGL-A protocol with record-level cut-off

K-Combination Approach for TPPGL-E

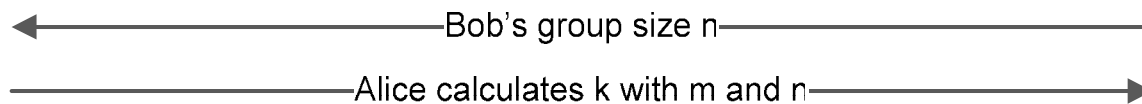
- Alice has a set of groups $\mathbf{R}=\{r_1,\dots,r_m\}$, and Bob has a set of groups $\mathbf{S}=\{s_1,\dots,s_n\}$. They negotiate a similarity threshold θ .
- Calculate the *minimum number of identical records* in \mathbf{R} and \mathbf{S} for them to be linked

$$\mathbf{SIM}(\mathbf{R}, \mathbf{S}) = k/(|\mathbf{R}|+|\mathbf{S}|-k) \geq \theta, \quad \text{so} \quad k = \left\lceil \frac{(m+n)\theta}{1+\theta} \right\rceil$$

- We enumerate all *k-combinations* of Alice's and Bob's group elements. \mathbf{R} and \mathbf{S} are linked iff there is at least one identical k-combination.

Input: Alice's group $\mathbf{R}=\{r_1,\dots,r_m\}$, Bob's group $\mathbf{S}=\{s_1,\dots,s_n\}$, and a pre-negotiate similarity threshold θ

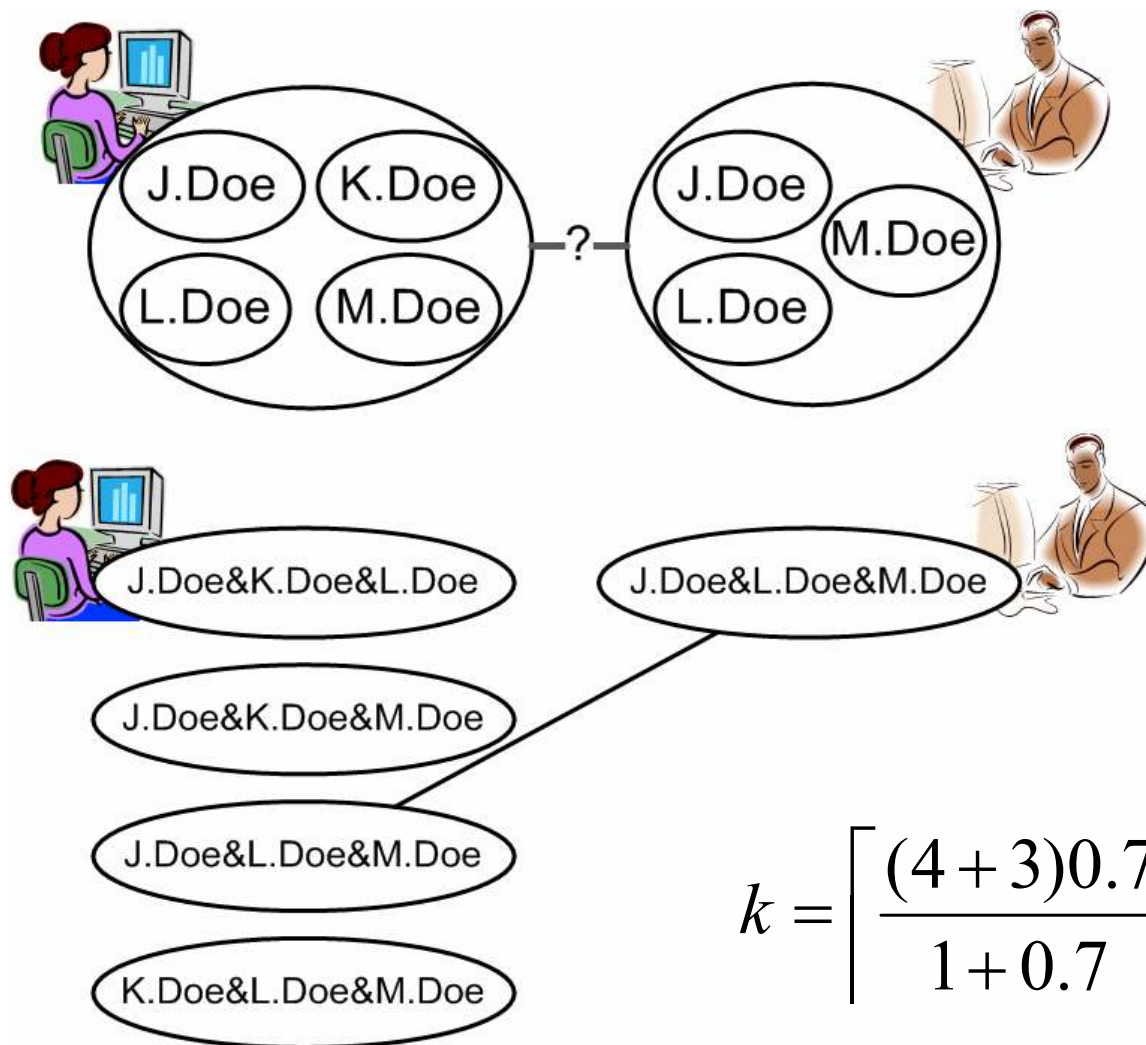
Protocol:



- Alice creates $p = C_k^m$ k -combinations and sort them: $\{A1, \dots, Ap\}$;
Bob creates $q = C_k^n$ k -combinations and sort them: $\{B1, \dots, Bq\}$;
 - Alice applies hash function to obtain: $\{h(A1), \dots, h(Ap)\}$;
Bob applies hash function to obtain: $\{h(B1), \dots, h(Bq)\}$;
 - Alice encrypts hashed k -combinations: $\{Er(h(A1)), \dots, Er(h(Ap))\}$;
Bob encrypts hashed k -combinations: $\{Es(h(B1)), \dots, Es(h(Bq))\}$
- Diagram illustrating the exchange of encrypted data:
- A horizontal arrow points from Alice to Bob, labeled $\{Er(h(A1)), \dots, Er(h(Ap))\}$.
 - A horizontal arrow points from Bob to Alice, labeled $\{Es(h(B1)), \dots, Es(h(Bq))\}, Es(\{Er(h(A1)), \dots, Er(h(Ap))\})$.
- Alice encrypts $\{Es(h(B1)), \dots, Es(h(Bq))\}$ with Er , and compares $Er(Es(h(B)))$ and $Es(Er(h(A)))$
 - If the intersection size is greater than 1, group similarity is greater than θ

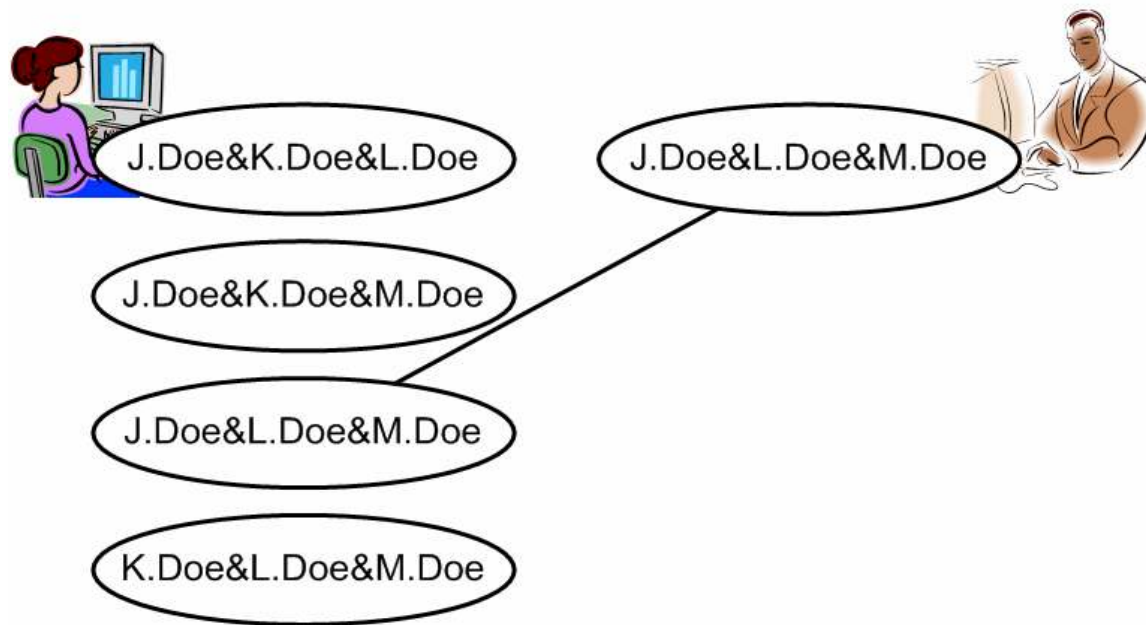
Result: Alice and Bob learn $|\mathbf{R}|$, $|\mathbf{S}|$, and if group similarity $> \theta$

K-Combination Approach Example



K-Combination Approach for TPPGL-E

- Problem?
 - Computation!



Homomorphic Encryption Approach for TPPGL-E

Input: Alice's group $\mathbf{R}=\{r_1,\dots,r_m\}$, Bob's group $\mathbf{S}=\{s_1,\dots,s_n\}$, and a pre-negotiate similarity threshold θ

Protocol:



← Bob group size n →
 ← k ; keys for homomorphic encryption →



- Alice constructs $R(x) = \prod (x - r_i)$ and computes coefficients α_u that $R(x) = \sum_{u=0}^{u=m} \alpha_u x^u$

← Alice encrypts the coefficients $\{E(\alpha_0), \dots, E(\alpha_m)\}$ and send to Bob →

- For each s_j , Bob evaluates the polynomial to get $\mathbf{Enc}(R(s_j))$, without decryption
- Bob chooses a random value γ , and a pre-set special value v . For each $\mathbf{Enc}(R(s_j))$, Bob computes $\mathbf{Enc}(\gamma \times R(s_j) + v)$.
- Bob chooses a random number kb , and injects kb number of $\mathbf{Enc}(v)$ into the set. Meanwhile, Bob also injects random number of random values into this set.

← Bob permutes the polluted set of $\mathbf{Enc}(\gamma \times R(s_j) + v)$ →

- Alice decrypts all items, and counts the number of v values: $kb + |\mathbf{R} \cap \mathbf{S}|$

← $\mathbf{Enc}(kb + |\mathbf{R} \cap \mathbf{S}|)$ →

- Bob calculates $\mathbf{Enc}(kb + |\mathbf{R} \cap \mathbf{S}|) - \mathbf{Enc}(kb + k) = \mathbf{Enc}(|\mathbf{R} \cap \mathbf{S}| - k)$, and then creates random number $\gamma' \ll N$, and $v' < \gamma'$

← $\mathbf{Enc}(\gamma' \times (|\mathbf{R} \cap \mathbf{S}| - k) + v')$ →

- Alice decrypts $m = \gamma' \times (|\mathbf{R} \cap \mathbf{S}| - k) + v'$, and output "YES" if $m < N/2$, or "NO" if $m > N/2$

Result: Alice and Bob learn $|\mathbf{R}|$, $|\mathbf{S}|$, and if group similarity $> \theta$

Group Linkage with Approximate Matching

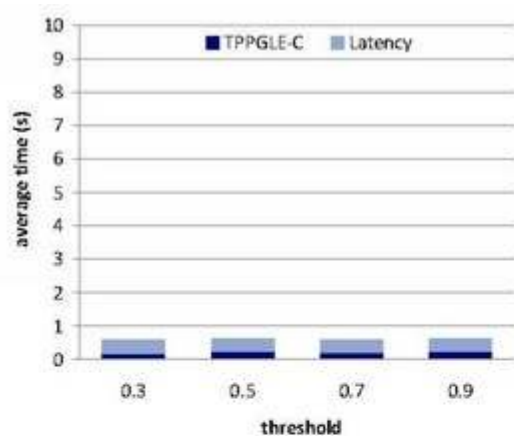
- Alice holds a group of records
- Bob holds a group of records
- Record level similarity: inner product with cut-off
- Group level similarity:

$$\begin{aligned}\mathbf{SIM}(R, S) &= \text{BMsim},\rho (R, S) \\ &= \min(m', n') / (|R| + |S| + \min(m', n'))\end{aligned}$$

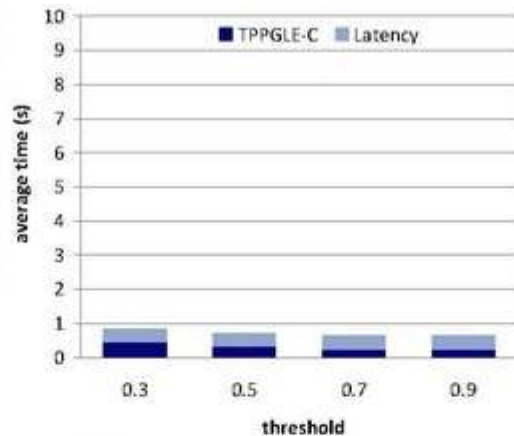
Experiment Results

- Three real data sets [Tang et. al., KDD 2009]
 - AN: a co-author network with 640,134 authors and 1,554,643 co-author relationships
 - CN: a paper citation network of 2,329,760 papers and 12,710,347 citations
 - MN: a movie network with 142,426 relationships
 - Generate synthetic groups
- Evaluate *end-to-end execution time* with varying *group-size* (with 5, 10, 15 records per group) and *threshold* $\theta(\in \{0.3, 0.5, 0.7, 0.9\})$

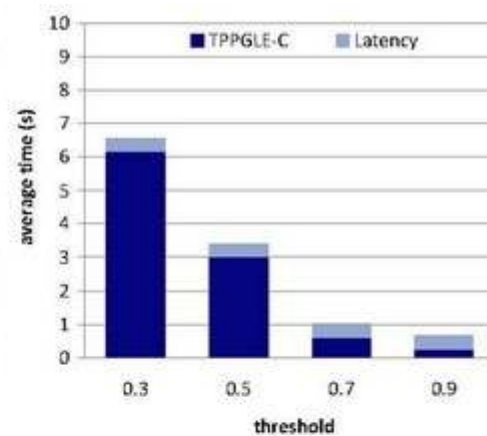
Average End-to-End Execution Time



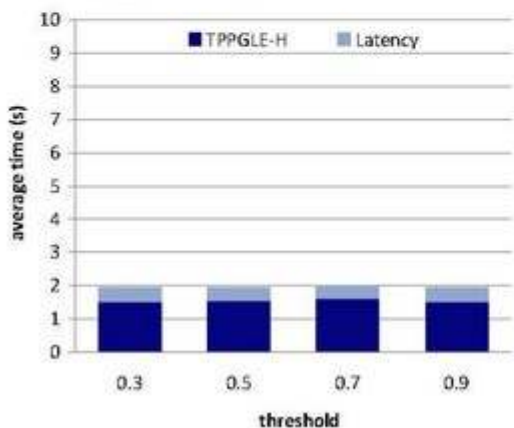
(a) Group size = 5.



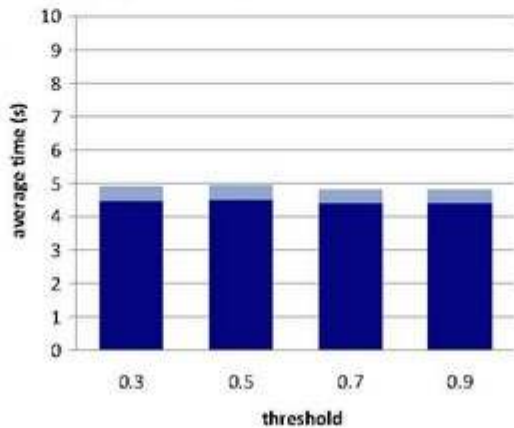
(b) Group size = 10.



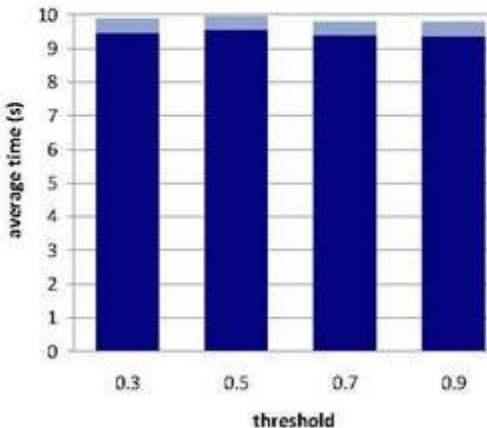
(c) Group size = 15.



(d) Group size = 5.



(e) Group size = 10.



(f) Group size = 15.



Thank You!

Questions?

