# IIT DBGroup
## Research Profile

**Boris Glavic**

# Outline

1) **Who am I?**

2) Provenance

3) Next-generation DB Engines

4) Foundations for Data Science
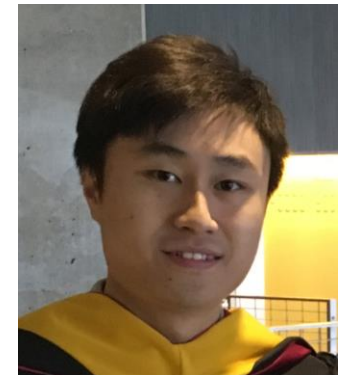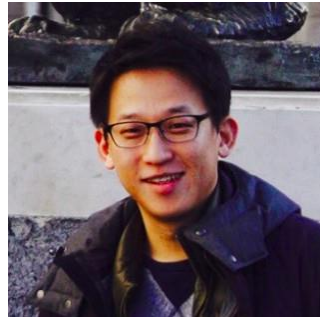
**2**

# Boris Glavic

I am a **database** guy!

Hi, I am **Boris Glavic,**
**Professor** in
**CS**

**3**

# IIT DBGroup

- ## **My research group**
  - ### [http://www.cs.iit.edu/~dbgroup/](http://www.cs.iit.edu/~dbgroup/)

**4**

# IIT DBGroup
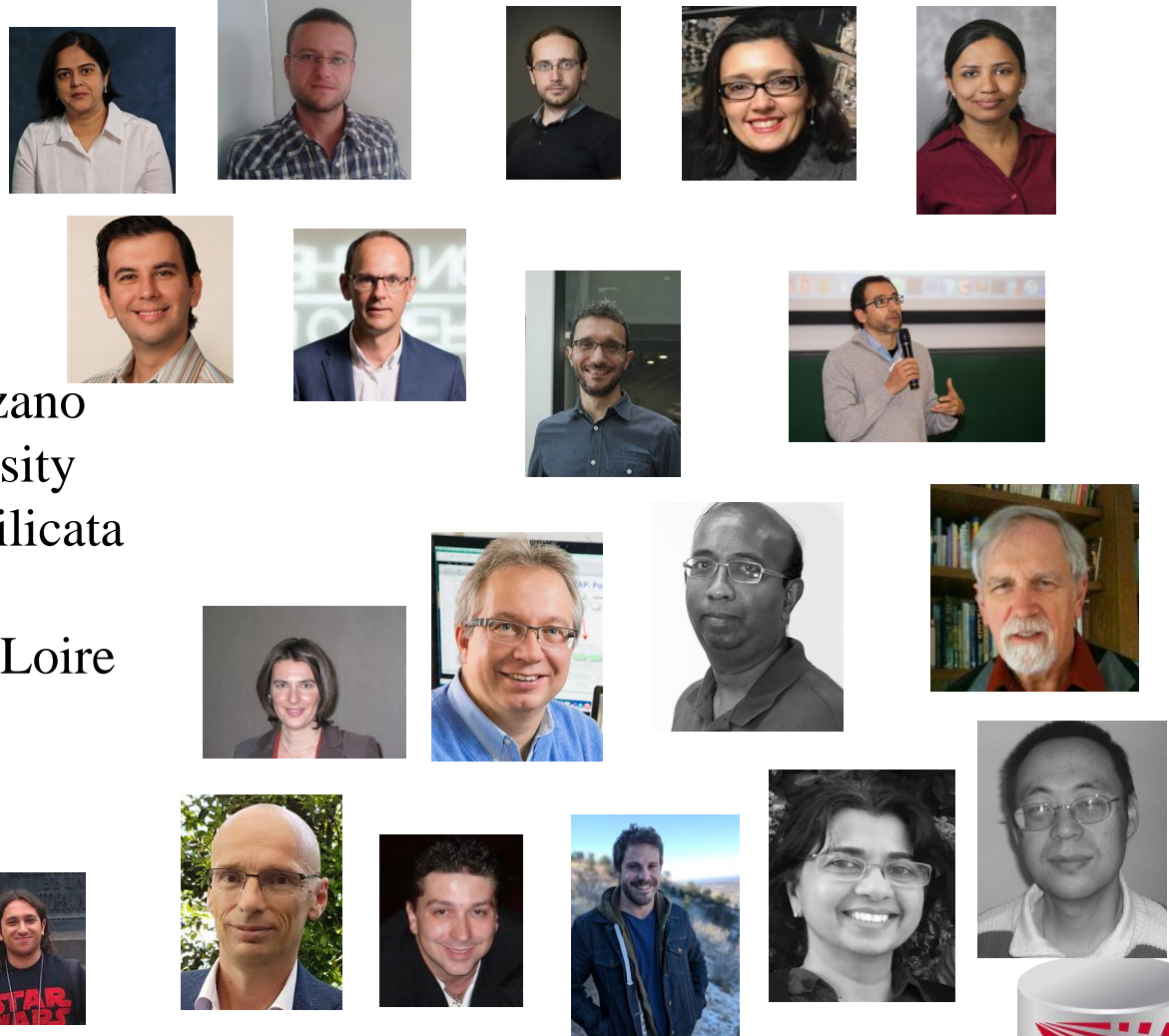
- **Collaborators**
  - Duke
  - SUNY Buffalo
  - UIUC
  - Oracle
  - NYU
  - Free University Bolzano
  - Northeastern University
  - Università della Basilicata
  - EURECOM
  - INSA Centre Val de Loire
  - University of Zurich
  - IIT
  - DePaul
  - **You?**

**5**

# Outline

1) Who am I?

2) **Provenance**

3) Next-generation DB Engines

4) Foundations for Data Science

**6**

# What the heck is provenance?

- **Provenance in art**
  - Record of the ownership (history) of a piece of art



Jan Van Eyck - Arnolfini Portrait

| 1523-4 | In another Mechelen inventory, a similar description, this time the name of the subject is given as "Arnoult Fin". |
|---|---|
| 1558 | In 1530 the painting was inherited by Margaret's niece Mary of Hungary, who in 1556 went to live in Spain. It is clearly described in an inventory taken after her death in 1558, when it was inherited by Philip II of Spain. A painting of two of his young daughters commissioned by Philip clearly copies the pose of the figures (Prado).[1] |

# What the heck is provenance?

- **Provenance in databases**

## Given a piece of data

- How do we know ...
  - which data it is derived from?
  - which transformations (SQL) where used to create it?
  - who created it?
  - ...

## Definition (Data Provenance)

Information about the origin and creation process of data.

## Example

result

|       | shop   | rev |
|-------|--------|-----|
| $t_1$ | Migros | 125 |
| $t_2$ | Coop   | 25  |

↑

```
SELECT shop,
       sum(price) AS rev
FROM sales, items
WHERE itemId = id
GROUP BY shop
```

↑                    ↑

sales

|       | shop   | itemId |
|-------|--------|--------|
| $s_1$ | Migros | 1      |
| $s_2$ | Migros | 3      |
| $s_3$ | Coop   | 3      |

items

|       | id | price |
|-------|----|-------|
| $i_1$ | 1  | 100   |
| $i_2$ | 2  | 10    |
| $i_3$ | 3  | 25    |

**8**

# Why should I care?

- **Auditing + Forensics**
- **Trust + Tracebility**
- **Proof of ownership**
- **Debugging**
- **Tracking uncertainty/probability**
- **Understanding data**
- **Relevance-based data management**
- **…**

**9**

- **Tracking provenance for updates+queries**
  - **reenactment**
- **Optimizing provenance capture and management**
- **Unifying provenance and missing answers**
- **Relevance-based data management**
- **GProM => Long term systems effort**

# Provenance for updates

- **Databases are typically are not static, but updated over time**
  - **Need to understand how …**
    - current state of the data was derived from old state
    - which operations affected the data
    - **=> provenance for updates!**

# Motivation

- **Understand how data is derived by concurrent transactions**

  - **Debugging transactions**
    - Akin to debugging concurrent programs

  - **Auditing data and Forensics**
    - Prove how this row in my database came to be?
    - Which operations of which transactions affected it?

  - **Data Curation and Integration**
    - From which source is this value derived?
    - Can I trust this result?
    - How did my cleaning heuristic interacts with my manual modifications

- **… this is post mortem provenance for transactional histories**

  - **Data dependencies** – this row is derived from these rows
  - **Operations** – these operations did affect this row
  - **Intermediate states** – state of relation R visible to operation

# Challenges

- **C1 - What is the right model for transaction provenance?**
  - Model data dependencies across database versions
  - Model control dependencies at granularity of updates
  - Account for low isolation levels
- **C2 - How to compute provenance for transactions efficiently?**
  - Efficient retrieval of provenance for a transactional history (or parts thereof)
  - Low overhead for operations if no provenance is requested
  - Tolerable storage requirements
  - Non-invasive possible (no changes to workload + system)?
- **C3 - How to represent and query provenance?**
  - Represent provenance in a human readable format
  - Provide powerful query capabilities for provenance

# Our approach

- **C1 – MV-semirings – a provenance model for queries and transactions**

  - A provenance model for transactions

  - Each row **is annotated with a symbolic expression** that encodes

    - From which previous row versions it was derived

    - How these row versions were combined by the derivation

    - Which updates of which transactions were involved in the derivation

  - Backwards compatible to **provenance semirings for queries**

    - the "most general" provenance model for a large class of queries

**14**

ILLINOIS INSTITUTE
OF TECHNOLOGY

**Account**

$C^1_{T_6,16}(U^1_{T_6,12}(C^1_{T_1,4}(I^1_{T_1,2}(x_1))))$

$C^2_{T_5,14}(U^2_{T_5,11}(C^2_{T_1,4}(I^2_{T_1,3}(x_2))))$

$C^3_{T_5,14}(U^3_{T_5,13}(U^3_{T_5,11}(C^3_{T_2,3}(I^3_{T_2,2}(x_3)))))$

| cust | typ | bal |
|------|-----|-----|
| Alice | Checking | -1100 |
| Alice | Savings | 1100 |
| Peter | Savings | 5390 |

**15**

# Our approach

- **C2 – Reenactment – replaying transactional workloads using queries**
  - Translate a **transactional history** (or parts thereof) into a **query** that is **equivalent under annotated semantics** (MV-semirings)
    - Produces the **same result**
    - Has the **same provenance** (data and control dependencies)
    - **Simulates interactions among transactions**
- **C3 – Implementing reenactment though SQL queries**
  - User-friendly **relational encoding** of our provenance model
  - Compile a **reenactment query into an SQL query** that produces this encoding
    - Uses **audit logging** and **time travel** (natively supported by many DBMS)
  - Make provenance requests an **SQL query feature**

# Advantages

- **Non-invasive:** requires no modifications to
  - DBMS
  - Transactional workload
- **Low overhead for your workload**
  - No provenance materialized upfront
  - Only overhead payed is ~20% for write only workloads
    - Time travel
    - Audit log
- **Efficient retroactive reconstruction of provenance (or hypothetical scenarios)**
  - One query -> optimize across statements and transactions
  - Replay does not result in disk writes
- **Fully integrated with SQL querying**
  - SELECT count(*)
  - FROM PROVENANCE OF TRANSACTION '1234454234'
  - WHERE ….

**17**

# Example application

- **Post-mortem debugger for transactions**



https://github.com/IITDBGroup/GProMTransactionDebugger

- **Generic middleware for implementing provenance functionality on top of DBMS**
  - **Query-to-query compiler**
    - Take queries with provenance features and compile them into SQL
  - **Modular system**
    - Pipelines consisting of multiple compilation steps
  - **Versatile**
    - All features mentioned in the beginning implemented in GProM
    - Useful beyond provenance: uncertainty, temporal queries …
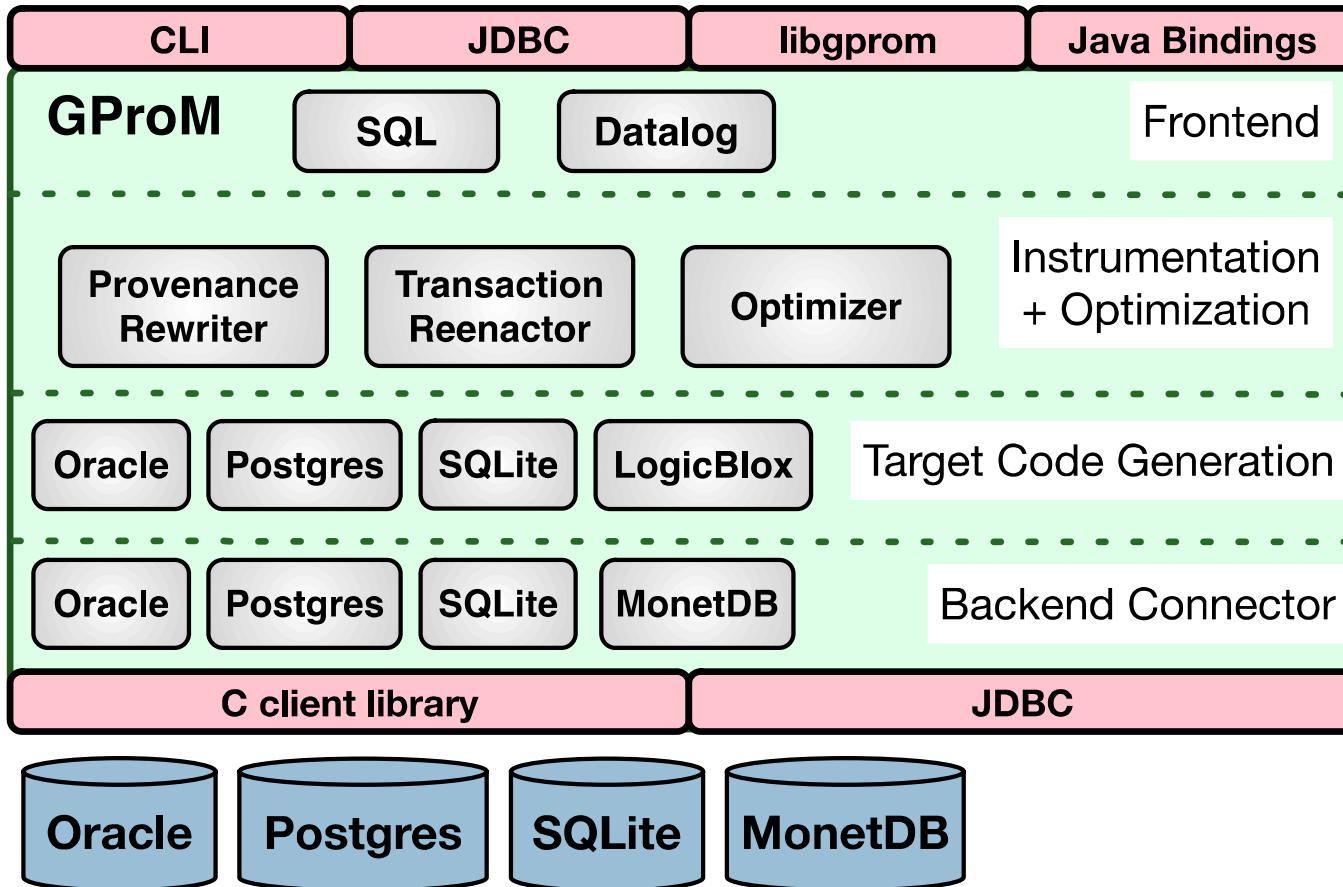
# GProM - Architecture

PROVENANCE OF
(SELECT * FROM ...

Q(X) :- R(X,Y).
WHY(Q(Peter))

https://github.com/IITDBGroup/gprom

| CLI | JDBC | libgprom | Java Bindings |
|-----|------|----------|---------------|

**GProM**

| SQL | Datalog | Frontend |
|-----|---------|----------|

| Provenance Rewriter | Transaction Reenactor | Optimizer | Instrumentation + Optimization |
|---------------------|----------------------|-----------|-------------------------------|

| Oracle | Postgres | SQLite | LogicBlox | Target Code Generation |
|--------|----------|--------|-----------|------------------------|

| Oracle | Postgres | SQLite | MonetDB | Backend Connector |
|--------|----------|--------|---------|-------------------|

| C client library | JDBC |
|------------------|------|

| Oracle | Postgres | SQLite | MonetDB |
|--------|----------|--------|---------|

**20**

# Outline

1) Who am I?

2) Provenance

3) **Next-generation DB Engines**

4) Foundations for Data Science
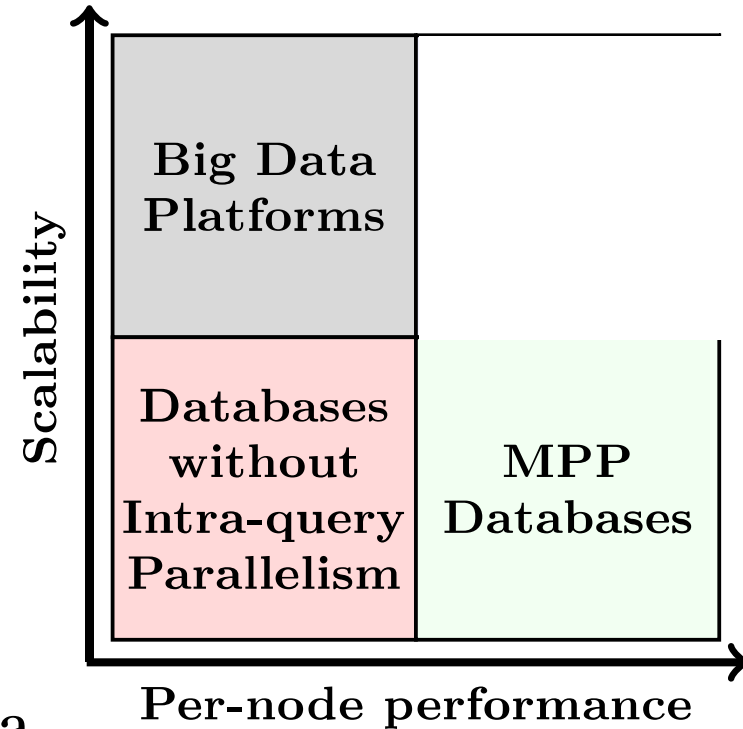
# Building Next-gen DB Engines

- **How to make databases more scalable without sacrifizing performance?**
  - **HRDBMS**
    - In-house distributed database
    - Scales like Big Data platforms, performs like databases
- **How to maximize performance?**
  - **NautDB** – Generality + Performance
    - Custom dataflow engine / kernel hybrid for compiled queries
  - Combine specialization techniques from the OS and DB community

**22**

# Scalable SQL Processing

- **Petabyte datawarehouses are becoming common**
  - eBay: 6.5PB (Greenplum) [2009]
  - eBay: 2.5PB (Teradata) [2009]
  - Facebook: 2.5PB (Hive) [2010]
  - Walmart: 2.5PB (Teradata) [2008]
  - Bank of America: 1.5PB (Teradata) [2008]
- **Meanwhile**
  - Traditional MPP DBMS are approaching their scalability limits
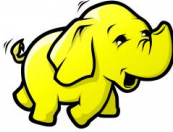    - e.g., Teradata (1024 nodes), Greenplum (1000 nodes), DB2 (1024 nodes)

# Dilemma

- **Big Data platforms provide scalability**
  - Hive
  - Spark SQL
  - Flink, AsterixDB, …

- **DBMS provide per-node performance**
  - DB2, Oracle SQL server
  - Teradata, Greenplum, Netazza
  - Hana, MonetDB, …

Scalability ↑

| Big Data Platforms | |
| --- | --- |
| **Databases without Intra-query Parallelism** | MPP Databases |

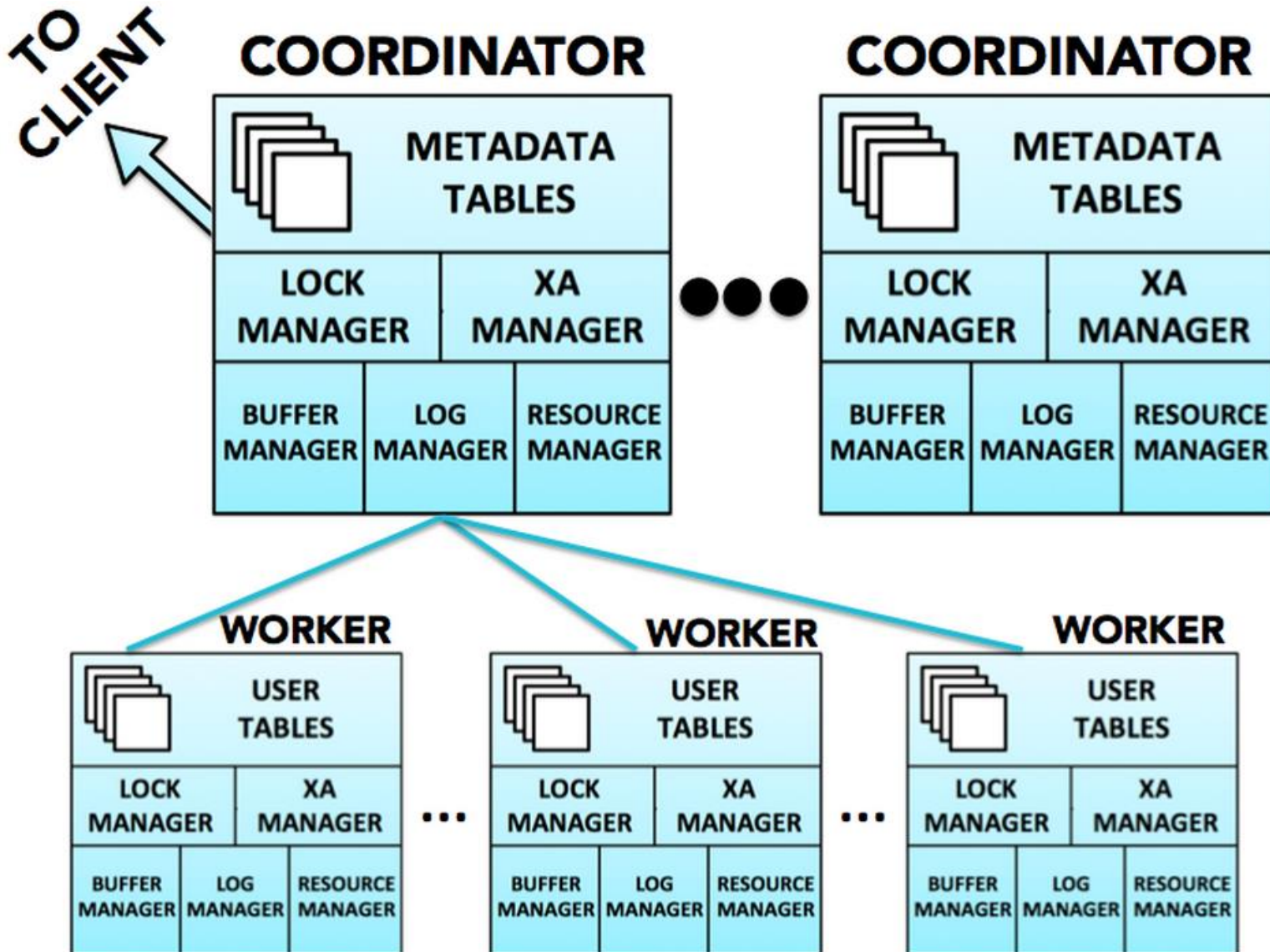Per-node performance →

**24**

# What is HRDBMS?

- **Distributed shared-nothing database**
  - Supports standard SQL queries
  - Support for serializable transactions
  - Runs on commodity hardware
  - Written in Java (~165K LOC)
- **Design goals**
  - **Scale** like an elephant (Hadoop)
    - i.e., Big data platform
  - **Perform** like a fruit (Greenplum)
    - i.e., traditional DBMS

# Approach – Pick the best

- **Combine the best of traditional disk-based DBMS and Big data platforms**

  - Distributed, asynchronous dataflow execution engine
  - Page-oriented storage
  - Buffer manager for in-memory cache
  - Pipelining
  - Intra and inter operator parallelism
  - Cost-based optimization based on statistics
  - Disk-resident index structures
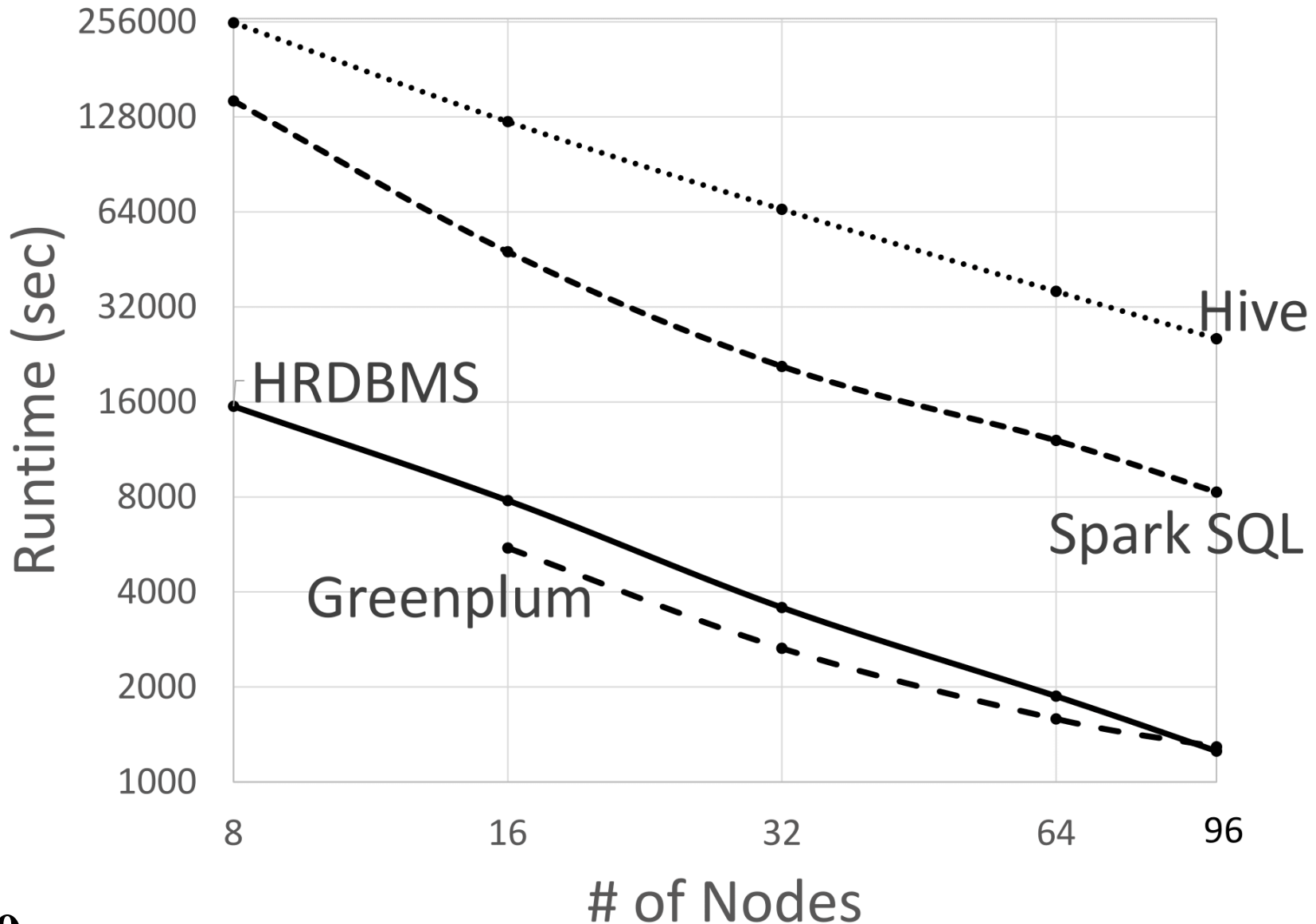  - Strong consistency (transactions)

27

- **Improve over both to avoid scalability and performance pitfalls**
  - Dedicated relational operator implementation in dataflow engine
  - Non-blocking shuffle with optional sort
  - Enforce scalable communication patterns
    - Limit #open-connections per node
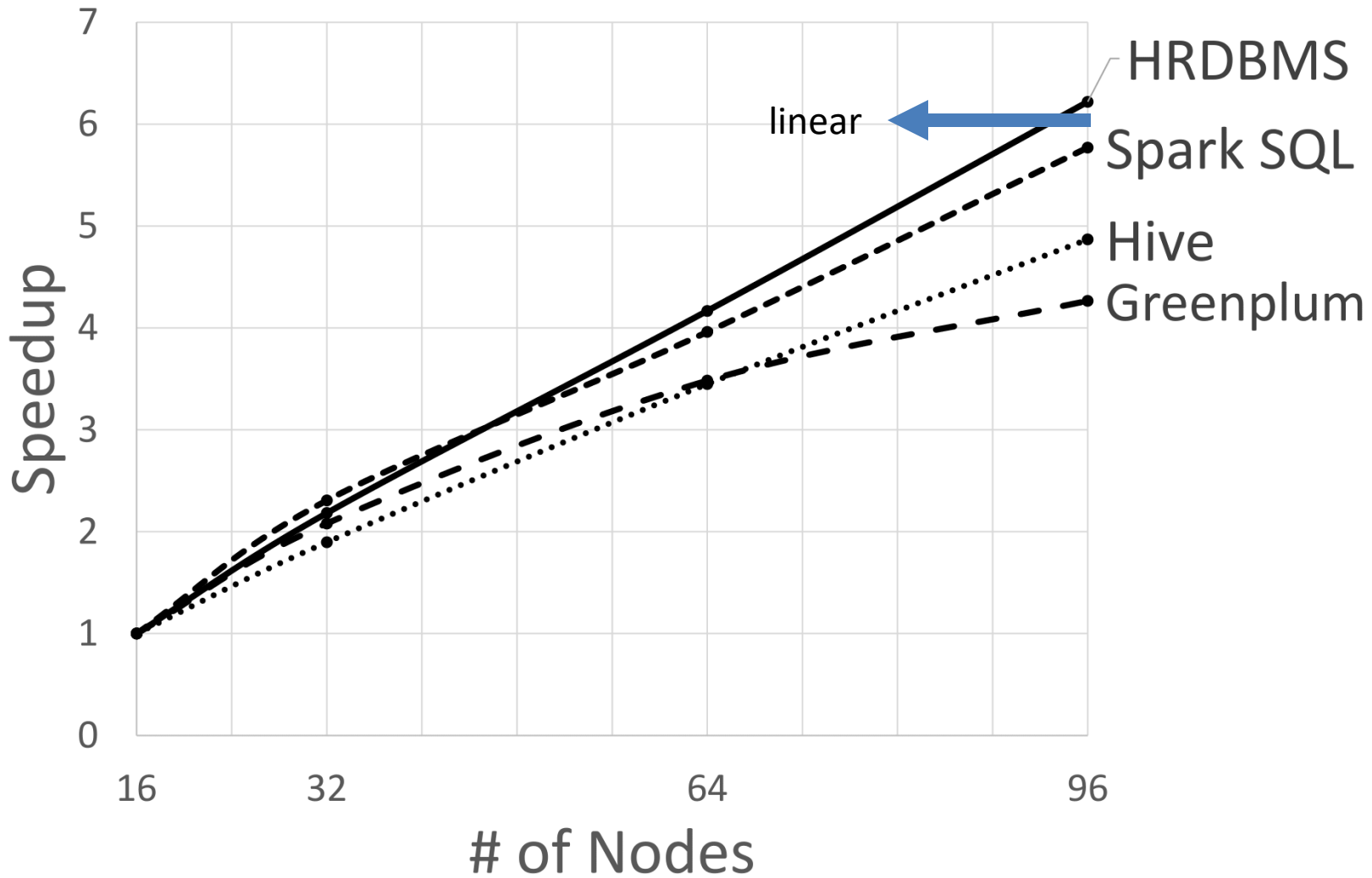  - Enforce data locality
  - Predicate cache for page-skipping

ILLINOIS INSTITUTE
OF TECHNOLOGY

- **Workload**
  - 1 TB TPC-H instance (SF=1,000)
- **Setup**
  - **Cooley cluster at Argonne (8 to 96 nodes)**
    - GPFS file system
  - **Nodes**
    - 2 x Intel Haswell E5-2620 v3 (12 cores total)
    - 386GB RAM (we use only 24GB)
- **Competitors**
  - Hive V1.2.1, Spark SQL V1.6.0, Greenplum V4.3.99

**Runtime (sec)** vs **# of Nodes**

- 256000, 128000, 64000, 32000, 16000, 8000, 4000, 2000, 1000 (y-axis)
- 8, 16, 32, 64, 96 (x-axis)

Hive

HRDBMS

Spark SQL

Greenplum

**30**

ILLINOIS INSTITUTE
OF TECHNOLOGY



HRDBMS

Spark SQL

Hive

Greenplum

linear

Speedup

# of Nodes

**31**

# Conclusions

- **Goal**
  - Scale like an elephant, perform like a fruit!

- **Assessment**
  - Goal achieved!

- **Take-away message(s)**
  - Don't disregard traditional methods
    - Stealing = good
  - Don't be afraid of new ideas
  - Systems research does the trick!
    - Build and evaluate large-scale systems!

**32**

# Questions?

- **Homepage:**

**Jason:** http://www.cs.iit.edu/~dbgroup/people/jarnold.php

**Boris:** http://www.cs.iit.edu/~glavic/

**Ioan:** http://www.cs.iit.edu/~iraicu/

- **HRDBMS Project**

http://www.cs.iit.edu/~dbgroup/research/hrdbms.php

- **Github**

https://github.com/IITDBGroup/HRDBMS

# Outline

1) Who am I?

2) Provenance

3) Next-generation DB Engines

4) **Foundations for Data Science**

# Data Science (what you think)

- **1. Collect data**
- **2. Ask questions**
- **3. Get answers (and more questions)**
- **4. Goto 2**

- **1. Collect data**

- **2. Spend weeks curating your data**

- **3. Discover data bugs (goto 2 until clean)**

- **3. Ask questions**

- **4. Discover more data bugs (goto 2)**

- **5. Get answers (and more questions)**

- **6. Realize that you cannot trust the answers because you are not sure whether curation introduced new errors => goto 2?????**

# Data Science (real world)

- **80% of time spend on data curation/collection/preparation**

- **This process is currently not supported as much as it should be**

- **Current tools do not automate subtasks (fix certain types of problems)**

- **Current tools do not help users to collaborate and track problems**

# My research

- **Data cleaning + integration**
  - There are many automated tools for data cleaning and curation
    - Help to evaluate them with **BART** and **iBench**

- **Uncertainty management**
  - How to track how uncertainty inherent in data or introduced by (semi-)automated curation through an analysis

- **Vizier**
  - Data curation and exploration platform

# Vizier

- **http://www.vizierdb.info/**
- **Data exploration and curation platform**
- **Multi-model Interface**
- **Best-Effort Data Ingest**
- **Data Debugging + Tracebility + Uncertainty Management**
- **Versioned Notebook**
- **Collaboration**
- **Scalable and Compatible**

# Multi-modal Interface

- **Notebook (think Jupyter)**
  - Python
  - SQL
- **Spreadsheet**
- **Graphs**







**40**

- **Full history of your notebook**
- **Branching, checkpointing, and time travel**
- **Share and release versions with full edit history**
- **Annotate data**
  - Annotations travel along with the data!

# Uncertainty Tracking

- **Automated curation is often heuristic in nature**
  - E.g., impute missing values by replacing them with values predicted by a classifier (ML model)

- **Heuristics make analysis results inherently untrustworthy!**

- **Vizier knows how automated curation methods introduce uncertainty**
  - Uncertainty is propagated through operations

**42**

# Conclusions

- **Brief overview of my research**
  - Want to know more? Just reach out
  - I will share these slides

- **Links**
  - https://github.com/IITDBGroup/
  - http://www.cs.iit.edu/~dbgroup/research/
  - http://cs.iit.edu/%7edbgroup/people/index.html
  - http://cs.iit.edu/%7edbgroup/publications/index.php

**43**