

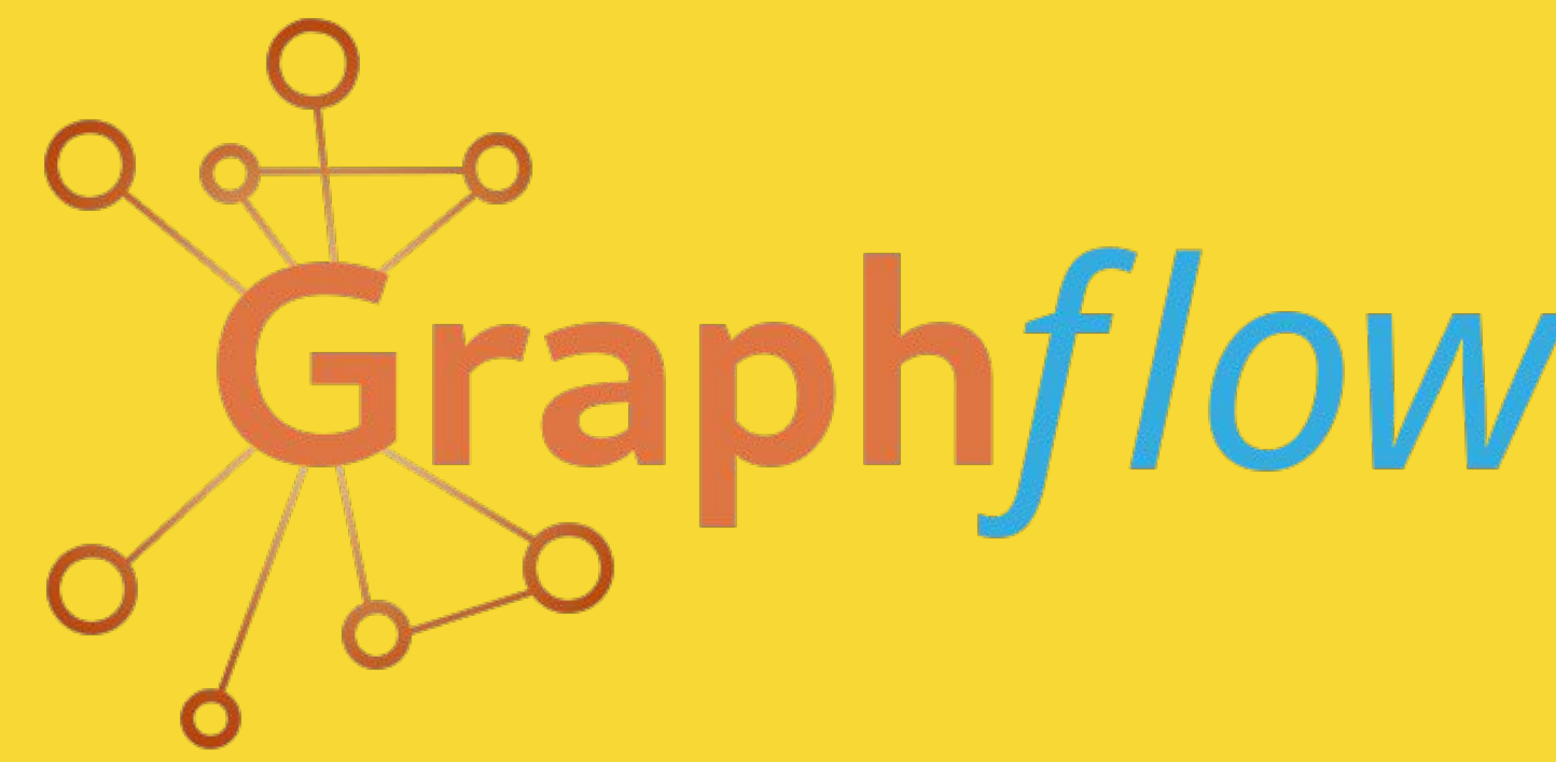
Columnar Storage and List-based Processing for Graph Database Management System

Pranjal Gupta, Amine Mhedhbi, Semih Salihoglu



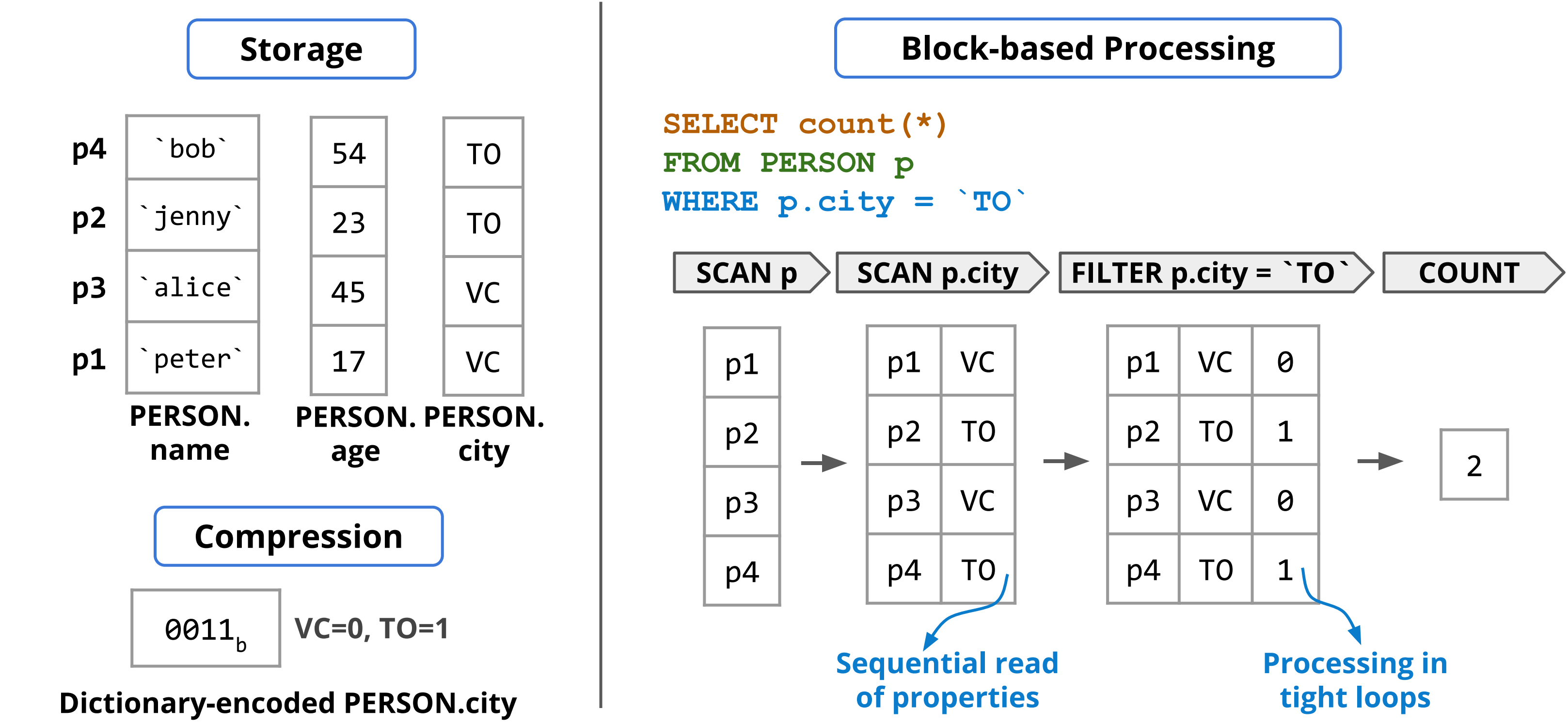
UNIVERSITY OF
WATERLOO

DSg Data
Systems
Group



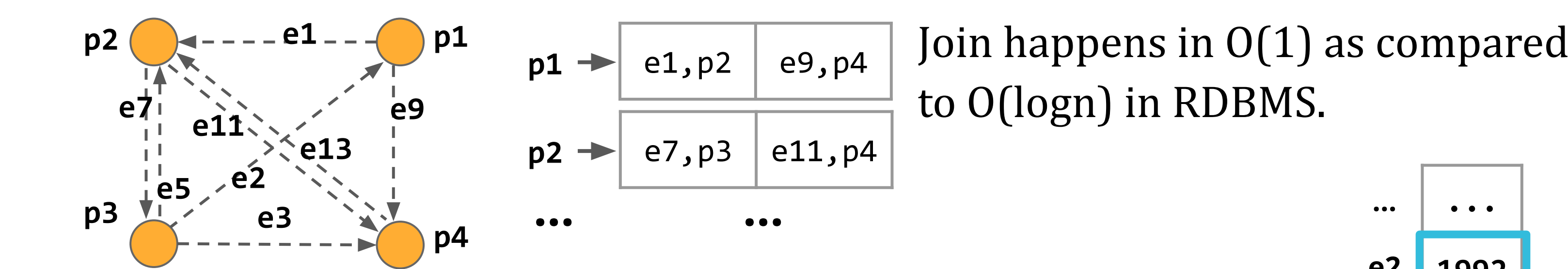
COLUMNAR RDBMSs

- Read-optimized systems, optimized for read-heavy workloads that involves reading and processing large chunks of data *sequentially*.



GDBMSs and GRAPH DATA

- Read-optimized systems, performant on many-to-many joins.
- Adjacency Lists for Fast Joins**

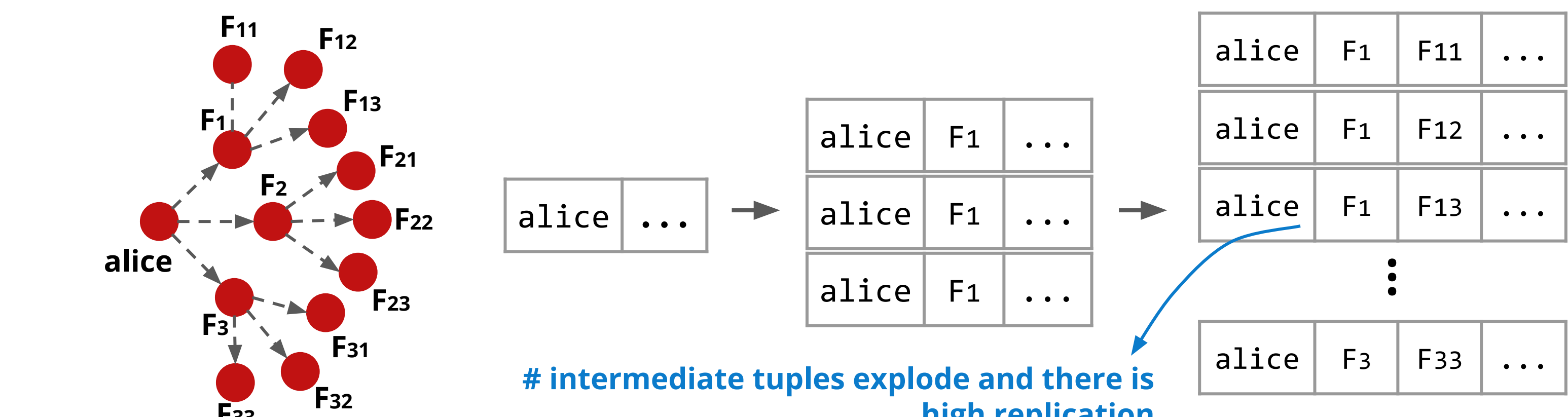


Random Data Access

```
MATCH (a:Person) -[e:Knows]->(b:Person)
WHERE a.name = 'alice'
RETURN b.age, e.since
```

SCAN a → FILTER a.name = 'alice' → JOIN b → SCAN b.age, e.since

Prominence of many-to-many relationships in graph data

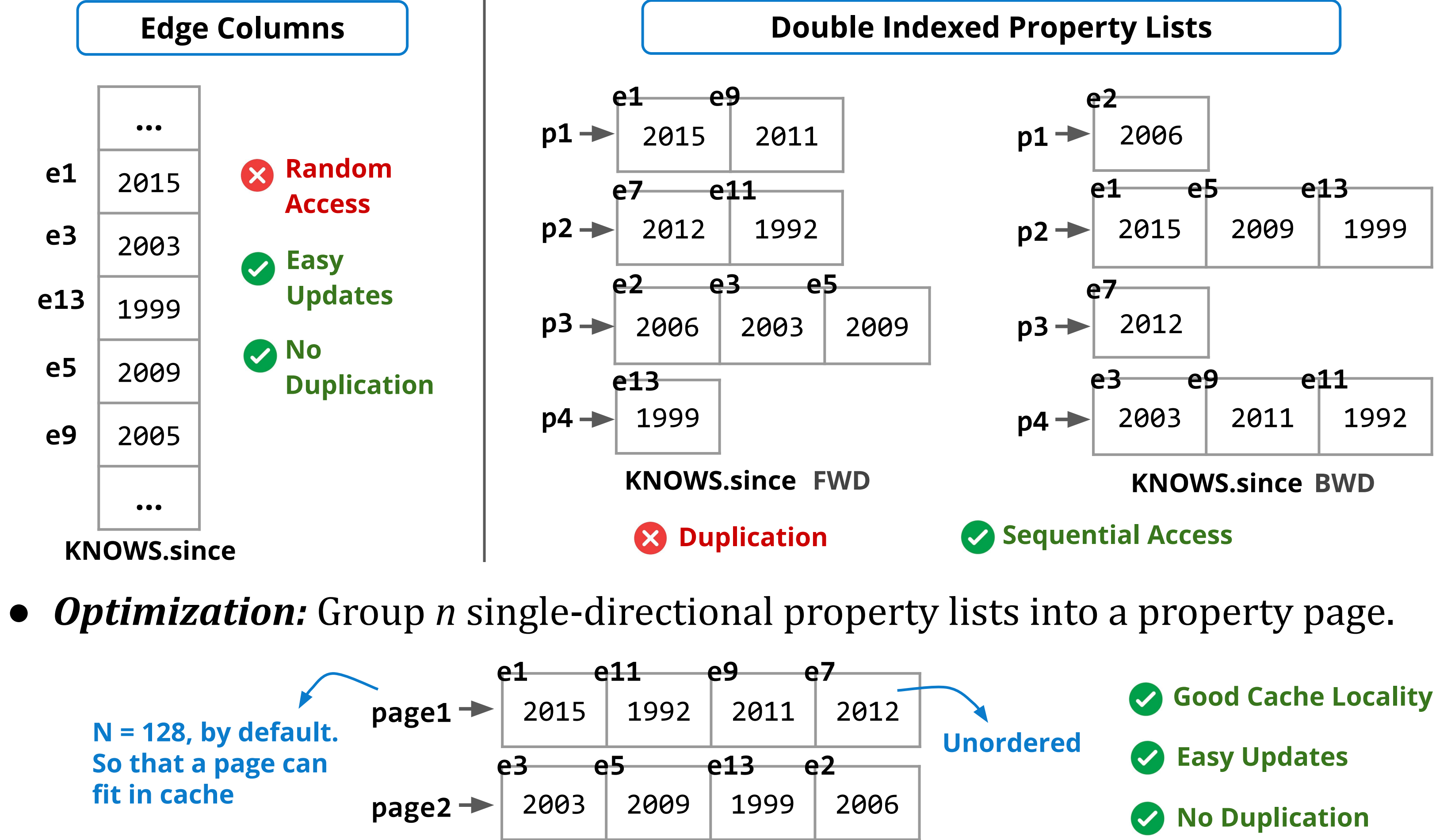


Partial Structuredness of graph data

- EDGE LABEL → SRC VERTEX LABEL and DEST VERTEX LABEL
- EDGE LABEL → Cardinality of edges
- EDGE/VERTEX LABEL → Properties keys

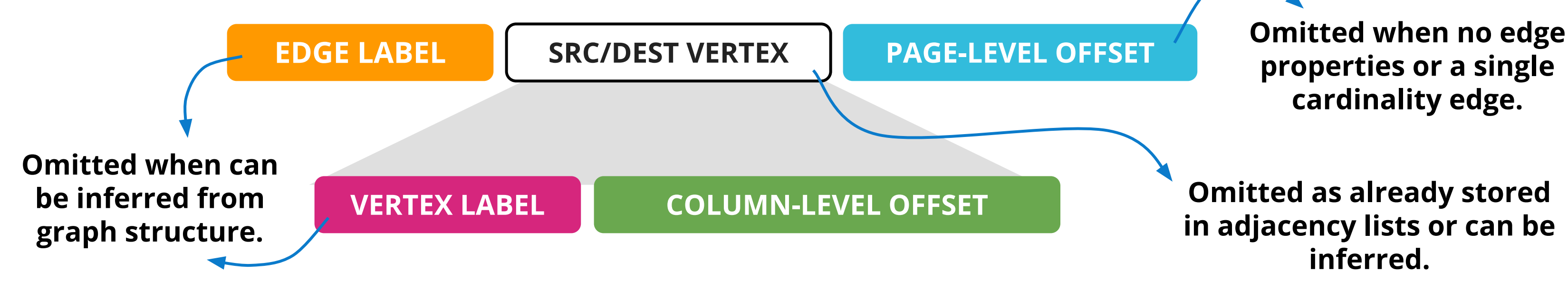
SINGLE-INDEXED EDGE PROPERTY PAGES

- Requirement:** Store and access the properties of edges sequentially in the order edges appear in the adjacency lists.

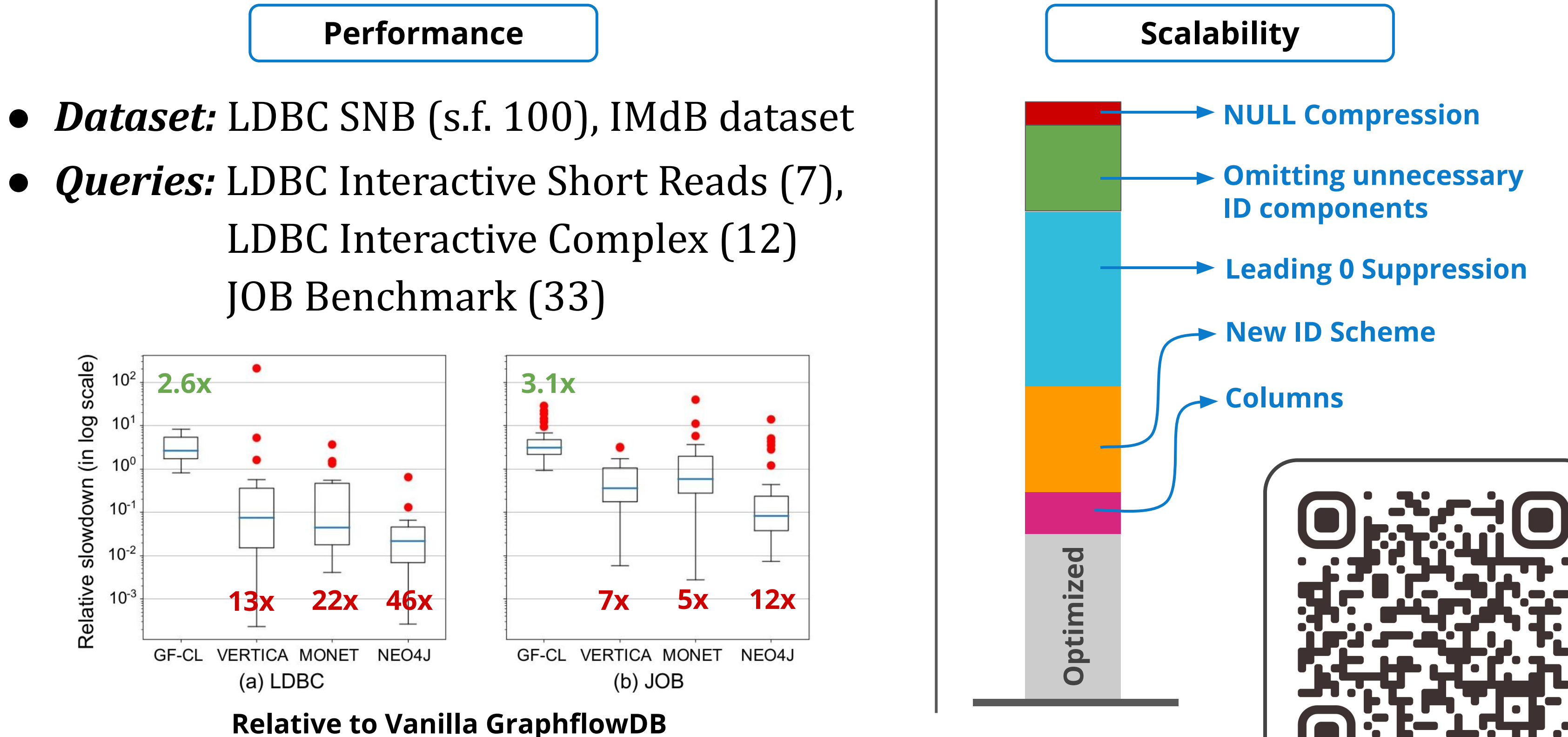


NEW ID SCHEME AND COMPRESSION

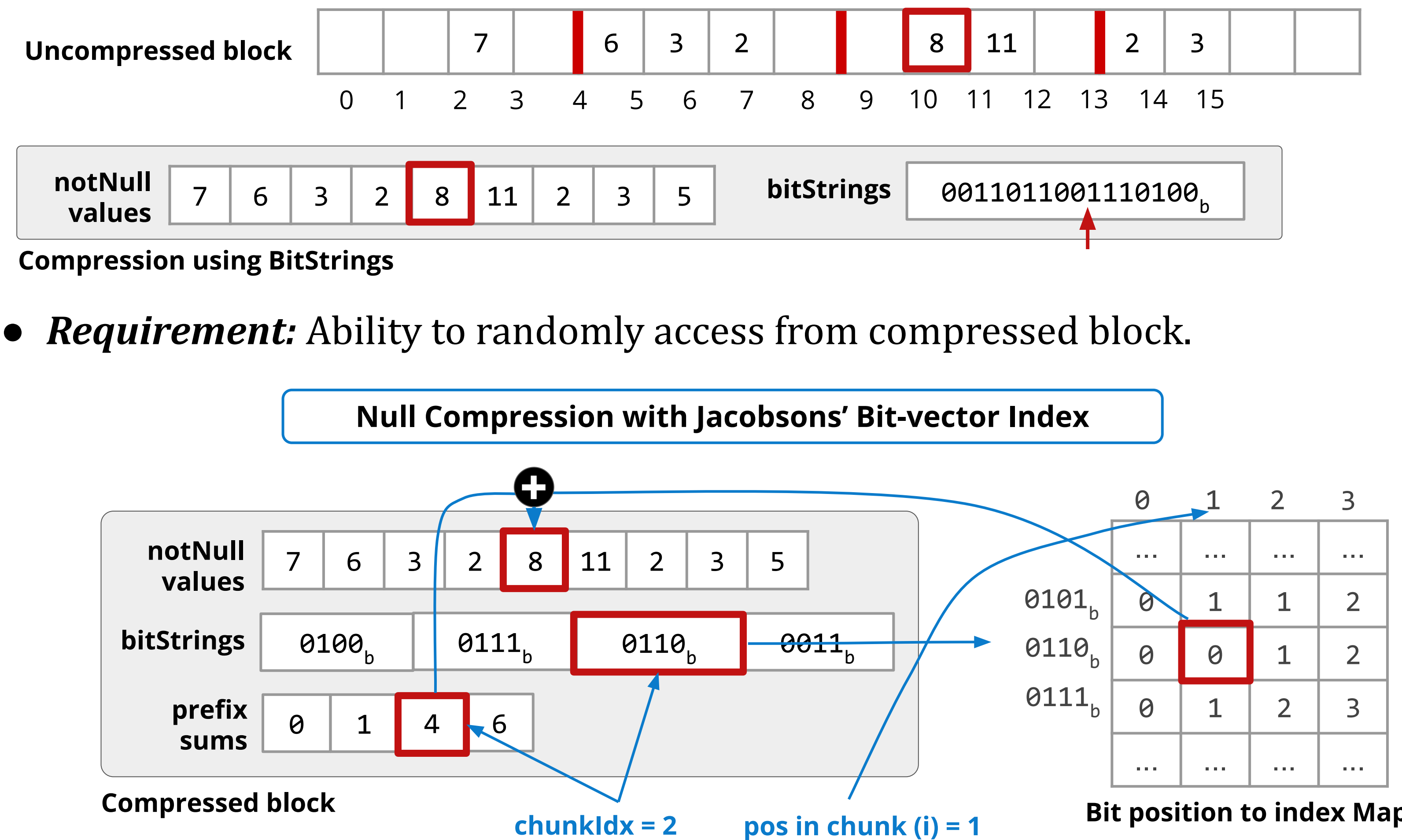
- Requirement:** To support random access to properties in Columns & Property Pages. Also to compress when storing in Adjacency Lists.



EVALUATION



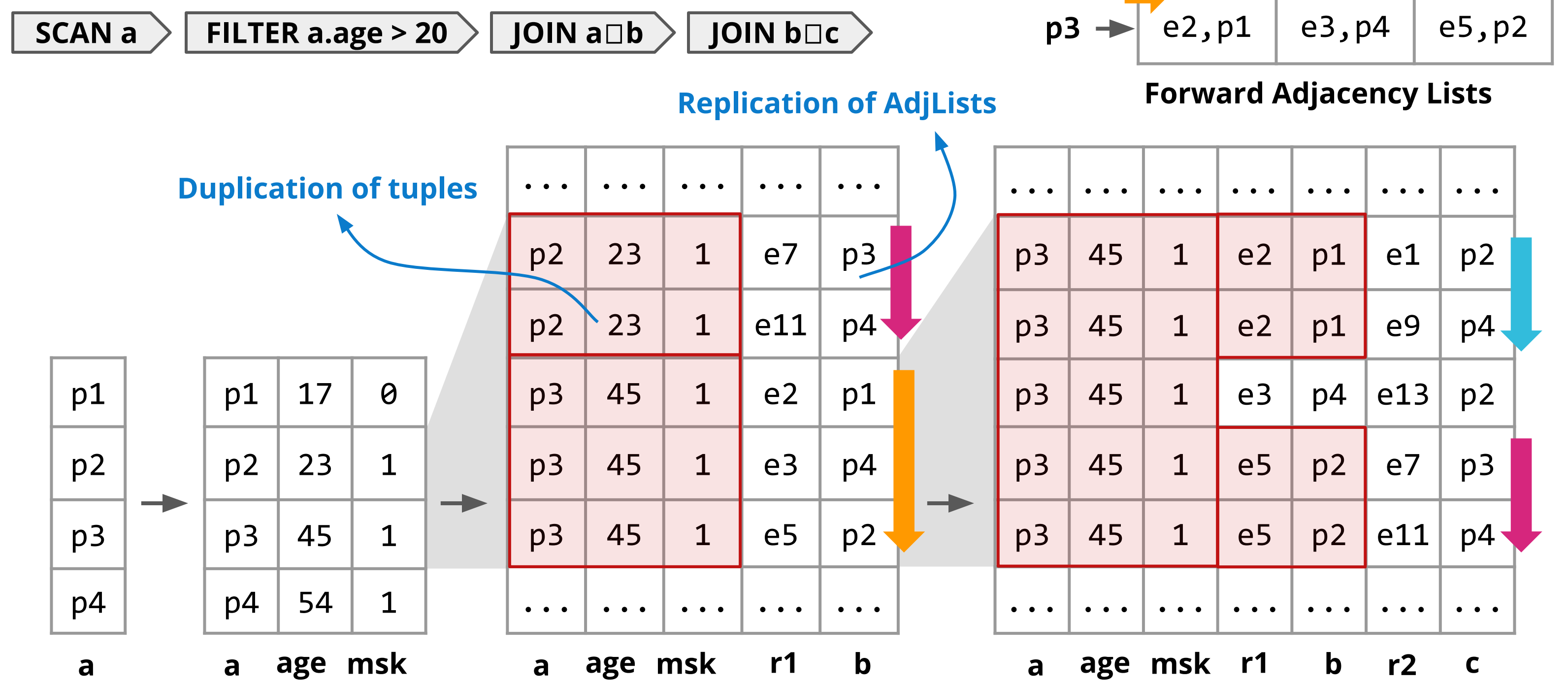
NULL COMPRESSION



LIST-BASED PROCESSOR

- Block-based Processor:** Not optimized for many-to-many joins

```
MATCH (a:PERSON) -[r1:FOLLOWS] -> (b:PERSON)
(b:PERSON) -[r2:FOLLOWS] -> (c:PERSON)
WHERE a.age > 20
RETURN ...
```



- Requirement:** Avoid duplication of data and materializing of Adjacency lists

Solution:

- Factorized Representation
- List Groups

No Fixed Size

Do not materialize AdjLists

Flattened or Unflattened

