

# Private Synopses-Driven Data Integration

Eros Fabrici @ Athena Research Center & UPC,  
eros.fabrici@athenarc.gr



## Overview

- Data Integration (DI) is a long-standing research topic
- Privacy in Data Integration is a more recent topic, where the focus is to make the process private across the different parties involved.
- The literature focused in Record Linkage.
- The objective of this PhD is to explore the usage of *sketches* and *differential privacy* for inferring the relationships between the data in a data-driven fashion to be applied on DI process.

## Privacy-Preserving Data Integration

The privacy threats in the DI lie in the first two phases. In particular, there is very limited research in privacy-preserving instance-based schema alignment. Privacy-Preserving Record Linkage (PPRL) is the most studied problem. Table 1 shows a categorization of the approaches proposed in the literature for PPRL.

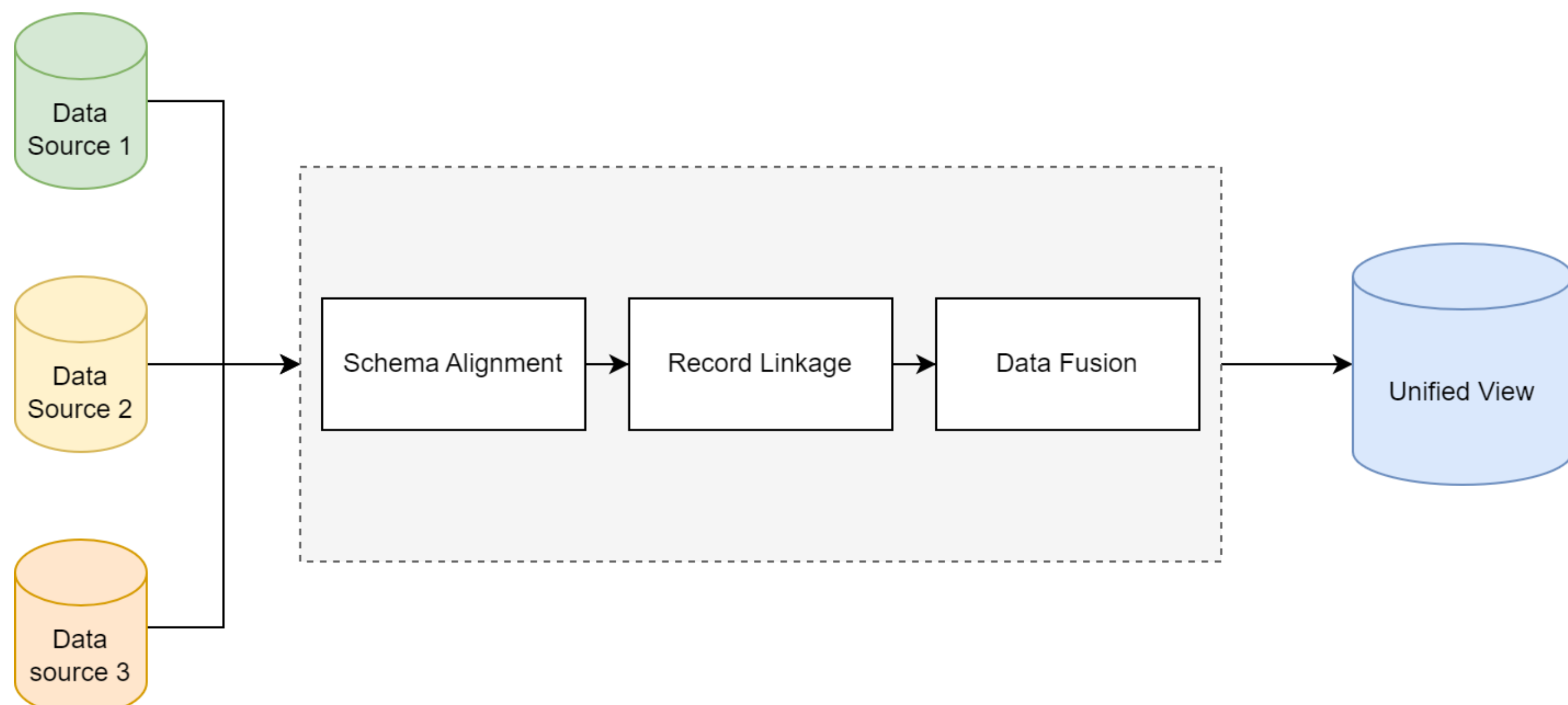


Figure 1: Data Integration Phases

## Differential Privacy

Differential Privacy is a technique for publicly sharing information about a dataset by describing the patterns of groups within the dataset while withholding information about individuals in the dataset. It is *de facto* standard in the field of privacy for ML, synthetic data and querying.

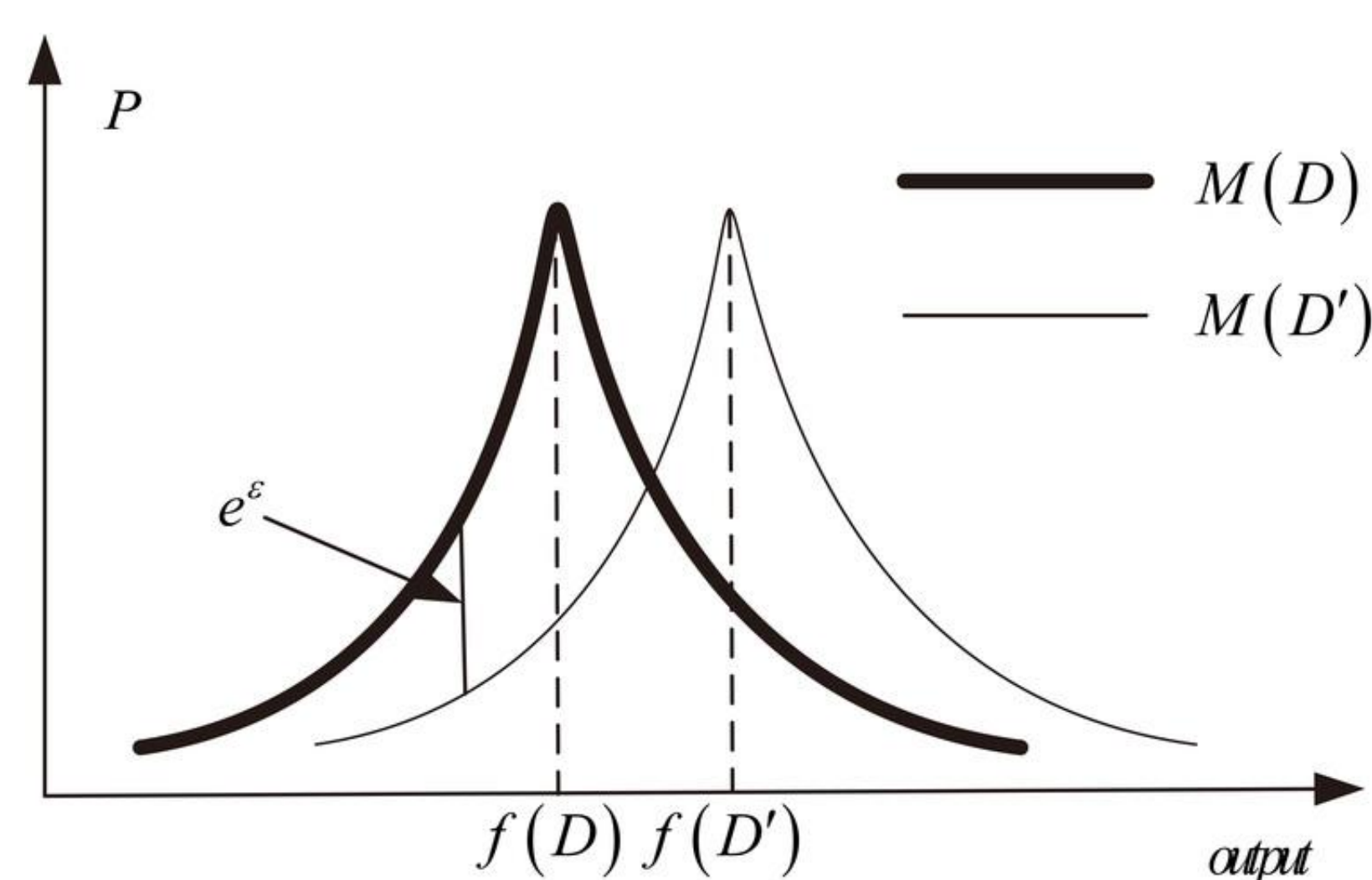


Table 1: Comparison of Private Record Linkage algorithms

	Paper	Privacy Guarantee		3 <sup>rd</sup> P. <sup>1</sup>	SA <sup>2</sup>	Matching function
		Blocking/Filtering	Matching			
Formal Privacy	[5, 32]	n.a.	DP <sup>3</sup> (RR <sup>4</sup> and Embedding)	✓	✗	Dice, E. <sup>5</sup>
	[2, 3, 12, 21]	n.a.	SMC <sup>6</sup> and Cryptography	✓	✗	TFIDF, Any, EM <sup>7</sup>
	[14, 17, 18, 27]	DP, <i>k</i> -anonymity	SMC	✓	✓	Dist. Based
Ad-Hoc Privacy	[10, 20, 31]	n.a.	Hashing (Phonetic, BF <sup>9</sup> )	✓	✗	EM, Dice
	[30, 34]	n.a.	Embedding (Complex P. <sup>8</sup> , SparseMap [16])	✓	✓	Dist. Based

<sup>1</sup> 3<sup>rd</sup> Party; <sup>2</sup> Schema-Aware; <sup>3</sup> Differential Privacy; <sup>4</sup> Randomized Response; <sup>5</sup> Euclidean Distance;  
<sup>6</sup> Secure Multiparty Computation; <sup>7</sup> Exact Matching; <sup>8</sup> Complex Plane; <sup>9</sup> Bloom Filter

## Sketches for estimating joint quantities

Sketches are probabilistic data structures to summarize the data and to obtain estimations on pre-defined queries (e.g. cardinality estimation) with theoretical bounds. The two main families of interests are *order-invariant hash based* for cardinality and joint quantities estimation and *linear sketches* for estimating distributional frequency.

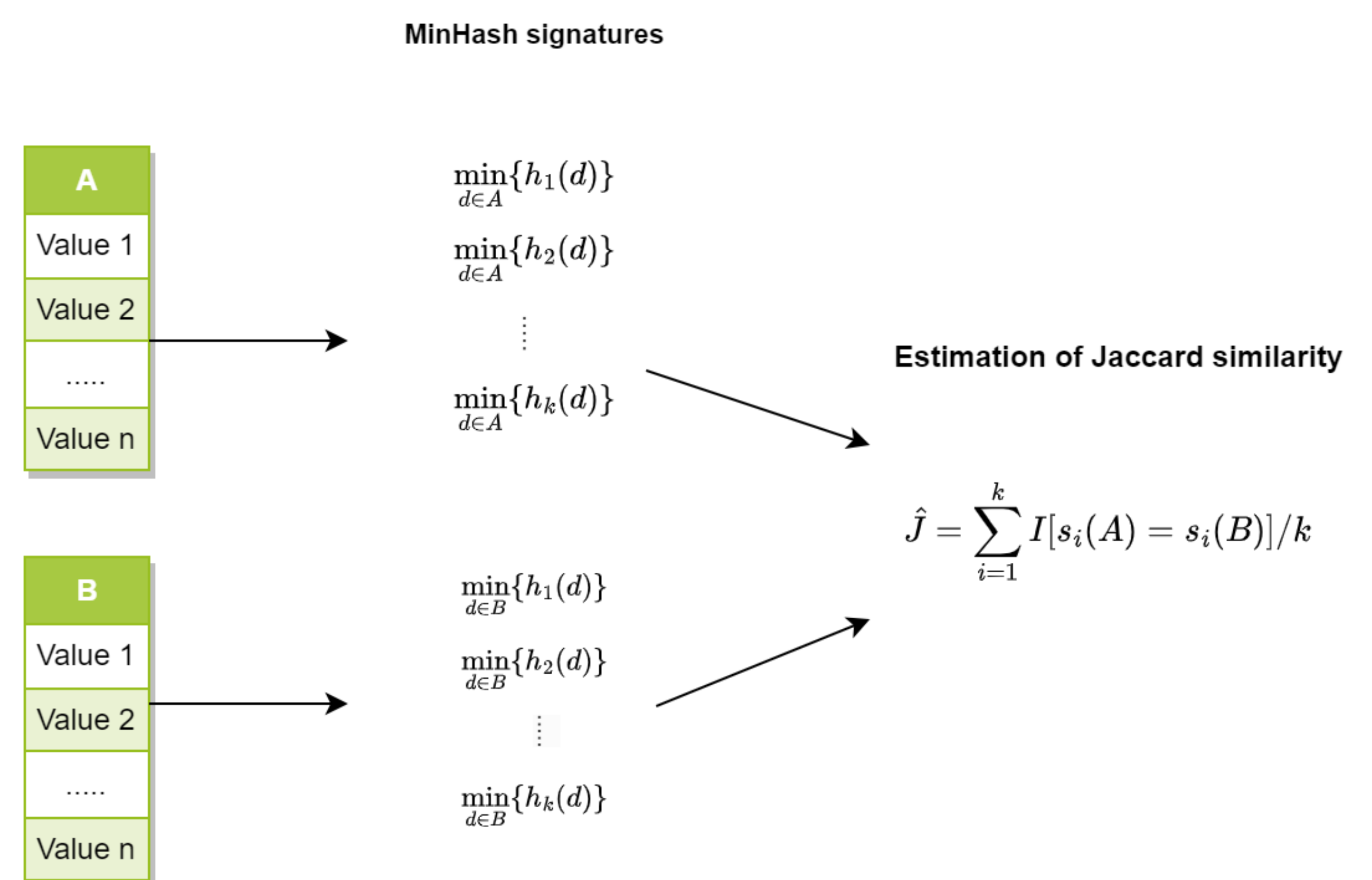


Figure 2: Estimating Jaccard similarity with MinHash

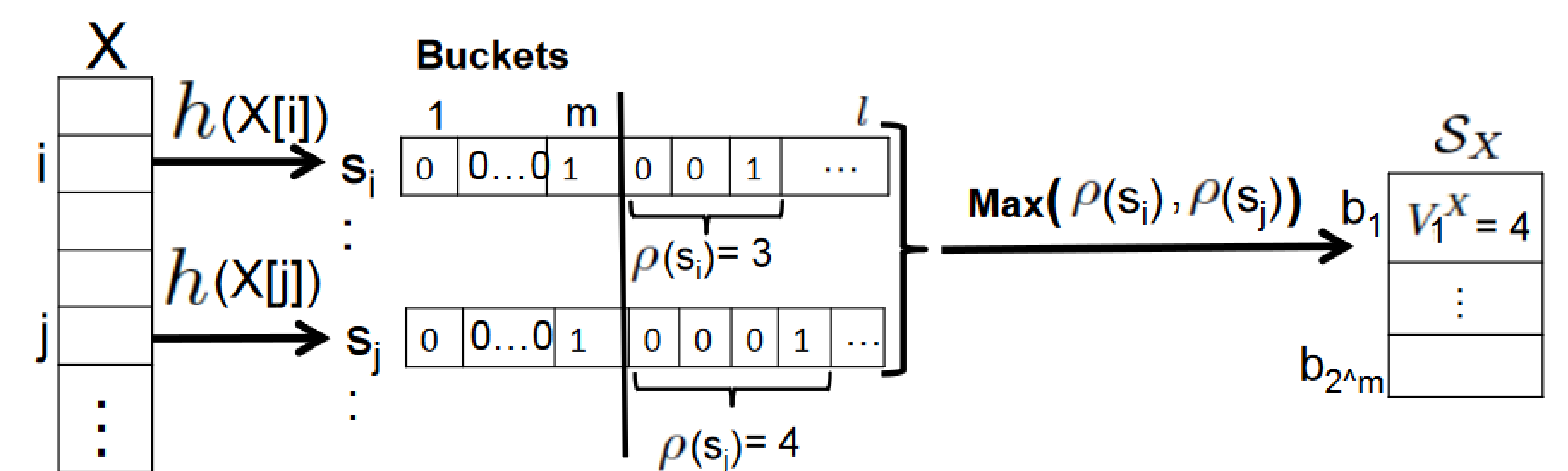


Figure 3: HyperLogLog, for cardinality and joint quantities estimations

## Current Directions

The current direction is to implement a framework for a fully data-driven, private DI, where state-of-the-art sketches will be tested. An empirical study will be carried out to find the best sketches for the purpose. Efficient privacy budget allocation tailored to the DI problem will be studied as well. The aim of the tool is to infer mappings across the datasets' attributes, while guaranteeing privacy.