

# Paljastavad küsimused ehk ligikaudne otsing

Sven Laur  
swen@math.ut.ee

Tartu Ülikool

# Road-map

- What is approximate matching?
- What can we do with approximate matching?
- There is no easy solution to secure approximate matching.
- Secure randomised affine transformation.
- More efficient settings
- There is no easy solution.
- Conclusion

# Motivation

How to sell the same thing twice? — Start a service!

- But collecting and systematising the information is expensive.
- And some services violate privacy.

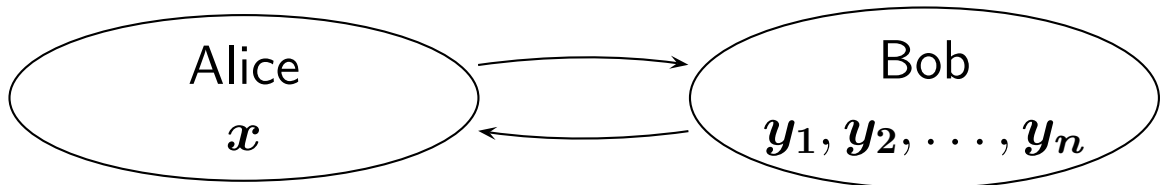
Usually privacy is not a main concern, except

- if queries reveal delicate information;
- if a leakage causes explicit economical expenses.

Onion-routing does not guarantee privacy, since

- a query itself can reveal the personality;
- a query itself can contain useful information.

# Approximate matching



$$d_0 = \min_i d(\mathbf{x}, \mathbf{y}_i)$$
$$i_0 = \operatorname{argmin}_i d(\mathbf{x}, \mathbf{y}_i)$$

## Security considerations

- Bob should learn nothing about query vector  $\mathbf{x}$ .
- Alice should learn nothing about the database, except the match distance  $d_0$  and the match number  $i_0$ .

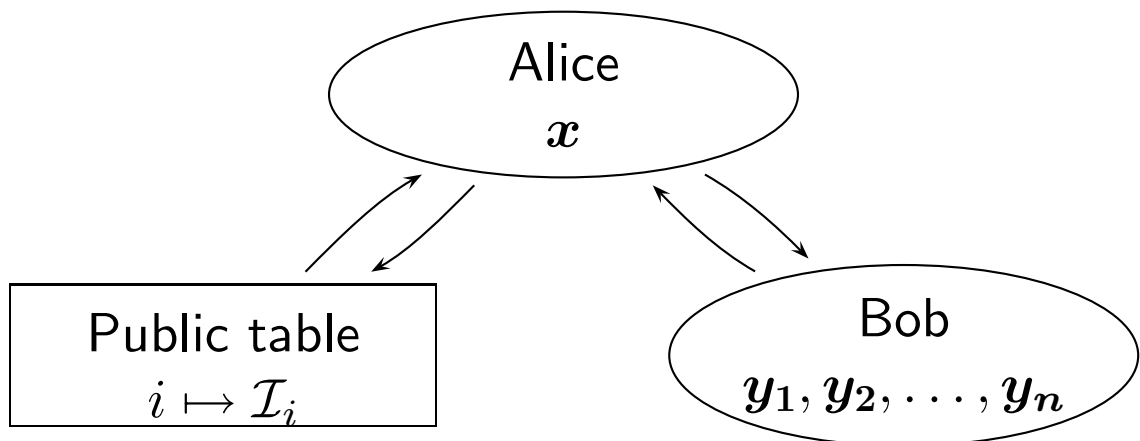
## Efficiency considerations

- Communication complexity should be **poly**( $\log n$ ).
- Computational complexity should be **poly**( $n$ ).

## Example application: Symptom-action type databases

Bob has accumulated knowledge in the following form.

- A symptom  $y_i$  and an appropriate action  $\mathcal{I}_i$ .
- Nearest neighbourhood search gives appropriate action for an unknown  $x$ .



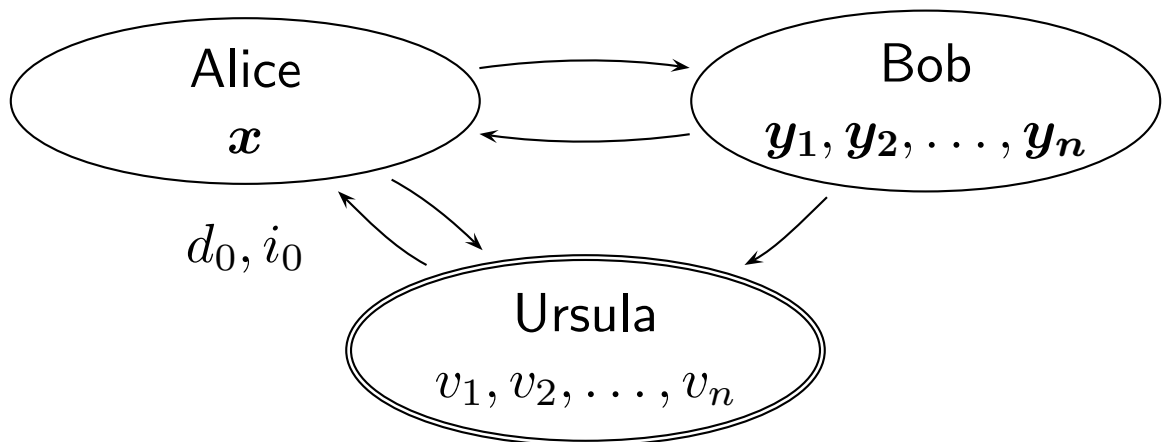
- The protocol does not require private information retrieval!
- An approximate matching might be more efficient!

## Further analysis

- Approximate matching requires a secure evaluation of minimum

$$d_0 = \min_i d(x, y_i).$$

- Currently, no efficient and cryptographically secure minimum finding protocols are known.
- To bypass the problem, we include a trusted third party Ursula.



- Naive implementations yield communication complexity  $\Theta(n)$ .

## From distance to scalar product

For the Euclidean distance, we must calculate

$$d_i = (\mathbf{x} - \mathbf{y}_i)^2 = \mathbf{x}^2 - 2\mathbf{x} \cdot \mathbf{y}_i + \mathbf{y}_i^2.$$

Transformation

$$\mathbf{x} = (x_1, \dots, x_m) \mapsto \mathbf{x}' = (-2x_1, \dots, -2x_n, 1),$$

$$\mathbf{y} = (y_1, \dots, y_m) \mapsto \mathbf{y}' = (y_1, \dots, y_n, y_1^2 + \dots + y_m^2)$$

gives

$$\mathbf{x}' \cdot \mathbf{y}'_i = -2\mathbf{x} \cdot \mathbf{y}_i + \mathbf{y}_i^2$$

and thus

$$d_0 = \mathbf{x}^2 + \min_i \mathbf{x}' \cdot \mathbf{y}'_i,$$

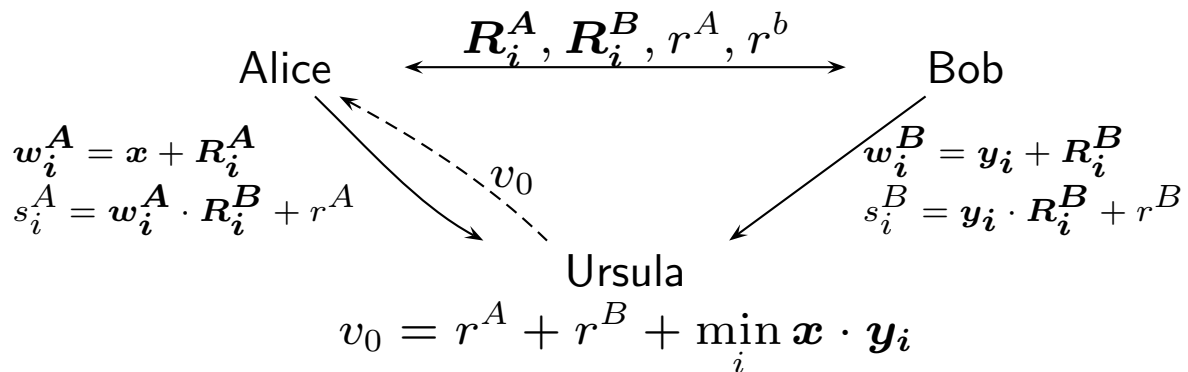
$$i_0 = \operatorname{argmin}_i \mathbf{x}' \cdot \mathbf{y}'_i.$$

# MinDASP protocol (Du and Atallah)

The protocol hinges on the fact that

$$\begin{aligned}
 & (\mathbf{x} + R_i^A) \cdot (\mathbf{y} + R_i^B) \\
 &= \mathbf{x} \cdot \mathbf{y}_i + \underbrace{R_i^A \cdot \mathbf{y}_i}_{s_i^B - r^B} + \underbrace{(\mathbf{x} + R_i^A) \cdot R_i^B}_{s_i^A - r^A}
 \end{aligned}$$

Therefore, we get a straightforward protocol



- Ursula gets additive shares  $v_i = \mathbf{x} \cdot \mathbf{y}_i - r^A - r^B$ .
- For one database element all values are perfectly masked, but this is not true for many items.



## The attack

Note that

$$s_i^B = R_i^A \cdot w_i^A + r^B = (w_i^A - x) \cdot w_i^B + r^B$$

and therefore

$$(w_i^B - w_1^B) \cdot x = w_i^A \cdot w_i^B - w_1^A \cdot w_1^B - (s_i^B - s_i^A).$$

Since Ursula knows everything except  $x$ , she can compose a system of linear equations  $Mx = z$ .

- We calculated the exact probability that  $Mx = z$  has a unique solution. The probability is too big.
- We offered two bug-fixes: a slight modification of the protocol and a secure scalar product protocol via homomorphic encryption.
- The latter reduces communication complexity more than four times, however the computational complexity rises.

## General result

It is unreasonable to assume that Ursula knows nothing about the database.

- Some database elements might be publicly known.
- During the longterm use of database, some vectors might leak.

The protocol should remain secure if less than  $\tau = \Theta(n)$  database elements are known to Ursula.

**Theorem 1.** *All protocols, where Ursula obtains additive shares  $v_i = \mathbf{x} \cdot \mathbf{y}_i + r$ , are insecure regardless of used sub-protocols.*

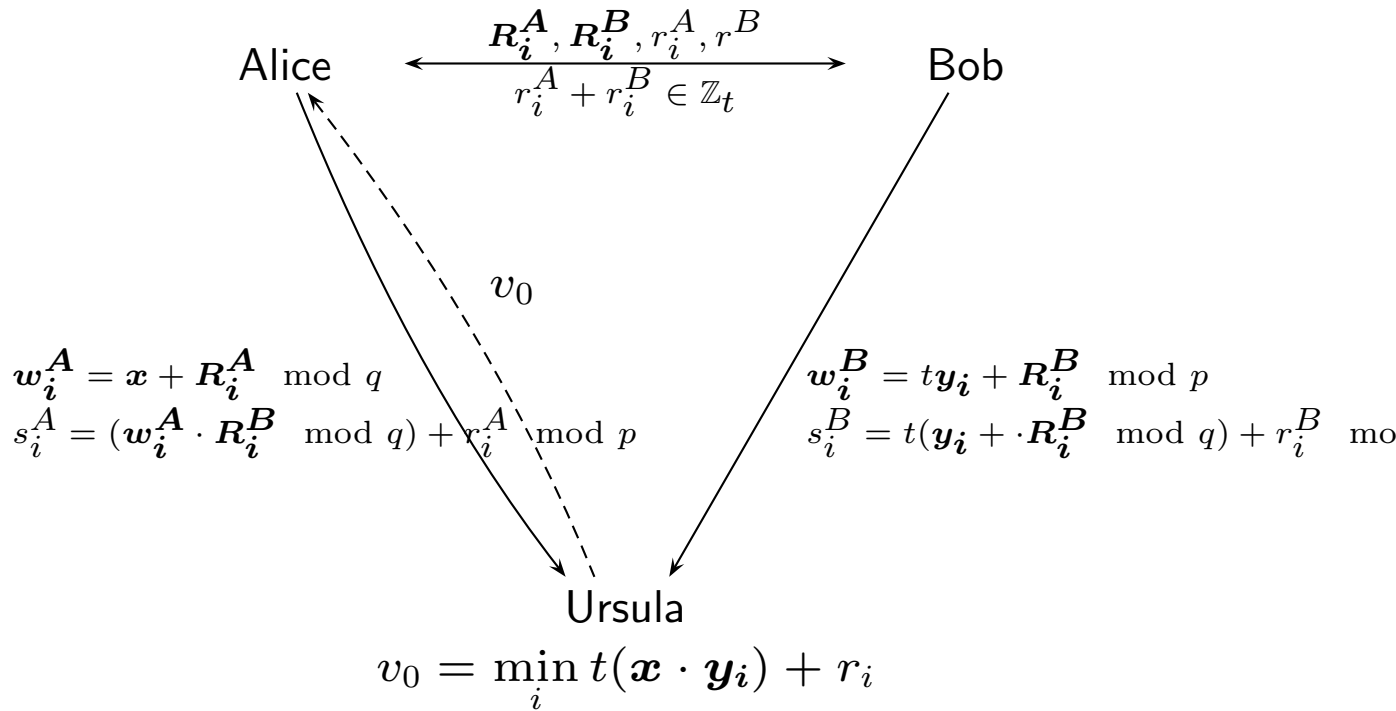
- If Ursula knows database vectors  $\mathbf{y}_1, \dots, \mathbf{y}_k$ , then she can construct a system of linear equations.
- Under the random database assumption the system has unique solutions with high enough probability.
- Otherwise the security depends on the specific database!

# The MinAffineSP protocol: The working principle

- Public parameters are integers  $p, q, t_{max}, t_{min}$ , so that  $t_{min} \approx 2^{160}$  and  $q(t_{max} + 1) < p$ .
- Alice and Bob jointly fix a random multiplier  $t \in [t_{min}, t_{max}]$ .
- Ursula obtains shares  $v_i = t(\mathbf{x} \cdot \mathbf{y}_i \bmod q) + r_i$ , where  $r_i \in \mathbb{Z}_t$ , that are randomly permuted.
- The smallest share  $v_0$  still corresponds to the minimal scalar product.
- Alice can eliminate the randomness

$$\min_i \mathbf{x} \cdot \mathbf{y}_i = \left\lfloor \frac{v_0}{t} \right\rfloor.$$

# The MinAffineSP protocol: The implementation



After some tedious calculations

$$v_i = (w_i^A \cdot w_i^B \pmod q) - s_i^A - s_i^B$$

$$= t(\mathbf{x} \cdot \mathbf{y}_i \pmod q) + r_i \pmod p.$$

# The MinAffineSP protocol: The preliminary security analysis

To break the protocol, Ursula must solve the system of equations

$$tx_i + r_i = z_i, \quad t \in [t_{min}, t_{max}],$$

where  $x_i \in \mathbb{Z}_q, r_i \in \mathbb{Z}_t$  are unknown.

- The values  $r_i$  are uniformly distributed for correct  $t$  and not generally uniformly distributed for  $t' \neq t$ .
- Small differences of  $v_i - v_j = \Delta x_{i,j}t + r_i - r_j$  can suggest the values for  $t$ . If  $\Delta x_{i,j} = 1$  then we have a slight probability peak at  $t$ .
- Linear combinations  $a_1v_1 + \dots + a_mv_m$  with small coefficients of  $a_i$  are more restrictive, since the true random term  $a_1r_1 + \dots + a_mr_m$  converges to the normal distribution.
- But the probabilistic reduction increases the uncertainty by  $m^{1/2}$  times.

# More efficient settings: Secure storage outsourcing problem

Consider a scenario

- First Alice outsources the database to Bob.
- Afterwards Alice needs match distances

$$\min_i (\mathbf{x}_j - \mathbf{y}_i)^2.$$

Security considerations

- Bob should learn nothing about database vectors  $\mathbf{y}_1, \dots, \mathbf{y}_n$  and query vectors  $\mathbf{x}_1, \dots, \mathbf{x}_k$ .
- The protocol should remain secure even if Bob knows  $\tau = \Theta(n)$  vectors  $\mathbf{x}_j$  or  $\mathbf{y}_i$ .

Efficiency considerations

- The protocol should have only  $\mathcal{O}(1)$  rounds.

## General result

Proposed solutions use additive sharing, that is Bob can calculate  $s_{ij} = (\mathbf{x}_j - \mathbf{y}_i)^2 + r_j$ .

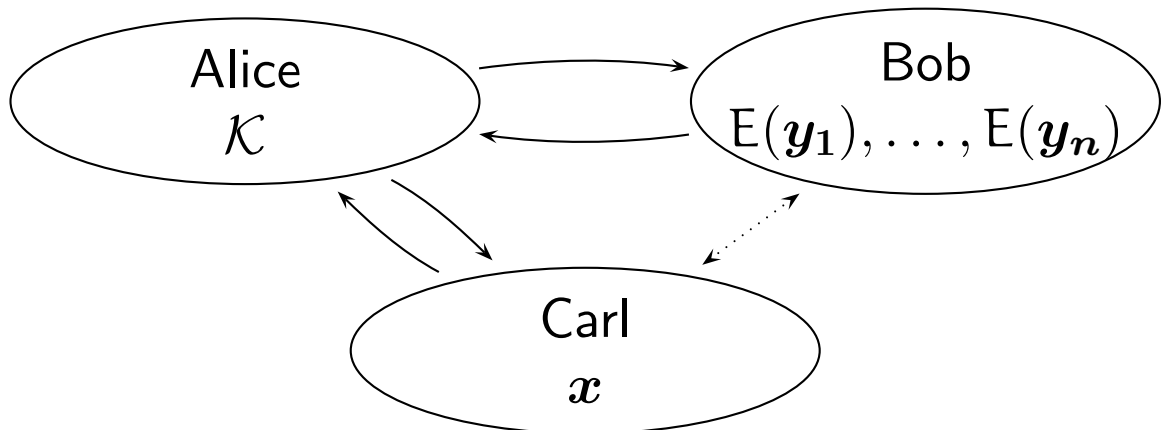
- The second security goal—a resistance against leaking vectors—is not satisfied.
- The difference matrix  $\Delta S = (\Delta s_{ij})$ , where

$$\Delta s_{ij} = s_{ij} - s_{i1} - s_{j1} + s_{11} = -2\Delta \mathbf{x}_j \Delta \mathbf{y}_i,$$

reveals linear information.

- If the database does not belong to the hyper-plane then the matrix columns reveal linear dependencies between query differences  $\Delta \mathbf{x}_j = \mathbf{x}_j - \mathbf{x}_0$ .
- The matrix rows reveal linear dependencies between the differences  $\Delta \mathbf{y}_i = \mathbf{y}_i - \mathbf{y}_1$ .

## More efficient settings: Secure storage and computing outsourcing problem



- First Alice outsources the database to Bob.
- Afterwards Carl composes queries with the help of Alice.
- Alice helps Carl to decode the answers.



# Practical requirements

## Security considerations

- Bob should learn nothing about the database vectors  $\mathbf{y}_1, \dots, \mathbf{y}_n$  and query vectors  $\mathbf{x}_1, \dots, \mathbf{x}_k$ .
- Alice should learn nothing about the query vectors.
- Collusion between Bob and Carl is allowed
- The protocol should remain secure even if Bob knows  $\tau = \Theta(n)$  vectors  $\mathbf{y}_i$ .

## Efficiency considerations

- The protocol should have only  $\mathcal{O}(1)$  rounds.

## General result

**Theorem 2.** *The protocols, where Bob obtains linear shares  $s_{ij} = \alpha_j(\mathbf{x}_j - \mathbf{y}_i)^2 + r_j$ , are insecure regardless of used sub-protocols.*

- The second security goal—resistance against leaking vectors—is not satisfied.
- Carl and Bob can restore originals of database elements by gradient search.
- More subtle attacks are possible.

## Conclusions: Curse of linearity

- It is hard to find the minimum when shares do not have the linear form  $\alpha x + r$ .
- But the linearity opens door to relatively simple attacks based on linear algebra.
- The linear transformation is not safe a way to hide data.