

DISCOVERING UNIQUE COLUMN COMBINATIONS ON DYNAMIC DATA

Motivation

Production of big datasets at very fast rates:

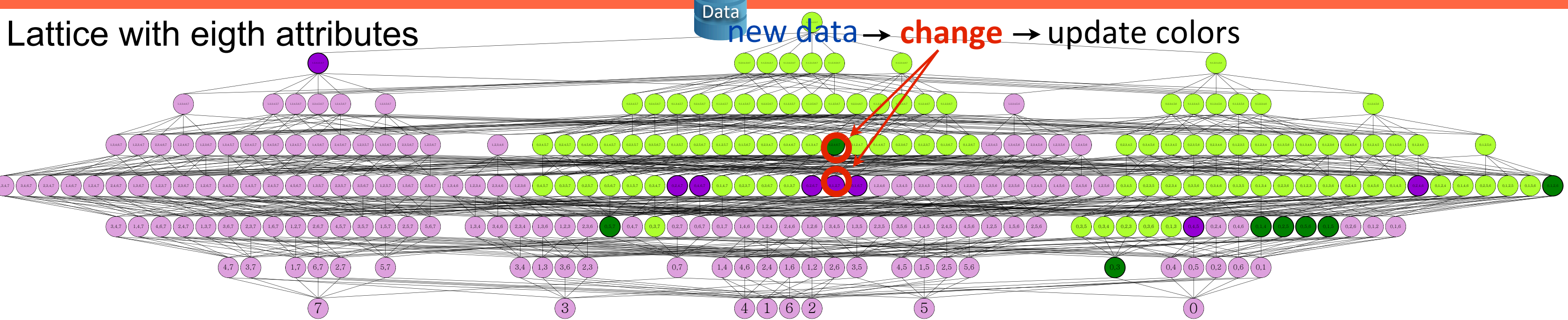
- Social networks
- Scientific applications
- Transactional applications
- ...

Knowing uniques is crucial for:

- Query optimization
- Anomaly detection
- Data modeling
- Indexing
- ...

Problem

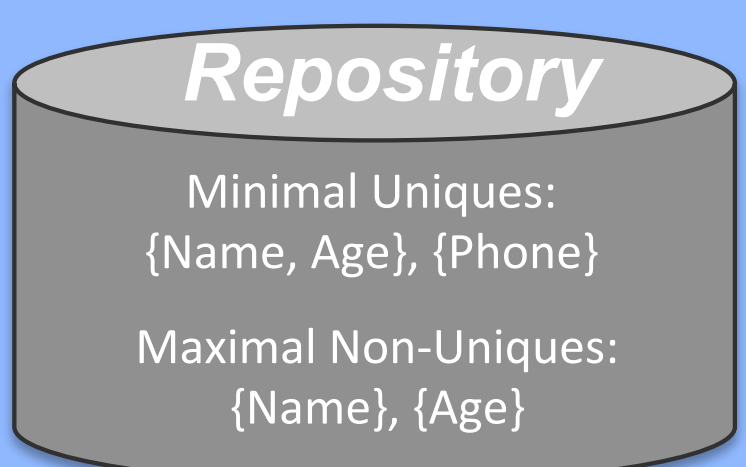
- Unique column combinations are quite often unknown in big datasets
- Exponential search space



- Finding unique column combinations is an NP-Hard problem
- Big datasets change quite often

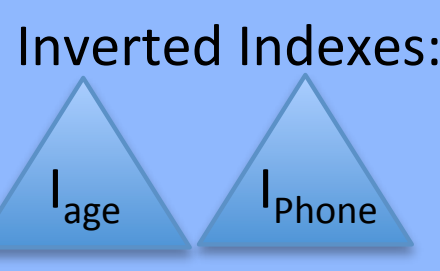
SWAN

Dealing with Inserts



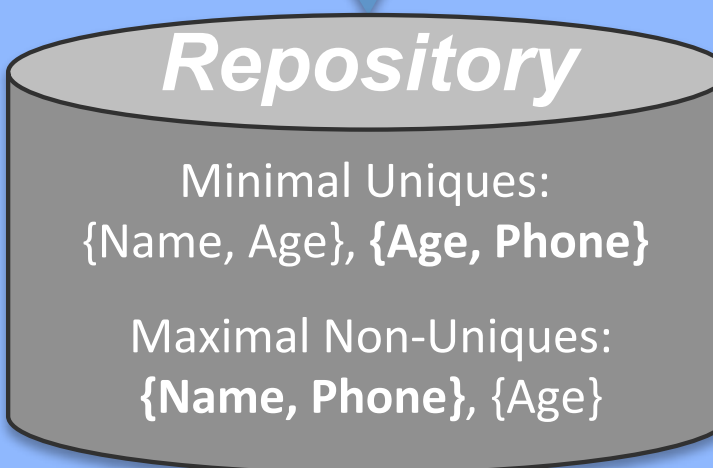
Tuple ID	Name	Phone	Age
(tuple ₁)	Lee	345	20
(tuple ₂)	Payne	245	30
(tuple ₃)	Lee	234	30

(insert ₁)	Payne	245	31
------------------------	-------	-----	----

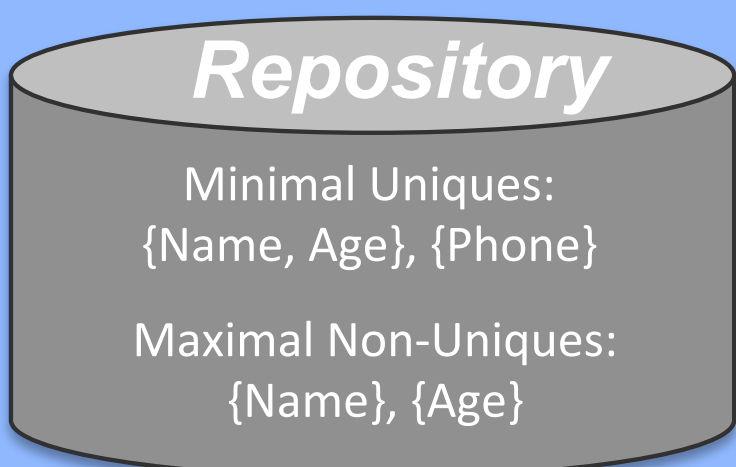


Retrieved tuples IDs with I_{Age} : {}
Retrieved tuples IDs with I_{Phone} : {(tuple₂)}

(tuple ₂)	Payne	245	30
(insert ₁)	Payne	245	31

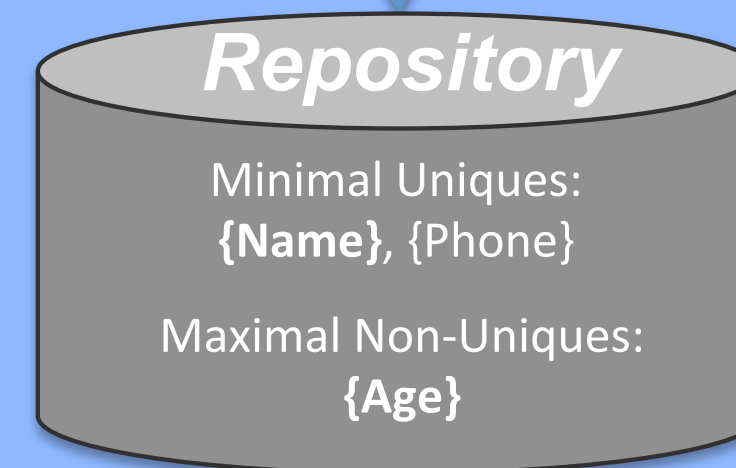
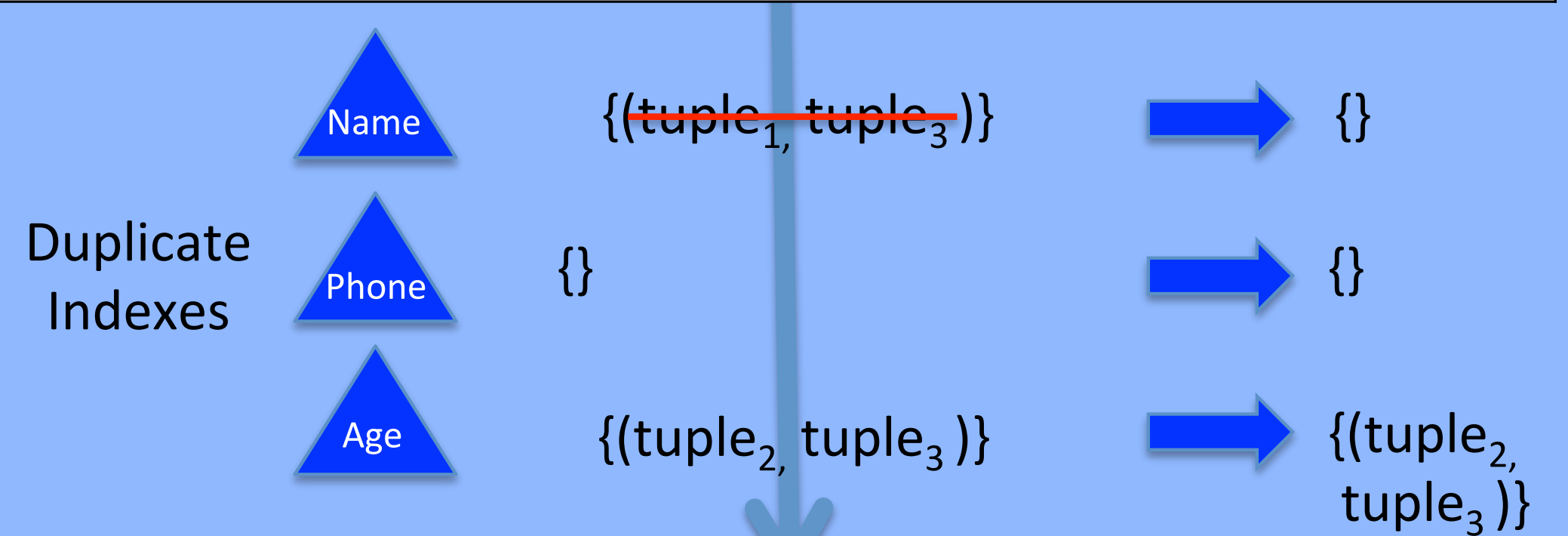


Dealing with Deletes



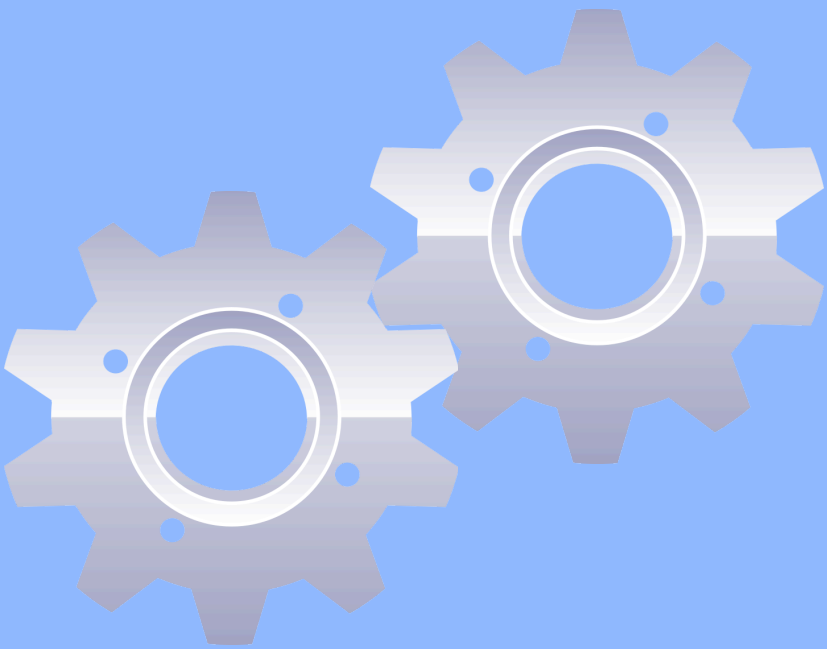
Tuple ID	Name	Phone	Age
(tuple ₁)	Lee	345	20
(tuple ₂)	Payne	245	30
(tuple ₃)	Lee	234	30

Delete:	(tuple ₁)	Lee	345	20
---------	-----------------------	-----	-----	----



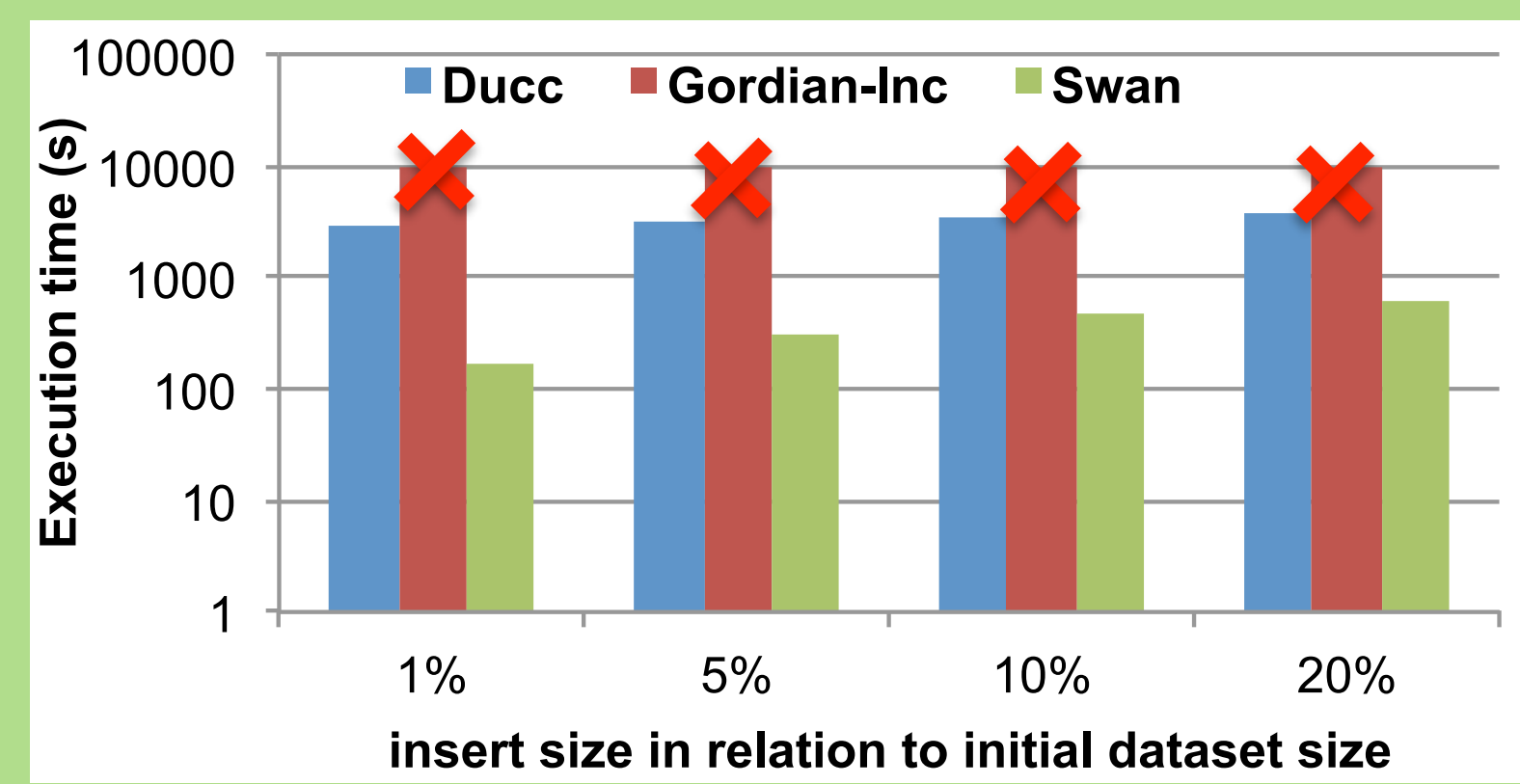
Swan Architecture

Index Creation Mechanism

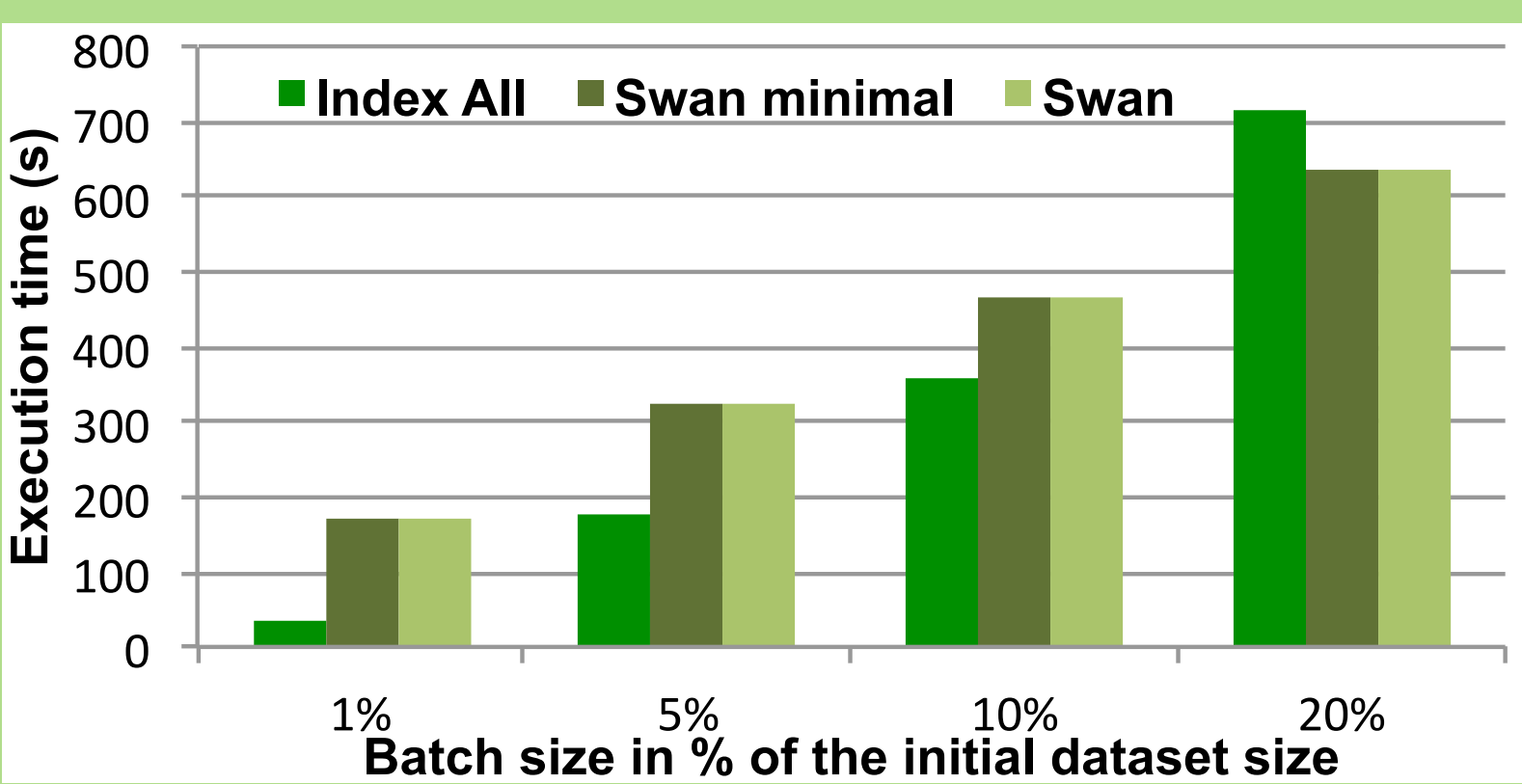


- Full scans for each change?
- Indexing all attributes?
- Minimal set of indexes:
 - cover all minimal uniques
- Speed-up indexes:
 - reducing number of tuples to retrieve

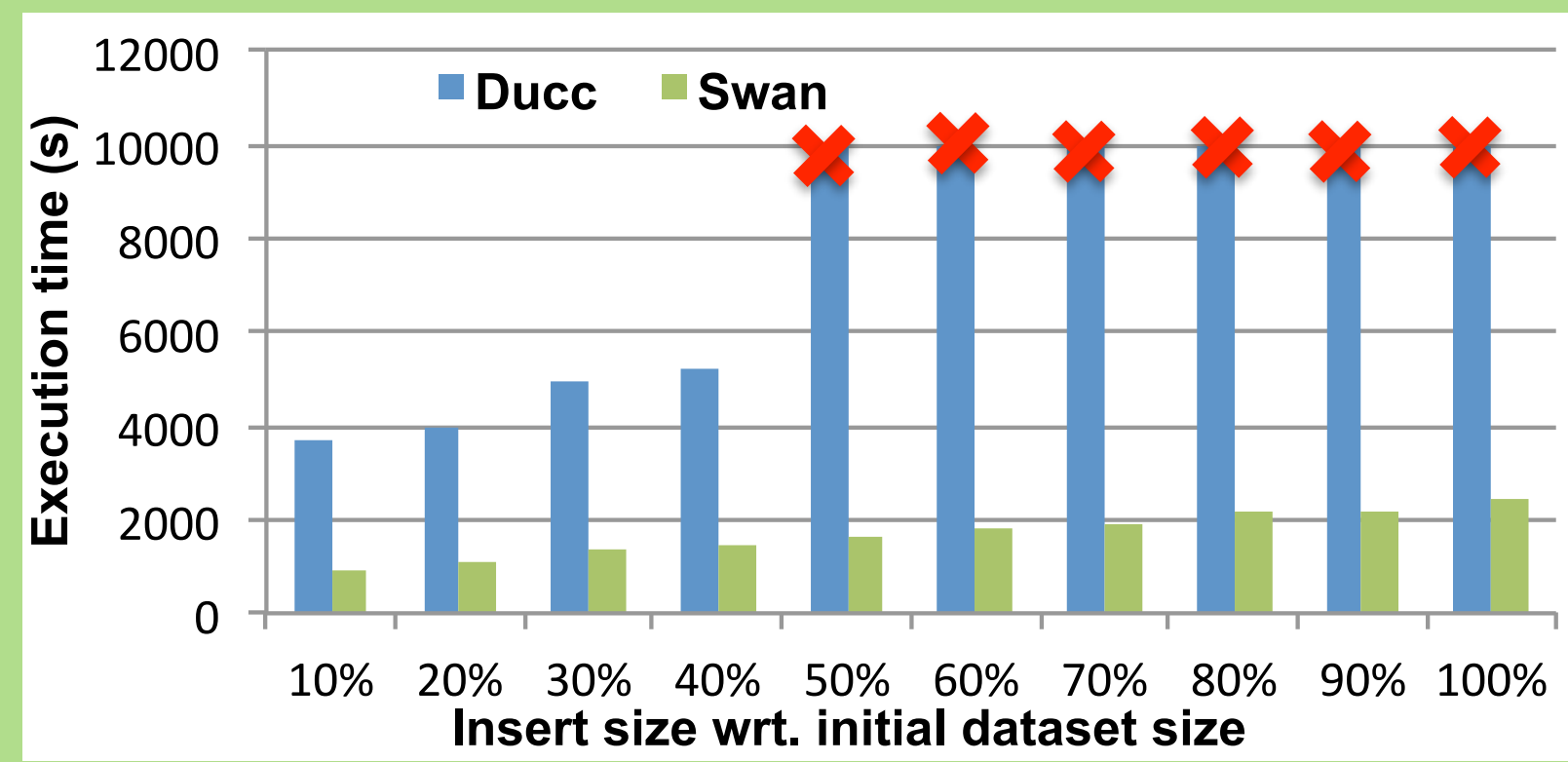
Results



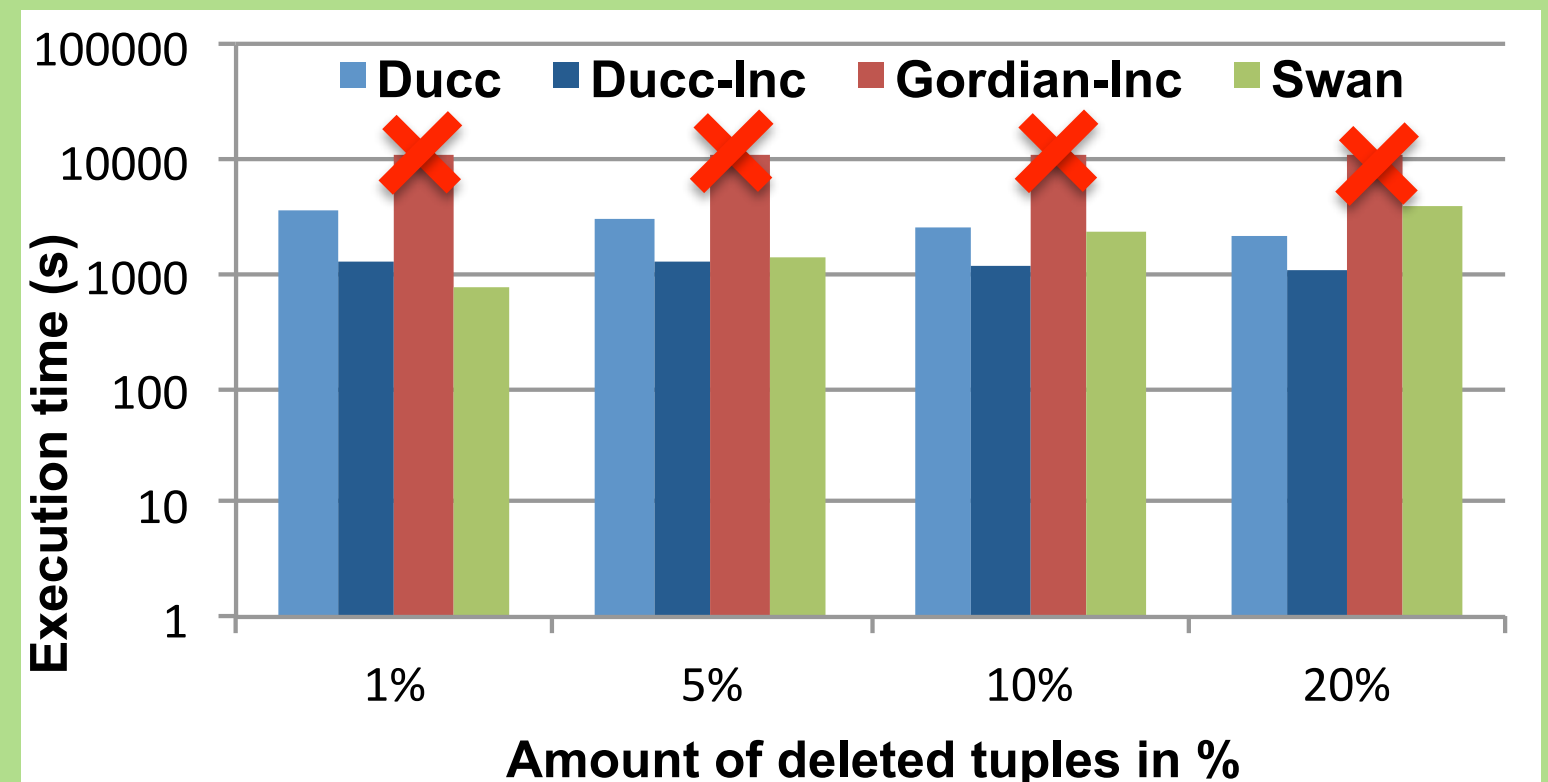
Dealing with inserts



Index evaluation



Stressing Swan



Dealing with deletes

5 millions rows, Lineitem TPC-H

Related Work

Gordian:

- Row-based approach
- Prefix-tree data organization

[Gordian: efficient and scalable discovery of composite keys. VLDB'06]

HCA:

- Column-based approach
- Histograms- and value-counting-based

[Advancing the discovery of unique column combinations. CIKM'11]

Ducc:

- Hybrid (row- and column-based) approach
- Depth-first + Random walk lattice traversal

[Scalable discovery of unique column combinations. VLDB'14]

Metanome

- It is a joint project between HPI and the QCRI
- It provides a fresh view on data profiling
- It aims at providing scalability for Big Data

Website: http://www.hpi.uni-potsdam.de/naumann/projekte/metanome_data_profiling.html