# TOWARDS INTERACTIVE
## DATA EXPLORATION

CARSTEN BINNIG

DATA MANAGEMENT LAB

Data-Driven Marketing

Data-Driven Science

Data-Driven Medicine

Data-Driven Production

**DATA EXPLORATION:**
Important first step to understand
**BIG DATA**

Airline Traffic Data

jeff@nytimes.com

Minority Report (2002)

DATA EXPLORATION VISION

4

**TODAY'S USER INTERFACES**

# TODAY'S USER INTERFACES

# … AND THE BIG DATA SYSTEMS?

# A TYPICAL EXPLORATION PIPELINE

**How do query interfaces need to change?**

*Vizdom (Visual Exploration)*
*DBPal / EchoQuery (NL Interface)*
*IDEBench (Benchmarking)*

**How do we enable more high-speed execution?**

*IDEA (Interactive Query Processing)*
*I-Store (Analytics on Modern Hardware),*
*XDB (Scalable Cloud Analytics)*

**How do we reduce data cleaning costs?**

*UnkownUnkowns (Data Quality)*
*IncMap (Schema Mapping)*
*Sherlock (Text Summerization)*

```
text = "Research has shown that it is often still
# insert your code here.. I suppose it's obvious
#text=text.replace("a","")
vowels=['a','e','i','o','u'];
for vowel in vowels:
    text=text.replace(vowel,"");
print(text)

Rsrch hs shwn tht t s ftn stll pssbl t ndrstnd tx
```

sales - Kladblok
Bestand  Bewerken  Opmaak  Beeld  Help
```
"Country","Salesperson","Order Amount","Quarter"
"UK","Smith",16753,"Qtr 3"
"USA","Johnson",14808,"Qtr 4"
"UK","williams",10644,"Qtr 2"
"USA","Jones",1390,"Qtr 3"
"USA","Brown",4865,"Qtr 4"
"UK","williams",12438,"Qtr 1"
"UK","Johnson",9339,"Qtr 2"
"USA","Smith",18919,"Qtr 3"
"USA","Jones",9213,"Qtr 4"
"UK","Jones",7433,"Qtr 1"
"USA","Brown",3255,"Qtr 2"
"USA","williams",14867,"Qtr 3"
"UK","williams",19302,"Qtr 4"
"USA","Smith",9698,"Qtr 1"
"USA","Jones",18978,"Qtr 2"
"UK","Brown",9080,"Qtr 4"
```

VISUAL INTERACTIVE DATA EXPLORATION

# vizdom

## Interactive Analytics through Pen and Touch

Andrew Crotty, Alex Galakatos, Emanuel Zgraggen, Carsten Binnig, Tim Kraska

Challenges & opportunities

# CHALLENGE: INTERACTIVITY

**Provide interactive response times for queries** even on very large data sets (e.g., <500ms)

## The Effects of Interactive Latency on Exploratory Visual Analysis

Zhicheng Liu and Jeffrey Heer

In this research, we have found that interactive latency can play an important role in shaping user behavior and impacts the outcomes of exploratory visual analysis. Delays of 500ms incurred significant costs, decreasing user activity and data set coverage while reducing rates of observation, generalization and hypothesis. Moreover, initial exposure to higher latency interactions resulted in reduced rates of observation and generalization during subsequent analysis sessions in which full system performance was restored.

# CHALLENGE: AD-HOC QUERIES

**Provide ad-hoc intuitive query interfaces** **AND no pre-defined static reports or low-level query interfaces**



**Census Data**

**Next steps:** influence of education, marital status, …?

# CHALLENGE: CONNECT & EXPLORE

**Users want to directly explore new data without waiting for data being loaded (or even cleaned) before**

# OPPORTUNITY: DECISION MAKING

**Exact results** are often not needed to make decisions



**Optimization:** Approximate results are good enough

# OPPORTUNITY: INCREMENTAL QUERIES



```
SELECT SUM(salary)
FROM census
GROUP BY salary-buckets
```

```
SELECT SUM(salary)
FROM census
WHERE age < 60
GROUP BY salary-buckets
```

**Optimization:** Reuse results / compute only the diff!

# OPPORTUNITY: THINK TIME

**User typically look at results for a significant amount of time ("think time") before executing next step**

**Optimization:** Speculative execution while user "thinks"

# HOW TO BUILD A BACKEND FOR VISUAL IDE?

# OUR APPROACH: IDEA

## IDEA = Interactive Data Exploration Accelerator



- Connect & Explore for new Data Sources
- **Progressive & Approximate Query Processing**
- **Incremental Query Building & Reuse**

*Andrew Crotty, Alex Galakatos, Emanuel Zgraggen, Carsten Binnig, Tim Kraska: The case for interactive data exploration accelerators (IDEAs). HILDA@SIGMOD 2016*

# BASIC IDEA OF AQP (FORM THE 90'S)

**Sales**

| Product | Amount |
|---------|--------|
| CPU | 1 |
| CPU | 1 |
| CPU | 2 |
| CPU | 3 |
| CPU | 4 |
| Disk | 1 |
| Disk | 2 |
| Monitor | 1 |

**SalesSample**

| Product | Amount |
|---------|--------|
| CPU | 1 |
| CPU | 2 |
| CPU | 3 |
| Disk | 2 |

SELECT SUM(Amount) FROM Sales WHERE Product = 'CPU'

Exact Answer:
1+1+2+3+4 = 11

Approx. Answer:
(1+2+3)*2= 12

# AQP: SPEED/ACCURACY TRADE-OFF

# IS CLASSICAL AQP GOOD ENOUGH?



No support for **incremental query building & reuse of intermediate results**

# OUR AQP FORMULATION

**Main idea:** results of prior approximate queries are represented as random variables X



$Pr(X=Male)=0.75$

$Pr(X=Female)=0.25$

**Enables reuse of approximate results with error bounds**

*Alex Galakatos, Andrew Crotty, Emanuel Zgraggen, Carsten Binnig, Tim Kraska: Revisiting Reuse for Approximate Query Processing. PVLDB 2017*

# AQP: RESULT REUSE

**Executed Interactions**



## Result Cache

| $P_{male}$ | $\{0.70, \varepsilon_1\}$ |
|---|---|
| $P_{female}$ | $\{0.30, \varepsilon_2\}$ |
| $P_{high}$ | $\{0.20, \varepsilon_3\}$ |
| $P_{low}$ | $\{0.80, \varepsilon_4\}$ |
| $P_{low|male}$ | $\{0.75, \varepsilon_5\}$ |
| $P_{low|female}$ | $\{0.92, \varepsilon_6\}$ |
| $P_{high|male}$ | $\{0.25, \varepsilon_7\}$ |
| $P_{low|female}$ | $\{0.08, \varepsilon_8\}$ |

# AQP: RESULT REUSE

salary

sex

Count

400M

400M

**Result Cache**

| | |
|---|---|
| P | {0.70, $\varepsilon$ } |
| | |
| | |
| | |
| | |
| $P_{low|female}$ | {0.02, $\varepsilon_6$} |
| $P_{high|male}$ | {0.25, $\varepsilon_7$} |
| $P_{low|female}$ | {0.08, $\varepsilon_8$} |

**Other Reuse Potentials**
- Law of total probability
- Inclusion-exclusion principle

**Bay**

$$P(A \mid B) = \frac{P(B \mid A) \cdot P(A)}{P(B)}$$

$$P_{male|high} = \frac{P_{high|male} * P_{male}}{P_{high}} \approx 0.88$$

# IDEA PERFORMANCE RESULTS

## Exploration Session (User Study)

| # | Query |
|---|---|
| #1 | sex |
| #2 | education |
| #3 | education **WHERE** sex='Female' |
| #4 | education **WHERE** sex='Male' |
| #5 | sex, education |
| #6 | sex **WHERE** education='PhD' |
| #7 | salary |
| #8 | salary **WHERE** education='PhD' |
| #9 | sex, salary |
| #10 | salary **WHERE** sex='Female' |
| #11 | salary |
| #12 | salary **WHERE** sex='Female' |
| #13 | salary **WHERE** sex<>'Female' |
| #14 | salary **WHERE** sex='Female' **AND** education='PhD', salary **WHERE** sex<>'Female' **AND** education='PhD' |
| #15 | age |
| #16 | salary **WHERE** 20<=age<40 **AND** sex='Female' **AND** education='PhD', salary **WHERE** 20<=age<40 **AND** sex<>'Female' **AND** education='PhD' |

**Evaluated Systems:**

- **MonetDB:** Analytical Column-Store

- **Online Aggregation** (From Hellerstein. 90's)

- **IDEA:** Our System

**Data:** 500M tuples

| Census | #1 | #2 | #3 | #4 | #5 | #6 | #7 | #8 | #9 | #10 | #11 | #12 | #13 | #14 | #15 | #16 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| MonetDB | 0.34 | 0.39 | 5.40 | 8.70 | 0.48 | 1.20 | 1.20 | 0.91 | 0.53 | 4.80 | 0.42 | 4.70 | 1.10 | 5.60 | 1.60 | 7.10 |
| Online Agg | 0.05 | 0.24 | 0.78 | 0.59 | 0.24 | 0.46 | 0.04 | 0.48 | 0.07 | 0.11 | 0.04 | 0.11 | 0.08 | 7.53 | 0.29 | 24.3 |
| IDEA | 0.09 | 0.29 | 0.42 | 0.00 | 0.00 | 0.00 | 0.09 | 0.12 | 0.00 | 0.17 | 0.00 | 0.00 | 0.00 | 0.48 | 0.37 | 2.87 |

# MANY OTHER CONSIDERATIONS

Natural Language Interfaces

Benchmarking

Complex Workloads (ML, Text, …)

Hardware Acceleration

…

# MANY OTHER CONSIDERATIONS

**Natural Language Interfaces**

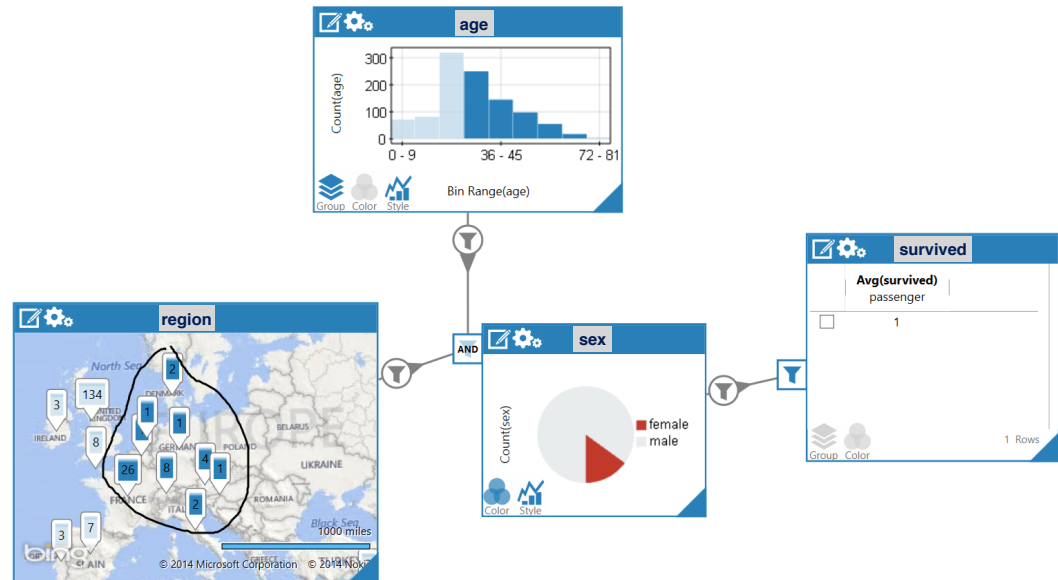Benchmarking

Complex Workloads (ML, Text, …)

Hardware Acceleration

…

# NL INTERFACE FOR DATABASES (NLIDB)

**NL Query:**

"How many older female people survived the sinking of the Titanic?"
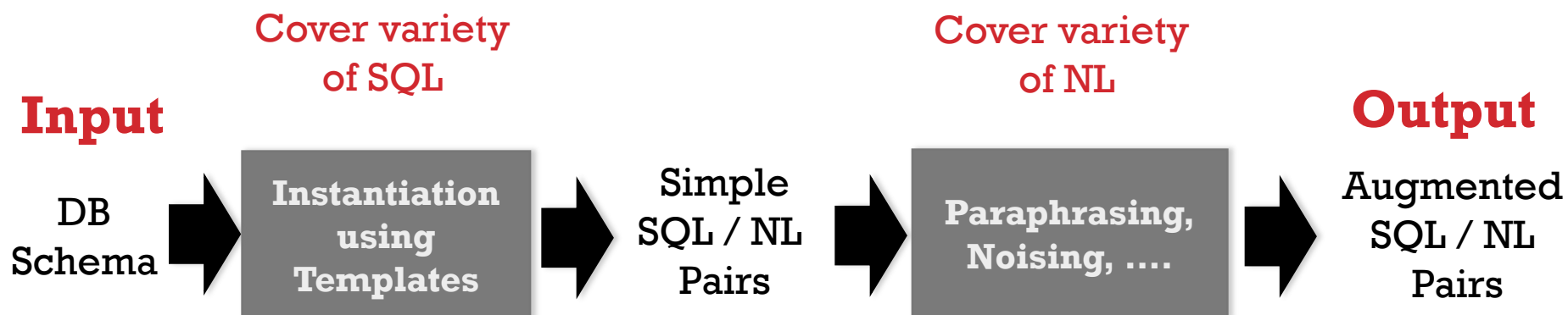
# ROBUST QUERY TRANSLATION?

## Language Translation Problem

Natural Language → [?] → SQL

How to get training data?

# OUR APPROACH: GENERATE TRAINING DATA

Cover variety
of SQL

Cover variety
of NL

**Input**

**Output**

DB
Schema

**Instantiation
using
Templates**

Simple
SQL / NL
Pairs

**Paraphrasing,
Noising, ….**

Augmented
SQL / NL
Pairs

**Distant Supervision:** Generate large (potentially noisy) training data instead of manually handcrafting it

*Fuat Basik, Benjamin Hättasch, Amir Ilkhechi, Arif Usta, Shekar Ramaswamy, Prasetya Utama, Nathaniel Weir, Carsten Binnig, Ugur Çetintemel: DBPal: A Learned NL-Interface for Databases. SIGMOD Conference 2018*

# OUR APPROACH: GENERATE TRAINING DATA

**Input**

DB Schema → **Generate using Templates** → SQL / Naïve NL Pairs → **Paraphrasing, Noising, …** → SQL / Augmented NL Pairs

**Output**

Cover variety of SQL

Cover variety of NL

---

*Template(s)*

SELECT <att>
FROM <table>
WHERE <filter>

Show me the <att>s of <table>s with <filter>?

*Naïve Corpus*

SELECT *name*
FROM *patient*
WHERE *diagnoses = fever*

Show me the names of patients with diagnoses fever?

*Paraphrasing*

Show me the names of patients diagnosed fever?

*Noising*

Show the names of patients with ~~diagnosed~~ fever?

| name | age | diagnoses |
|------|-----|-----------|
| Carsten | 39 | fever |
| Emilie | 8 | flu |
| Frederik | 4 | fever |

*Patient Table*

…

…

**Millions of different NL/SQL pairs**

# EXPERIMENTAL EVALUATION

**Evaluated Systems**

- **NaLIR:** Rule-based NLIDB (Best Paper VLDB 2015)

- **Neural Semantic Parser (NSP):** Neural Machine Translation (supervised learning -> **manual effort per database schema**)

- **DBPal:** Our Approach (distant supervision -> **NO manual effort** per database scheme)

|  | Patients | GeoQuery |
|---|---|---|
| NaLIR (w/o feedback) | 15.60% | 7.14% |
| NaLIR (w feedback) | 21.42% | N/A |
| NSP++ | N/A | **83.9%** |
| NSP (template only) | 10.60% | 5.0% |
| DBPal (w/o augmentation) | 74.80% | 38.60% |
| DBPal (full pipeline) | **75.93%** | 55.40% |

# DBPAL IN ACTION



http://titanx.smn.cs.brown.edu:8888/#/patients

# MANY OTHER CONSIDERATIONS

Natural Language Interfaces

**Benchmarking**

Complex Workloads (ML, Text, …)

Hardware Acceleration

…

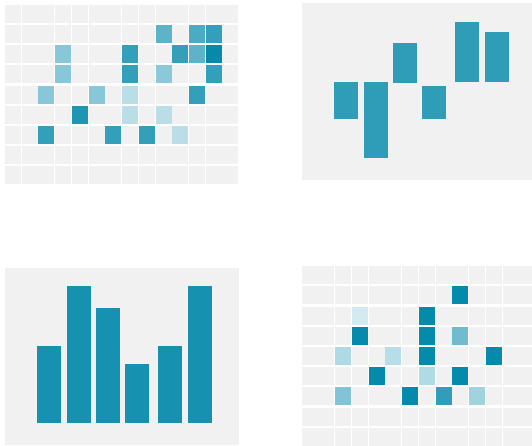# BENCHMARKING IDE (BIDE)

**DB Community**

▸ existing benchmarks (e.g., TPC-H)

▸ static, non-incremental queries

▸ report the performance of single queries

**Viz Community**

▸ user studies for evaluation

▸ use insight-based metrics

▸ costly, difficult to compare/reproduce

**BIDE**

▸ simulates common user behavior (think time,

  incremental queries, ...)

▸ introduces new metrics to better measure „interactivity"

**Website:** http://idebench.github.io
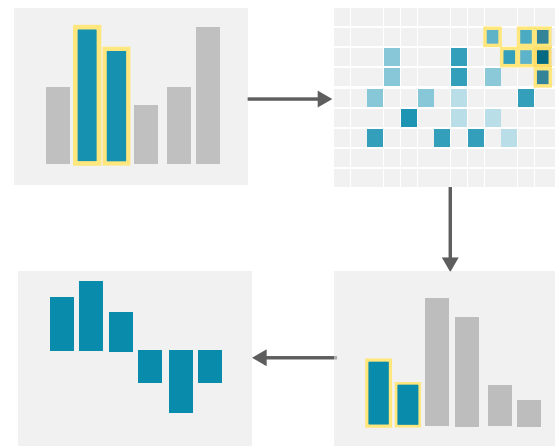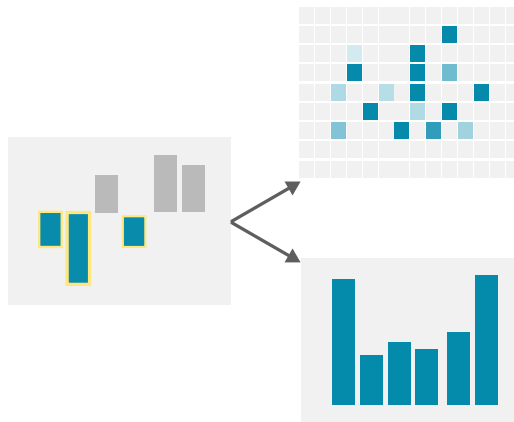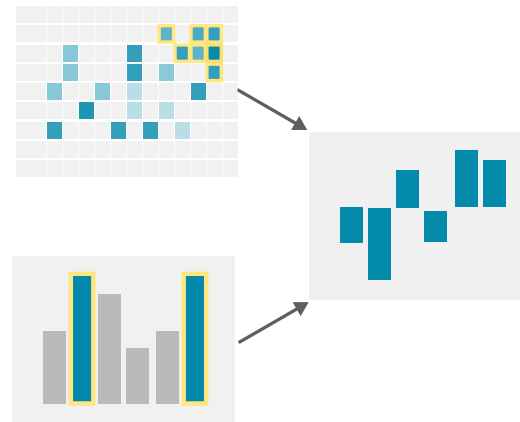
# BENCHMARK WORKLOAD



a) Independent Browsing

b) Sequential Linking

c) 1:N Linking

d) N:1 Linking

# WORKLOAD: OTHER DIFFERENCES

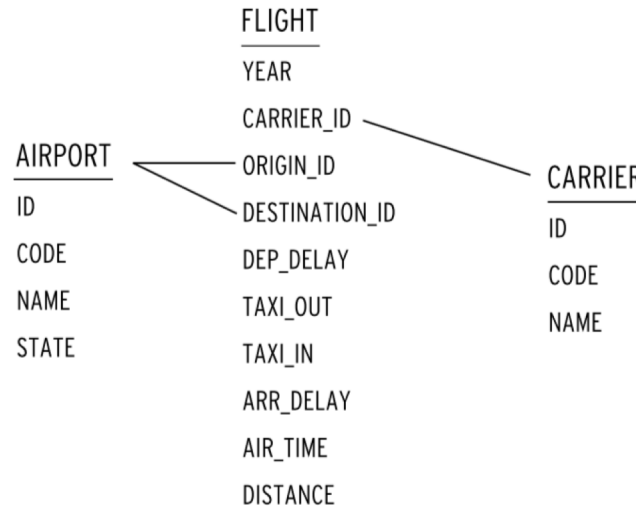**Multiple Concurrent Queries (triggered by one UI interaction)**



**Visualization-Specific Functions (e.g., Binning, Cross-Filter)**

**Other Parameters (Think Time, …)**

# BENCHMARK DATA SETS

**IDEBench comes with real-world data sets (e.g., Airline)**

FLIGHT
YEAR
CARRIER_ID
ORIGIN_ID
DESTINATION_ID
DEP_DELAY
TAXI_OUT
TAXI_IN
ARR_DELAY
AIR_TIME
DISTANCE

AIRPORT
ID
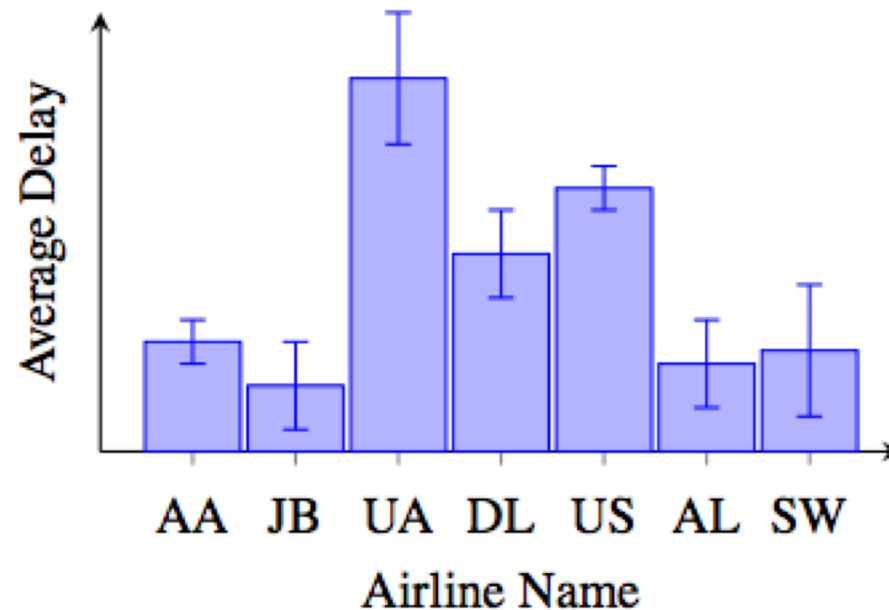CODE
NAME
STATE

CARRIER
ID
CODE
NAME

## Data Generator:

- **Supports different schema variants (normalized vs. denormalized)**
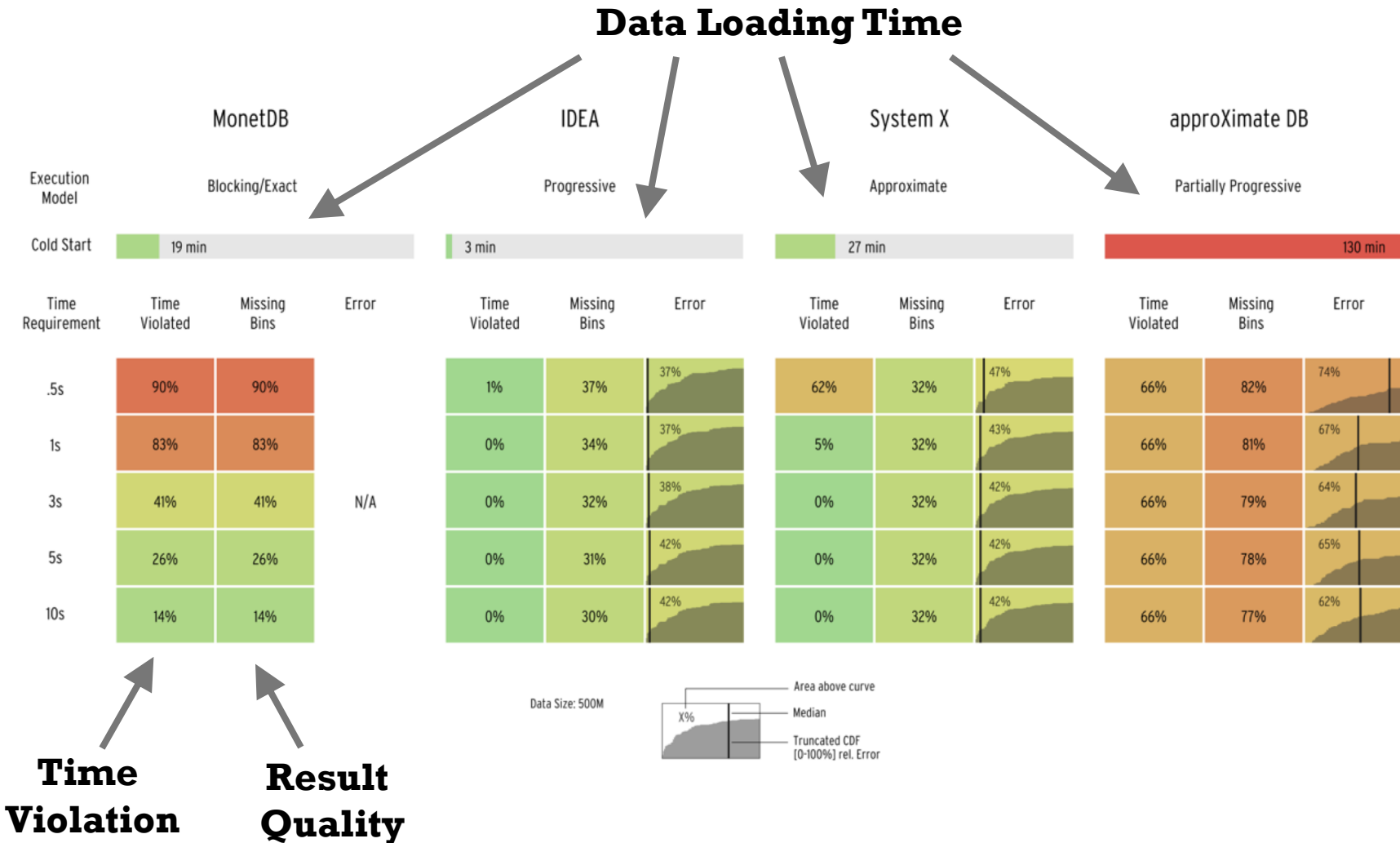
- **Can be used to scale-up and down data sets**

# BENCHMARK METRICS

**Result quality (error, completeness) and time** are both important metrics for such a benchmark



**Main Idea:** Capture quality of result after time t

# BENCHMARKED SYSTEMS

| Classical Analytical DBMSs | Approximate DBMSs | Specialized IDE Engines |
|---|---|---|

*System X*



*System Y*

# REPORTED RESULTS

**Data Loading Time**

**Time Violation**

**Result Quality**

# SUMMARY

**Interactive Data Exploration is challenging**

**We need to rethink the full data exploration stack**

- Query Interfaces

- Query Execution

- Cleaning / Loading

**Other Considerations:**

- Complex Workloads (ML, Text, …)

- Hardware Acceleration

- …

# COLLABORATORS