

Register automata as a model for local computation in distributed query evaluation

Tony Tan
National Taiwan University

Collaboration with: Jonny Daenen, Frank Neven, Frederic Servais, Jan Van den Bussche (Hasselt University), Nicole Schweikardt (Humboldt University, Berlin), Stijn Vansummen (Free University of Brussels)

Popular distributed systems

Systems in which parallelisation is achieved via key-value paradigm.

- Apache Hadoop, which is based on Google's map-reduce.

[Yahoo's Hadoop; Dean and Ghemawat OSDI'04]

- Spark.

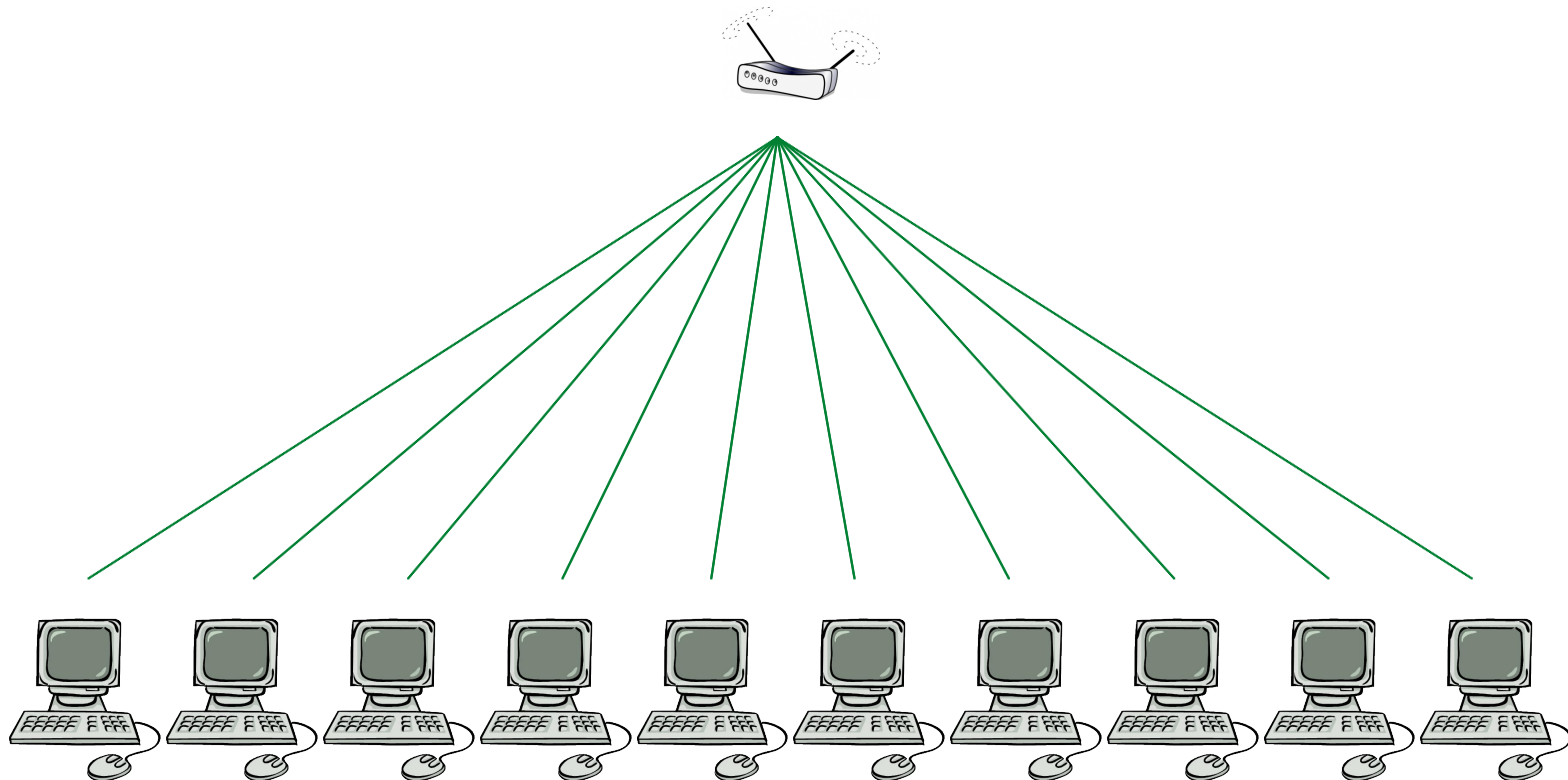
[Zaharia, et. al. NSDI'12; Zaharia's PhD thesis 2014]

Hadoop and Spark

Some reasons for their popularity:

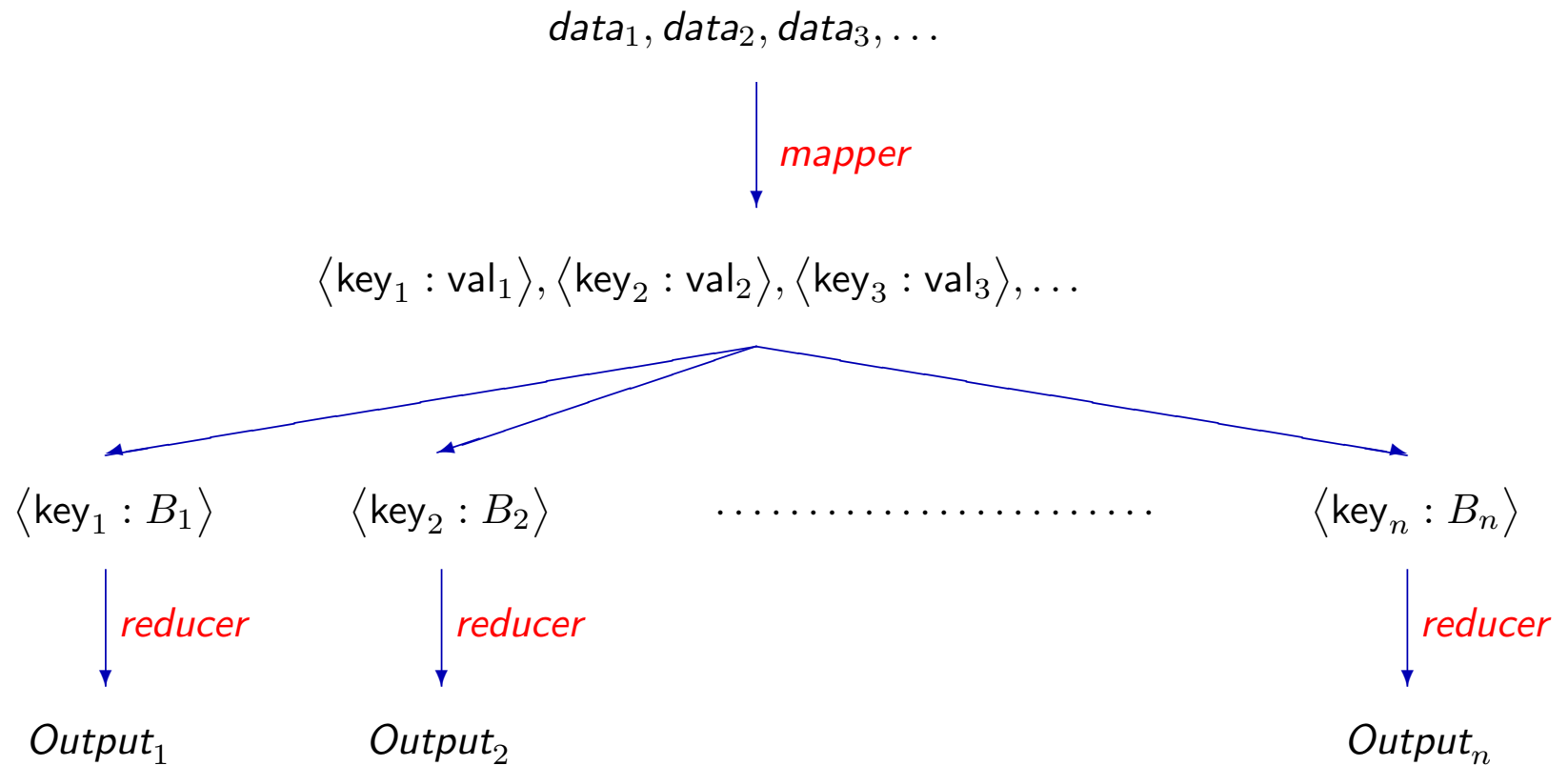
- They are free.
- They take care of the communication.
- Programmers only specify the local computation.
- Do not require special network architecture.

Hadoop and Spark

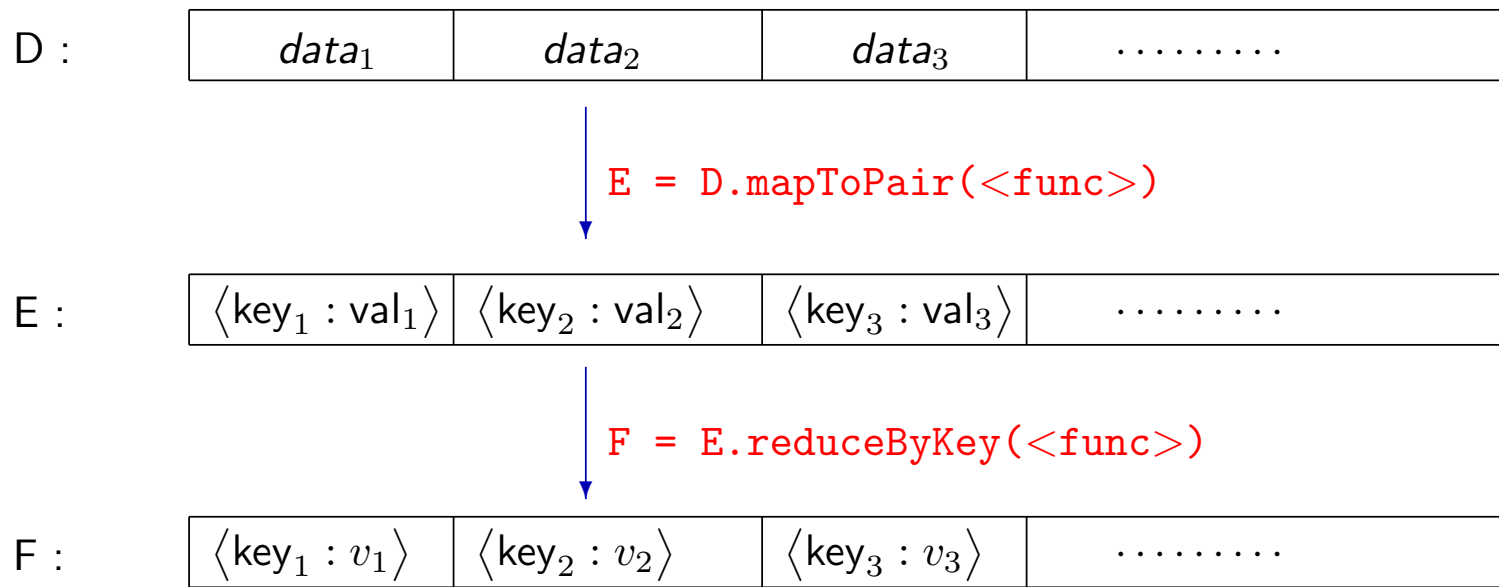


Protocol: Local computation, communication, local computation, communication, . . .

Key-value paradigm in Hadoop



Key-value paradigm in Spark using RDD (Resilient Distributed Datasets)



Brief review on relational algebra (core of SQL)

- *Union*: $A \cup B$.
- *Intersection*: $A \cap B$.
- *Difference*: $A - B$.
- *Projection*: $\pi_{i_1, \dots, i_k}(A)$.
- *Selection*: $\sigma_{i=j}(A)$.
- *Semijoin*: $A \bowtie_{i=j} B$.

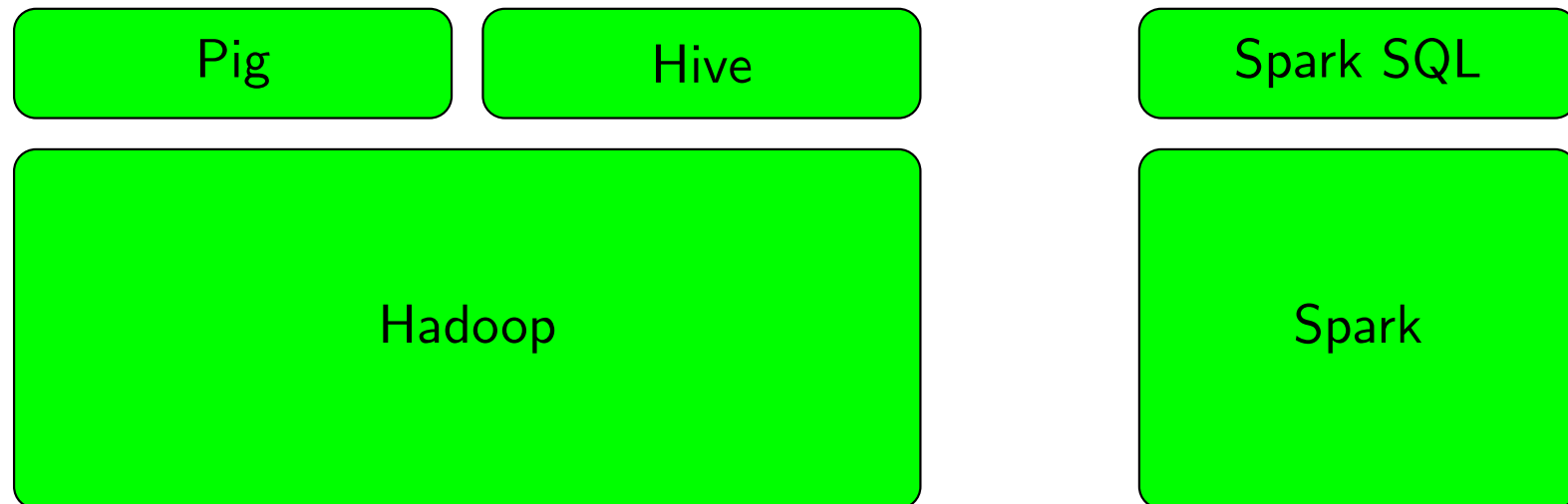
$$A \bowtie_{i=j} B := \left\{ \bar{a} \in A \mid \text{there is } \bar{b} \in B \text{ such that } a_i = b_j \right\}$$

- *Join*: $A \Join_{i=j} B$.

$$A \Join_{i=j} B := \left\{ (\bar{a}, \bar{b}) \in A \times B \mid a_i = b_j \right\}$$

Database engines on top of Hadoop and Spark

- Pig [Gates, et. al. VLDB'09; Olston, et. al. SIGMOD'08]
- Hive [Thusoo, et. al. ICDE'10]
- Shark and Spark SQL [Xin, et. al. SIGMOD'13 & '15]



Queries on Pig/Hive/Spark SQL

In general, users don't need to mention:

- the number of available servers,
- how to distribute and partition the data.

Users write queries as if on a single processor.

Pig/Hive/Spark SQL

Question: What are the necessary and sufficient *local* computation to evaluate relational algebra?

Main theme of the talk

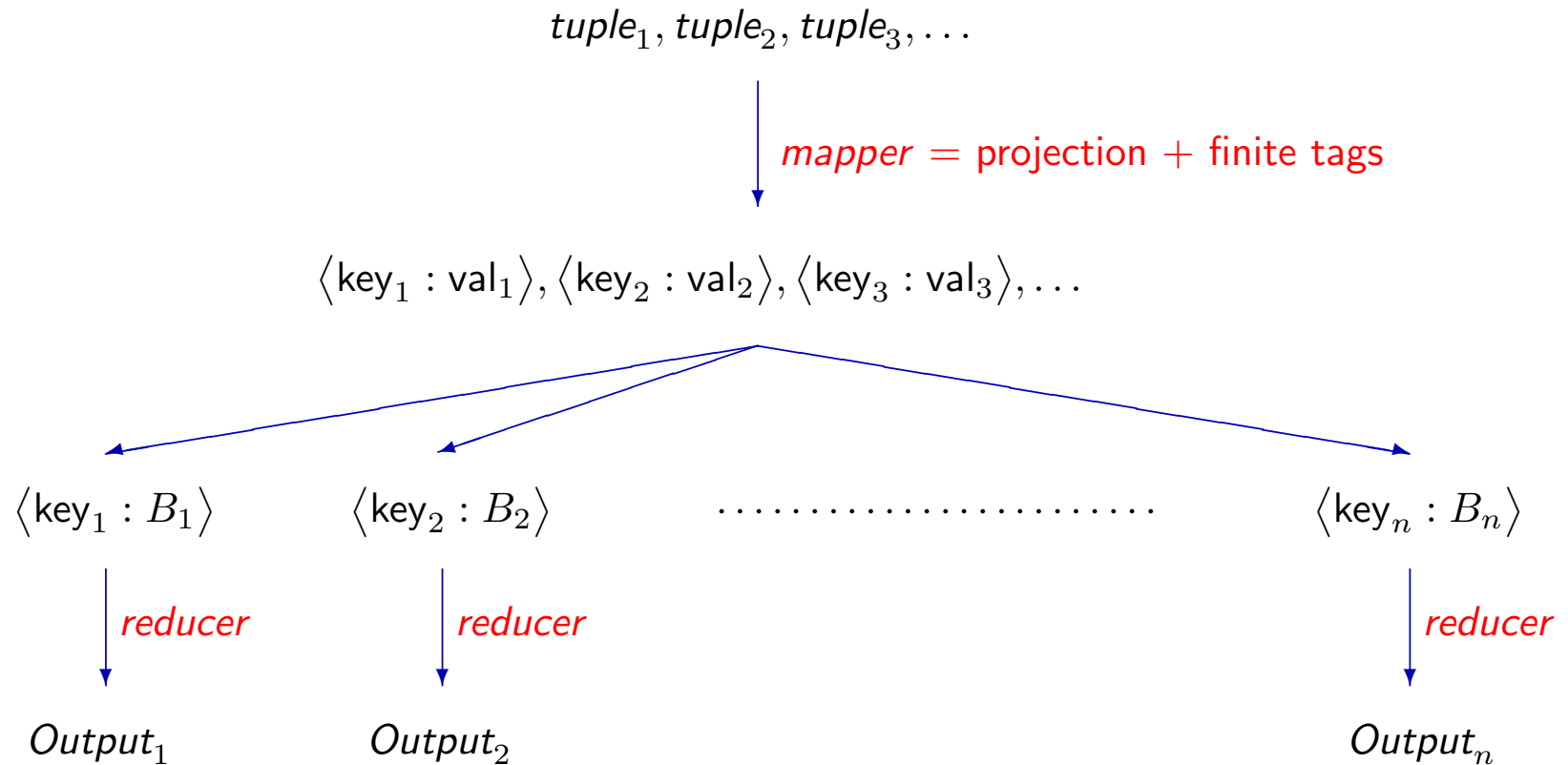
Three models of computation in key-value paradigm.

- DSA (Distributed Streaming with register Automata)
- DST (Distributed Streaming with register Transducers)
- DSTJ (Distributed Streaming with register Transducers and Joins)

The mappers are “generic mappers,” i.e., projection + finite tags.

The reducers are various models with limited computation that process the values in streaming fashion using finite memory.

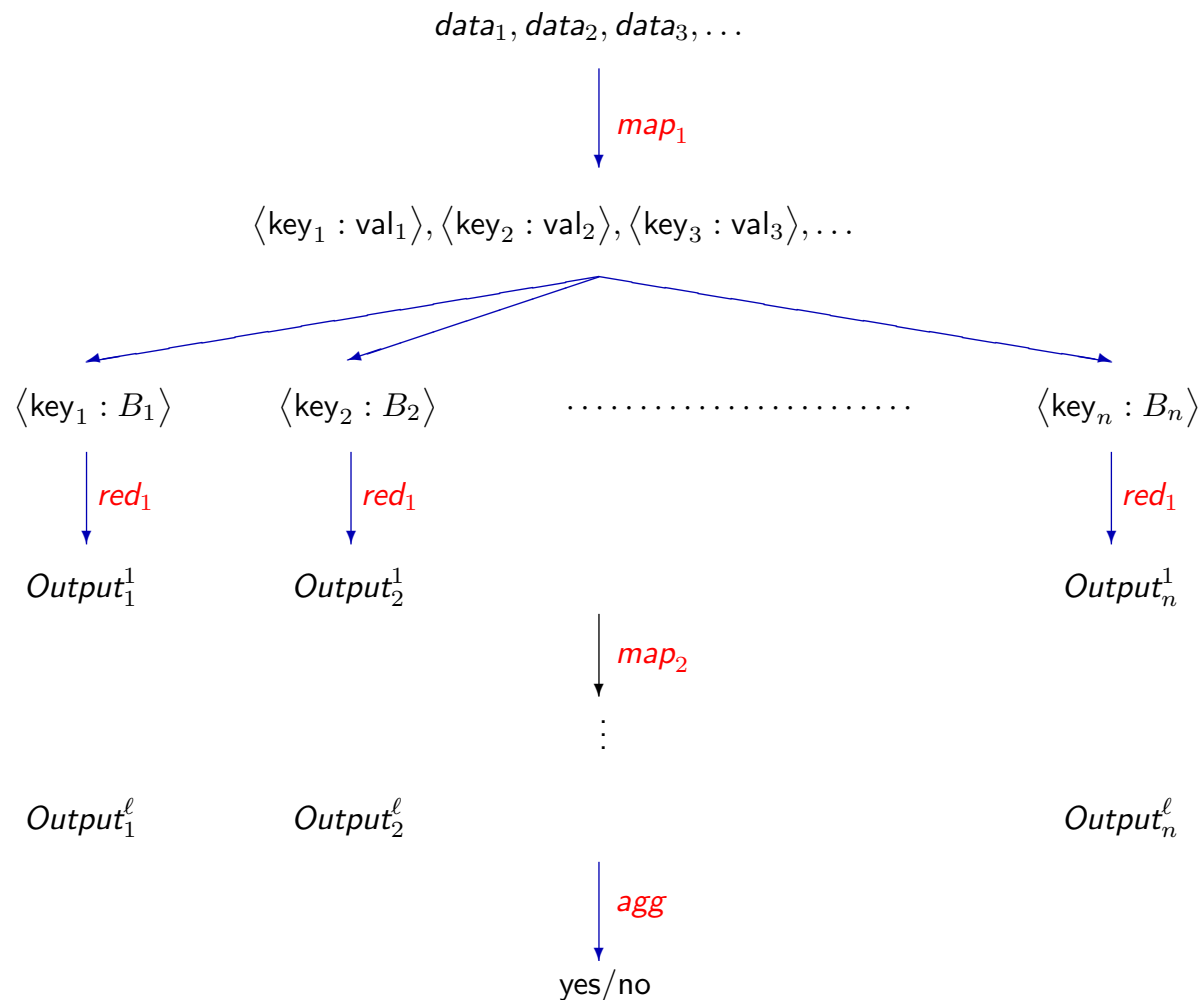
DSA, DST, DSTJ



An example: $R(1, 2) \mapsto \langle S(1) : T(1, 2) \rangle$

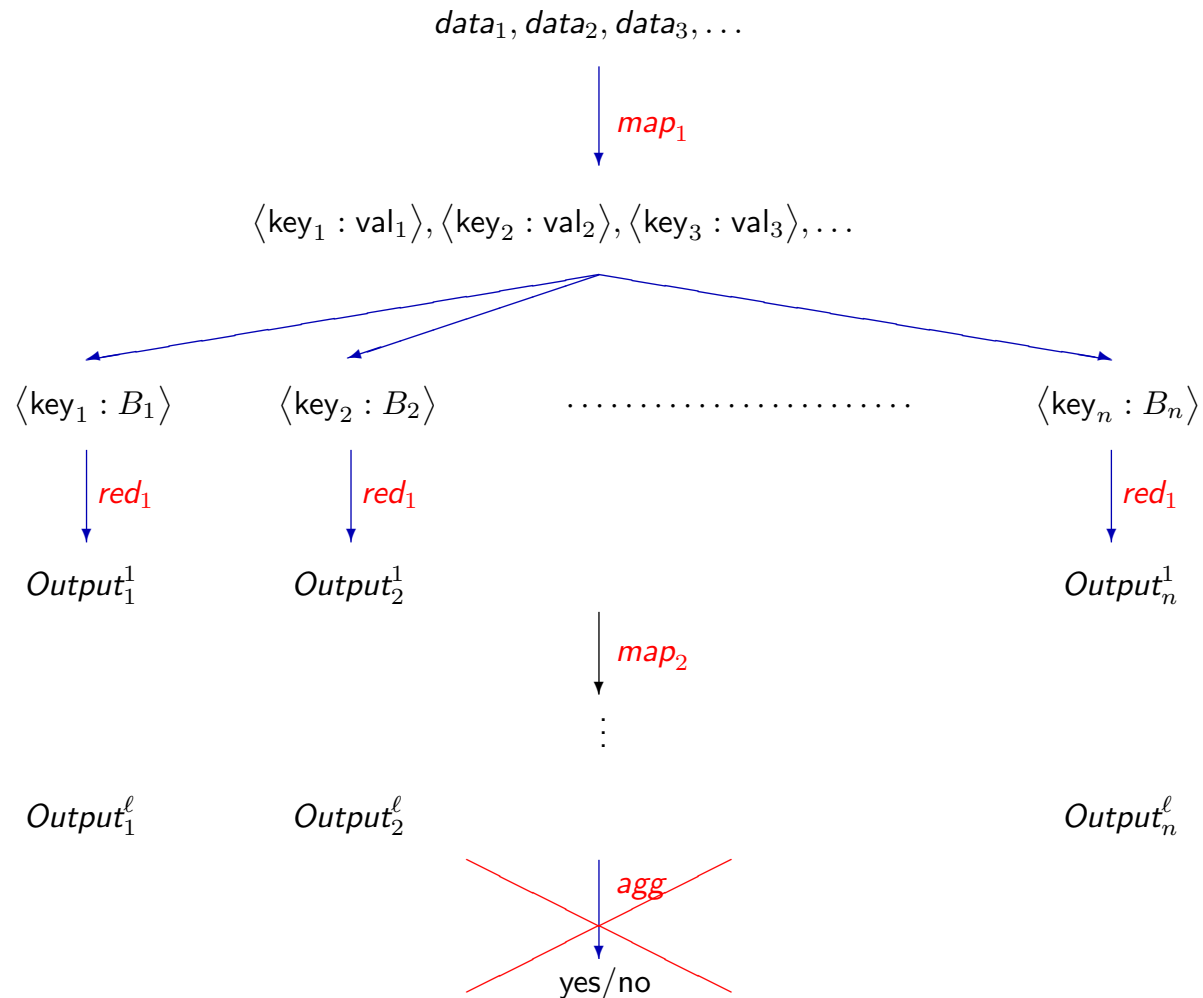
Boolean DSA, DST, DSTJ:

$(map_1, red_1, map_2, red_2, \dots, map_\ell, red_\ell, agg)$



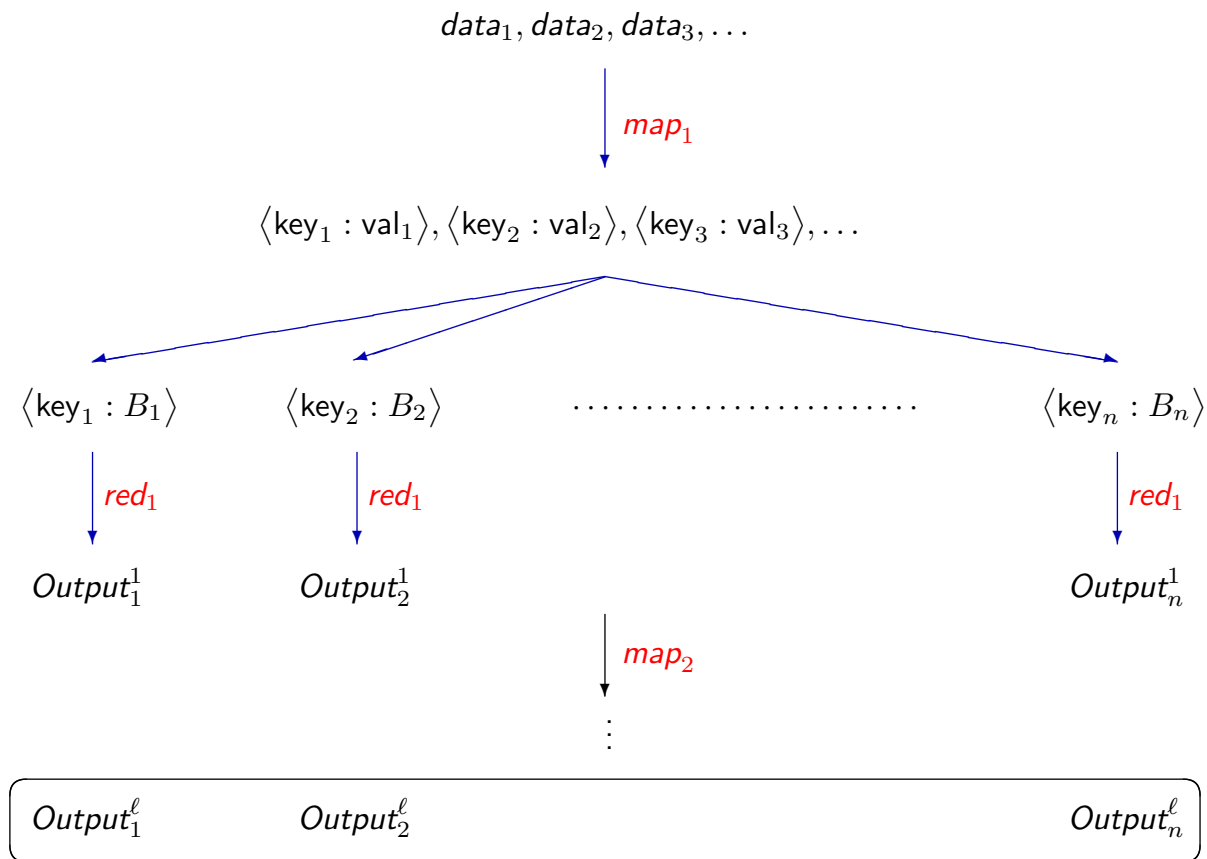
Non-Boolean DSA, DST, DSTJ:

$(map_1, red_1, map_2, red_2, \dots, map_\ell, red_\ell, \text{~~agg~~)$



Non-Boolean DSA, DST, DSTJ:

$(map_1, red_1, map_2, red_2, \dots, map_\ell, red_\ell)$



The big picture for DSA, DST and DSTJ



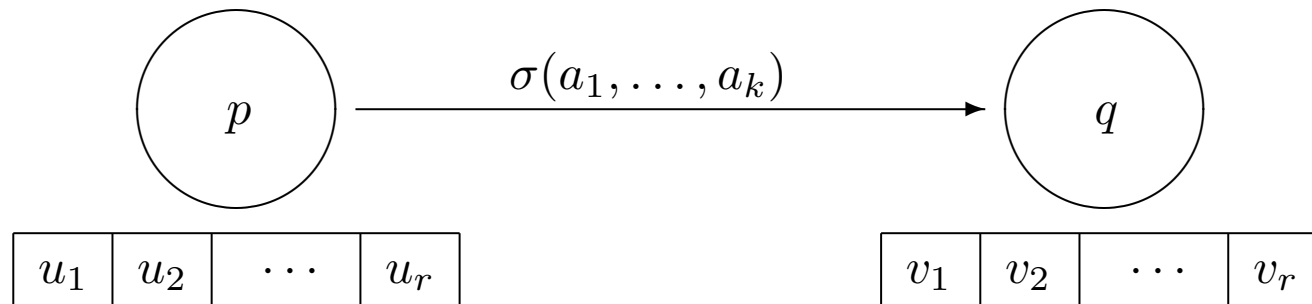
For Boolean case: DSA captures FO over bounded degree databases.

DSA: $(map_1, red_1, map_2, red_2, \dots, map_\ell, red_\ell, agg)$

- map_1, \dots, map_ℓ are generic mappers, i.e., projections + finite tags.
- red_1, \dots, red_ℓ are deterministic register automata “with output.”
- agg is deterministic register automaton.

Register automata

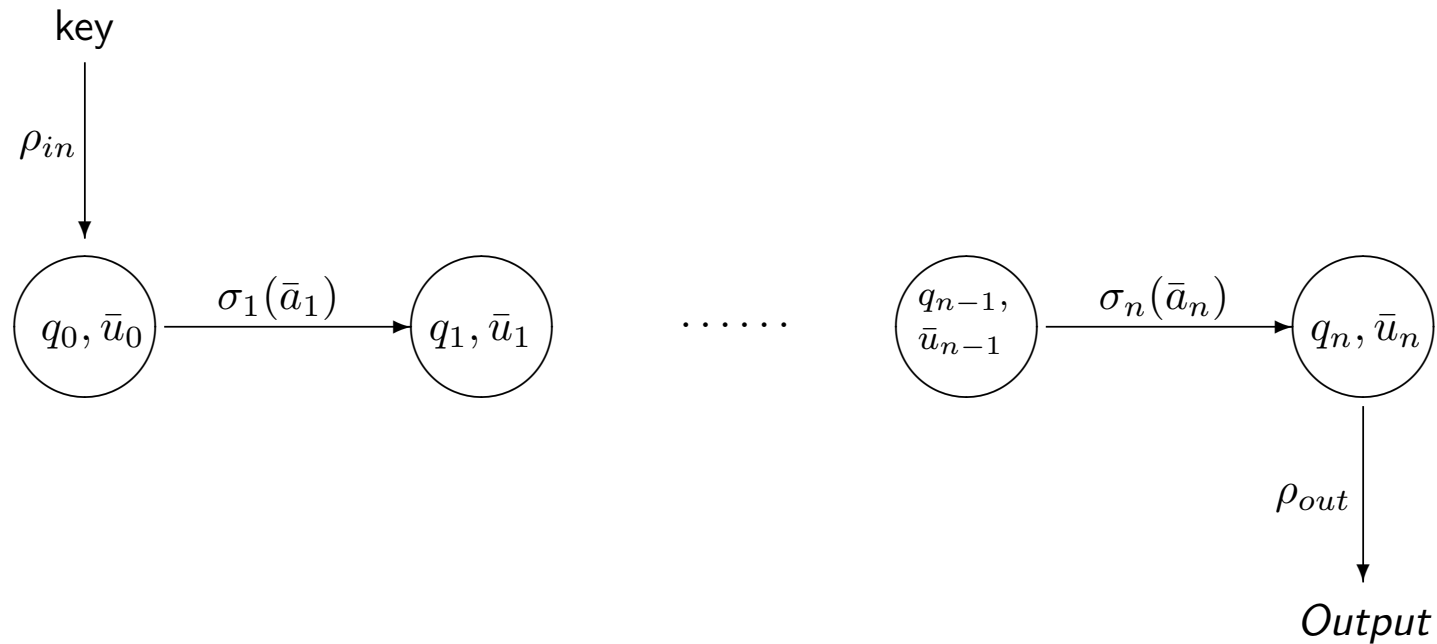
Finite state automaton with a fixed number of registers.



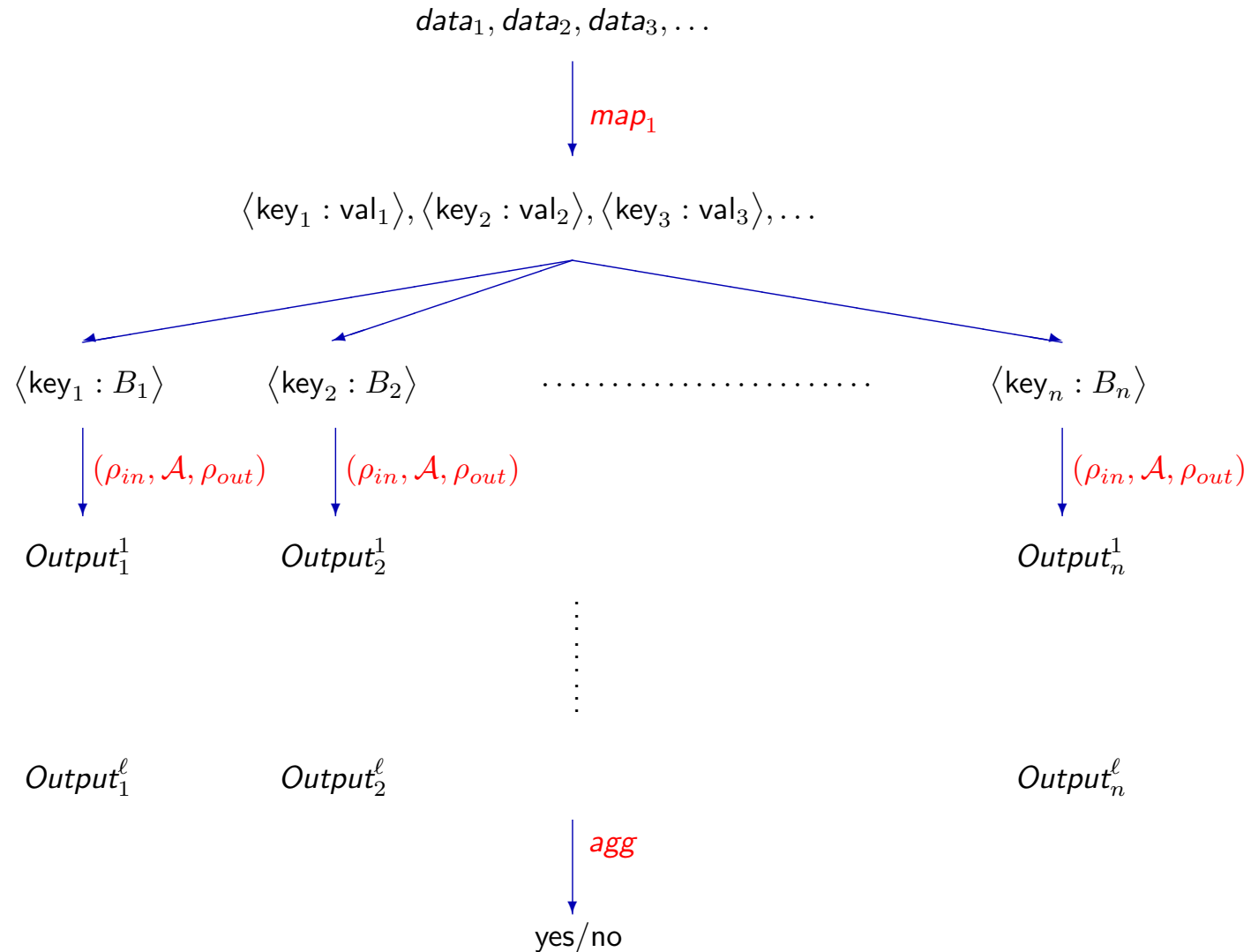
Originally studied by Francez, Kaminski and Shemesh.

DSA: The reducer $red_i = (\rho_{in}, \mathcal{A}, \rho_{out})$

On input $\langle \text{key} : B \rangle$, where $B = \{\sigma_1(\bar{a}_1), \dots, \sigma_n(\bar{a}_n)\}$



Boolean DSA: $M = (map_1, red_1, \dots, map_\ell, red_\ell, agg)$



DSA & Relational Algebra

Theorem. *Non-Boolean DSA can perform intersection, union, set difference, projection and selection.*

Theorem. *For every FO sentence (rel. algebra) φ and an integer m , there is DSA $M_{\varphi,m}$ such that for every database DB of degree $\leq m$, $M_{\varphi,m}$ accepts DB if and only if $\text{DB} \models \varphi$.*

- It still holds for FO with modulo counting ($\exists^{i \bmod m} x$).
- The proof is via Hanf locality.

DSA cannot detect the existence of walk

Definition. ℓ -WALK = $\{G \mid G \text{ contains a walk of length } \geq \ell\}$.

Theorem. *There is no DSA that accepts 5-WALK.*

Theorem. *On bounded degree graphs, there is ℓ -round DSA that accepts 2^ℓ -WALK.*

Theorem. *There is no ℓ -round DSA that accepts $2^{\ell+1}$ -WALK on graphs with degree at most 2.*

The big picture for DSA

DSA

DST

DSTJ

U

Rel. alg. without join

Semi-join algebra

Relational algebra

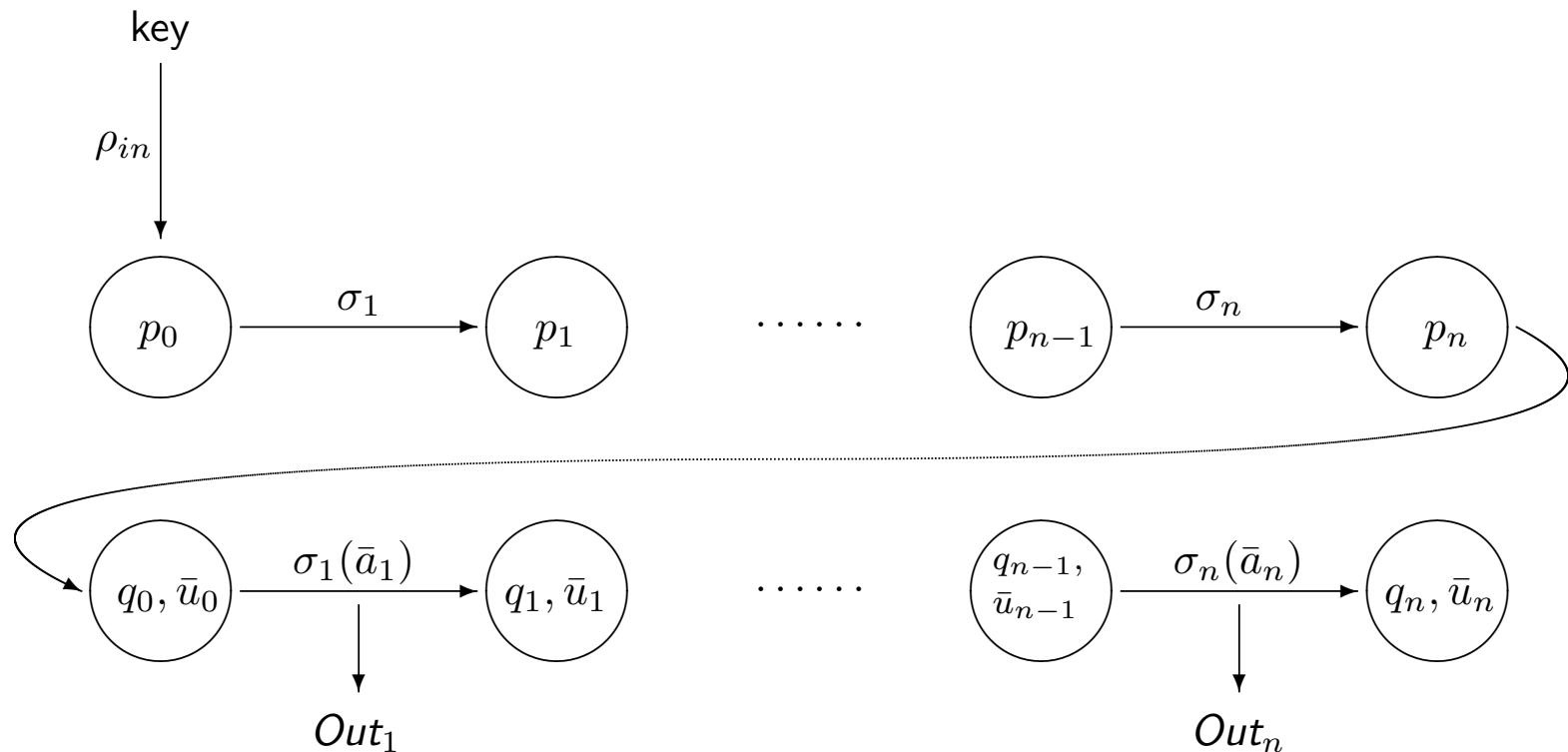
(FO over bounded deg. DB)

DST: $M = (map_1, red_1, map_2, red_2, \dots, map_\ell, red_\ell, agg)$

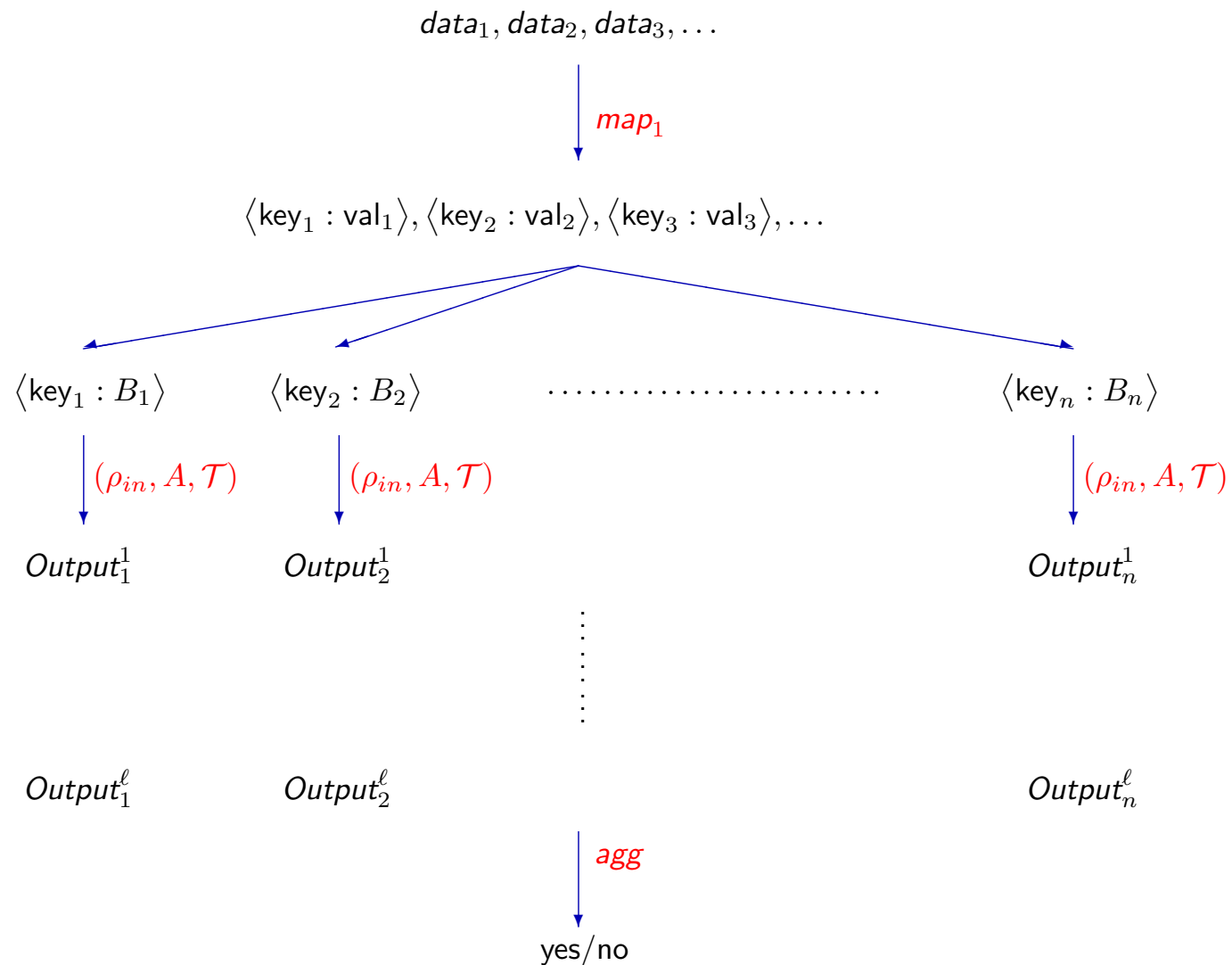
- Each map_i is a generic mapper.
- Each red_i is finite state automaton and deterministic “register transducer.”
- agg is deterministic register automaton.

DST: The reducer $red_i = (\rho_{in}, A, \mathcal{T})$

On input $\langle \text{key} : B \rangle$, where $B = \{\sigma_1(\bar{a}_1), \dots, \sigma_n(\bar{a}_n)\}$



Boolean DST: $M = (map_1, red_1, \dots, map_\ell, red_\ell, agg)$



DST & semi-join algebra

Theorem. *Non-Boolean DST can perform intersection, union, set difference, projection, selection and semi-join, i.e. semi-join algebra.*

DST and walks on graphs

Theorem.

- For every $\ell \geq 0$, there is an ℓ -round DST that accepts (2ℓ) -WALK.
- For every $\ell \geq 0$, there is no ℓ -round DST that accepts $(2\ell + 2)$ -WALK.
- $(\ell + 1)$ -round DSTs are strictly stronger than ℓ -round DSTs.

Corollary. DSTs are strictly more expressive than DSAs.

Recall that there is no DSA that accepts 5-WALK.

DST & TRIANGLE

Definition. $\text{TRIANGLE} = \{G \mid G \text{ contains a triangle}\}$.

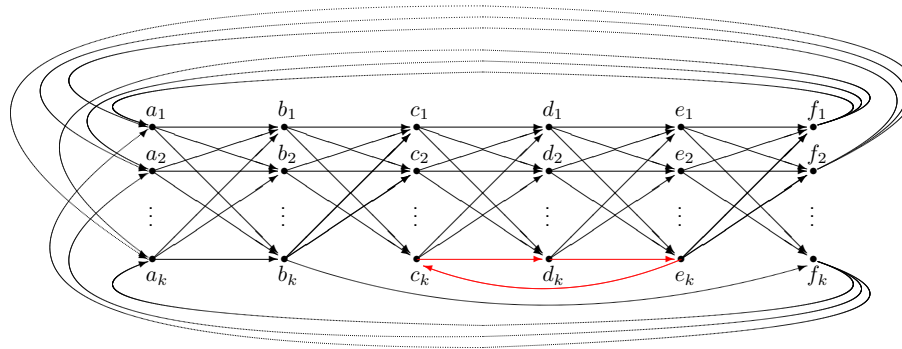
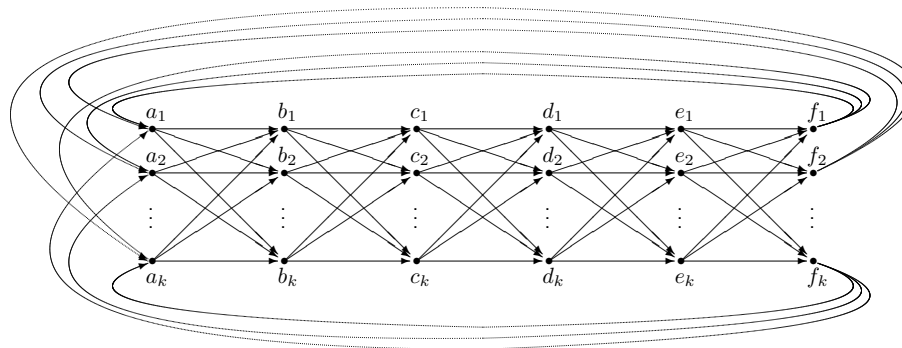
Theorem. *There is no DST that can accept TRIANGLE.*

Reason:

- The output of DST is linear in the size of the input database.
- The number of triangles in can be $m^{3/2}$, where m is the number of edges.

DST cannot accept TRIANGLE

For every DST M , there is k such that M cannot differentiate:



Big picture for DSA, DST

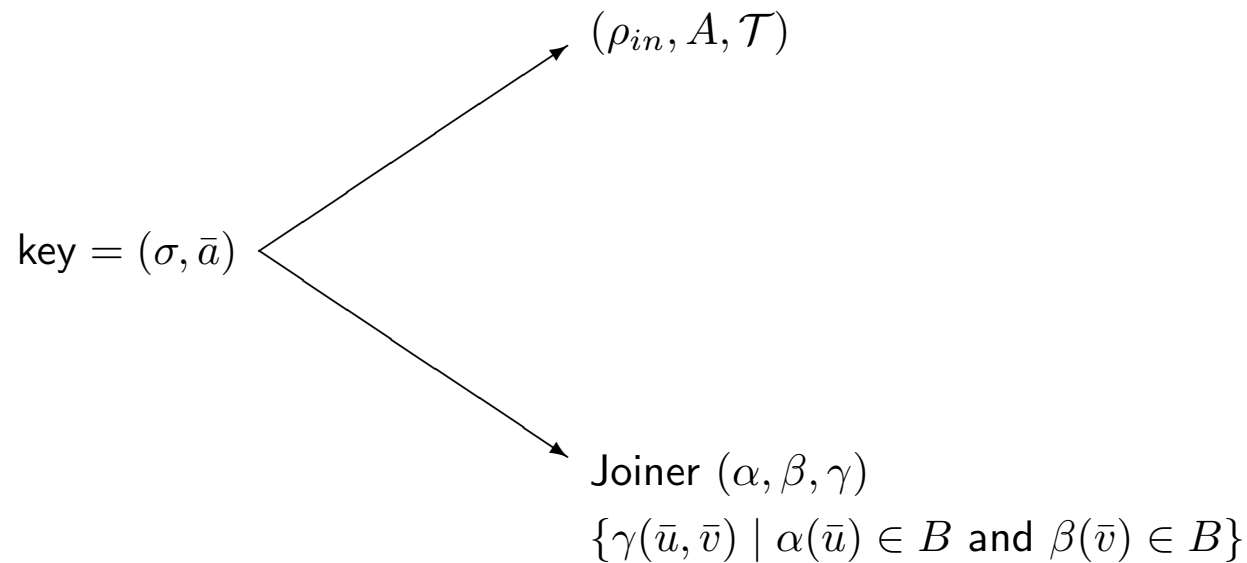


DSTJ: $M = (map_1, red_1, \dots, map_\ell, red_\ell, agg)$

- Each map_i is a generic mapper.
- Each red_i is either:
 - $(\rho_{in}, A, \mathcal{T})$, as in DST, or
 - a **joiner** to perform Cartesian product.
- agg is deterministic register automaton.

DST: The reducer red_i

On input $\langle \text{key} : B \rangle$



DSTJ & relational algebra

Theorem. *Non-Boolean DSTJ can perform intersection, union, set difference, projection, selection, semi-join and equi-join.*

Theorem. *There is a 2-round DSTJ that accepts TRIANGLE.*

Corollary. *DSTJs are strictly stronger than DSTs.*

DSTJ and cycles in graphs

Definition. ℓ -CYCLE = $\{G \mid G \text{ contains a cycle of length } \ell\}$.

Theorem.

- For every $\ell \geq 0$, there is an ℓ -round DSTJ that accepts (2^ℓ) -CYCLE.
- For every $\ell \geq 0$, there is no ℓ -round DSTJ that accepts $(2^{\ell+1})$ -CYCLE.
- $(\ell + 1)$ -round DSTJs are strictly stronger than ℓ -round DSTJs.

Big picture for DSA, DST, DSTJ



Applicability

We built a software that implements the “algorithmic idea” of DSTJ.

- His name is GUMBO.
- He is built on top of Hadoop.
- He performs parallel multiple semijoins. In fact, he can perform guarded fragment queries.
- His performance is comparable with Pig and Hive.
- He has additional novel strategies that make him more effective in evaluating multi-semijoins.
- He can be found here:
<https://zenodo.org/record/51517#.V7hxt5N96fU>

References

- Gumbo: Guarded Fragment Queries over Big Data.
EDBT 2015 (demo).
Jonny Daenen, Frank Neven, Tony Tan.
- Distributed Streaming with Finite Memory.
ICDT 2015.
Frank Neven, Nicole Schweikardt, Frédéric Servais, Tony Tan.
- Parallel Evaluation of Multi-Semi-Joins.
Jonny Daenen, Frank Neven, Tony Tan, Stijn Vansummeren.
PVLDB 9(10): 732-743 (2016)