

# Data-Intensive Scalable Science: Beyond MapReduce

Bill Howe, UW

...plus a bunch of people





**<http://escience.washington.edu>**



escience institute

Search

[Advanced Search](#)

Web [+ Show options...](#)

Results 1 - 10 of about 908,000 for **escience institute** (0.40 seconds)

Did you mean: [e science institute](#)

[University of Washington eScience Institute](#)

The University of Washington eScience Institute was formed in January 2008 by the Vice President for UW Technology and the Provost, in partnership with key ...

[escience.washington.edu/](#) - [Cached](#) - [Similar](#) -

[University of Washington eScience Institute](#)

For more information about the Institute, or for help with your data-intensive research, please contact: [info escience.washington.edu](#) ...

[escience.washington.edu/content/](#) - [Cached](#) - [Similar](#) -

[+ Show more results from escience.washington.edu](#)

[National e-Science Centre](#)

e-Science Hub Kelvin Building Glasgow G12 8QQ United Kingdom Tel: 0141 330 8606. Fax:

0141 330 4913 email, [e-Science Institute](#) 15 South College Street ...

[Contacts](#) - [e-Science Events](#) - [About NeSC](#) - [Technical Papers](#)

[www.nesc.ac.uk/](#) - [Cached](#) - [Similar](#) -

[ALL RESULTS](#)**RELATED SEARCHES**[Science Museum](#)[Franklin Institute Science Museum](#)[Search Institute](#)[Science Center](#)[Skin Science Institute](#)[Reproductive Science Institute](#)[Health Sciences Institute](#)[Brick Computer Science Institute](#)**SEARCH HISTORY**[escience institute](#)[escience institute  
washington](#)

Were you looking for: [science institute](#)

[ALL RESULTS](#)1-10 of 87,700,000 results · [Advanced](#)**[Institute of Noetic Sciences: Home Page](#)**

IONS is a nonprofit membership organization that conducts and sponsors leading-edge research into the potentials and powers of consciousness: including perceptions, beliefs ...

[www.noetic.org](#) · [Cached page](#)**[Virginia Institute of Marine Science - Home](#)**

An **institution** that is uniquely prepared to educate the highly qualified researchers, resource managers, and educators needed for the future. Gloucester Point, Virginia.

[www.vims.edu](#) · [Cached page](#)**[ICSI: International Computer Science Institute | Berkeley, California](#)**

International Computer **Science Institute** at the University of California, Berkeley. Research groups and technical reports.

[http.icsi.berkeley.edu](#) · [Cached page](#)**[National Institute for Discovery Science \(NIDS\): Science institute ...](#)**

The National **Institute** for Discovery **Science** (NIDS) is a privately funded **science institute** engaged in research of UFOs, animal mutilations, and other related anomalous phenomena.

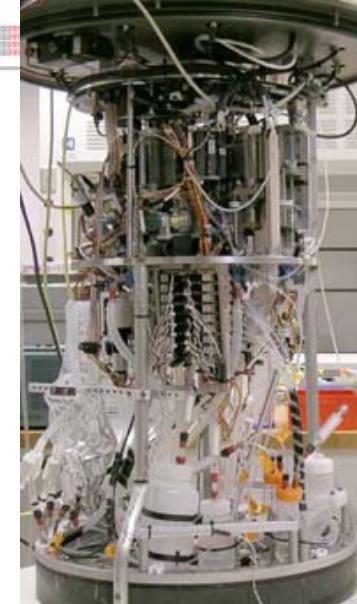
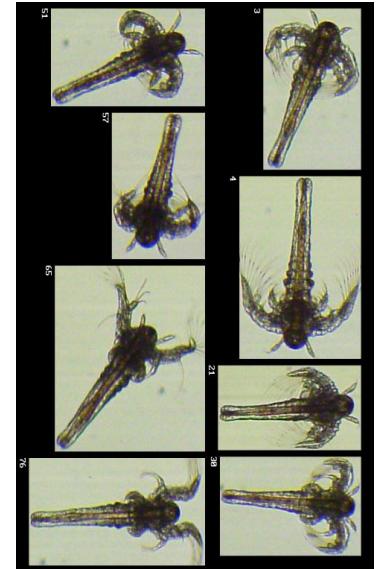
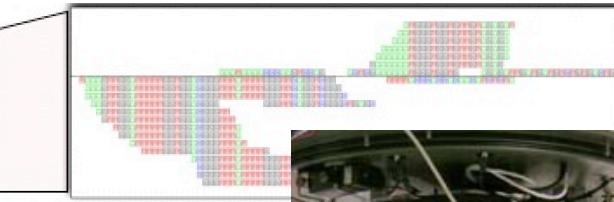
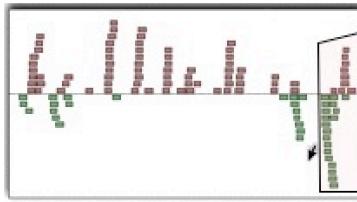
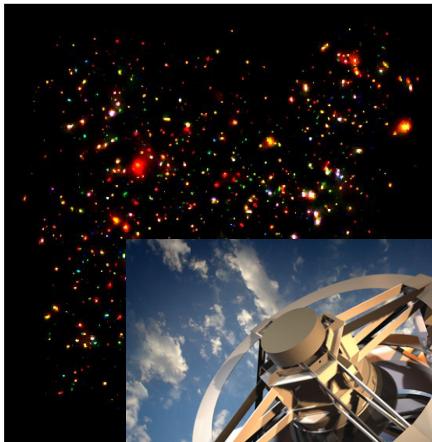
[www.nidsci.org](#) · [Cached page](#)

# Science is reducing to a database problem

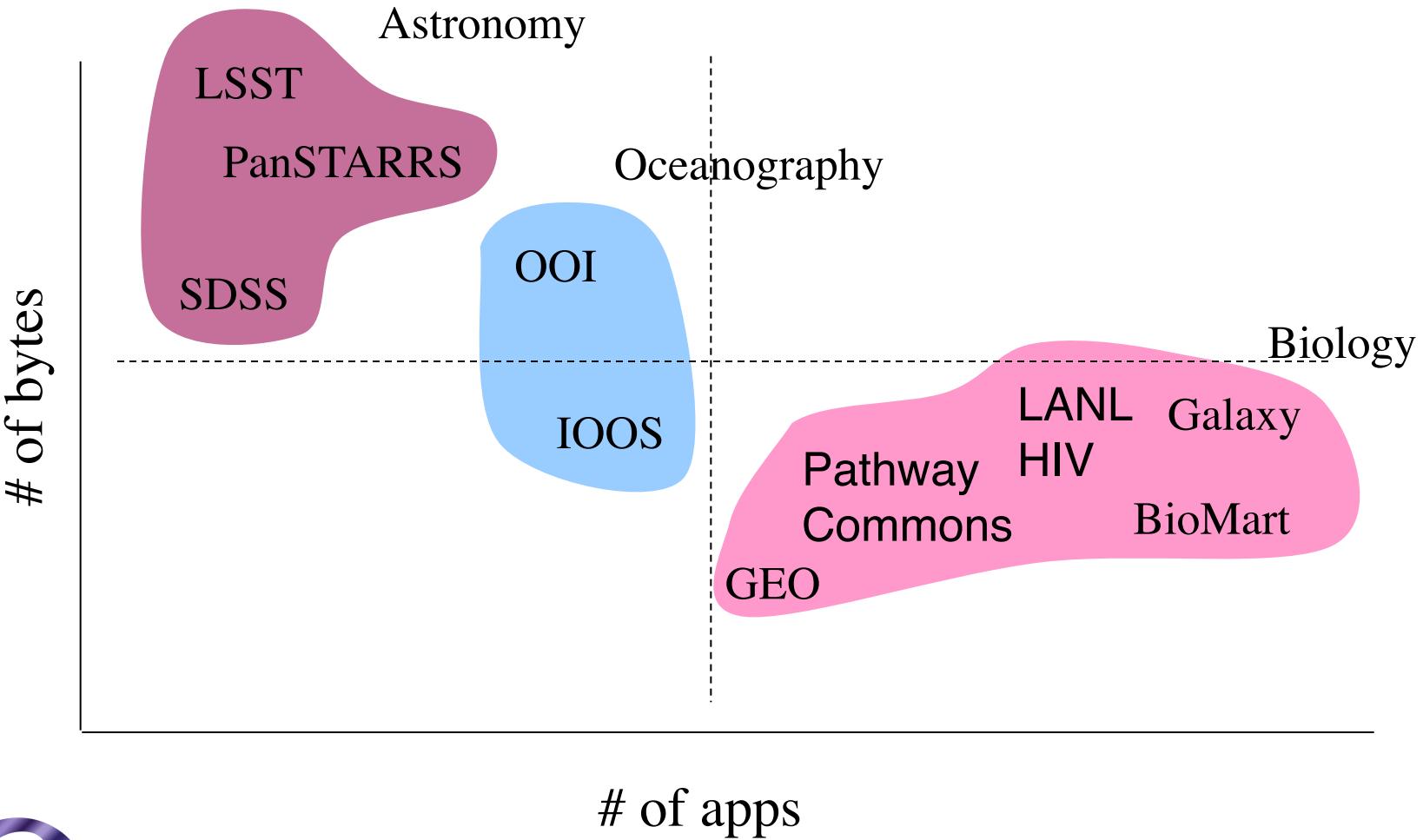
**Old model: “Query the world” (Data acquisition coupled to a specific hypothesis)**

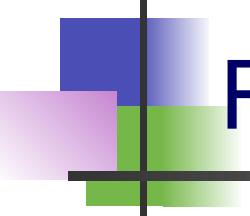
**New model: “Download the world” (Data acquired en masse, in support of many hypotheses)**

- Astronomy: High-resolution, high-frequency sky surveys (SDSS, LSST, PanSTARRS)
- Oceanography: high-resolution models, cheap sensors, satellites
- Biology: lab automation, high-throughput sequencing,



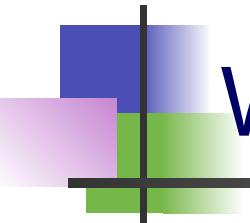
# Two dimensions





# Roadmap

- Introduction
- Context: RDBMS, MapReduce, etc.
- New Extensions for Science
  - Spatial Clustering
  - Recursive MapReduce
  - Skew Handling

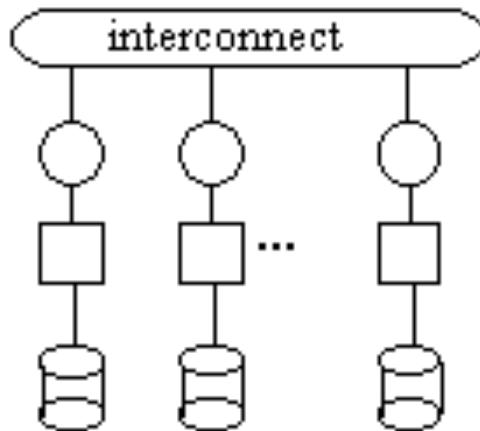


# What Does Scalable Mean?

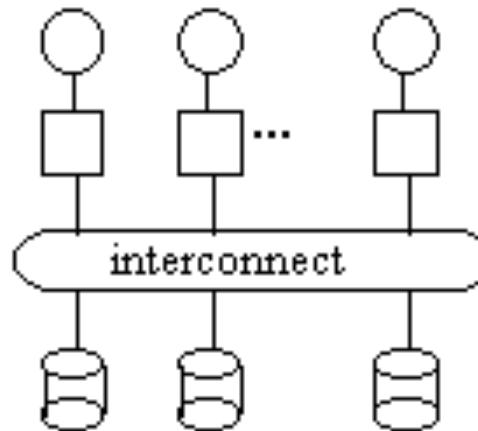
- In the past: Out-of-core
  - “**Works even if data doesn’t fit in main memory**”
- Now: Parallel
  - “**Can make use of 1000s of independent computers**”

# Taxonomy of Parallel Architectures

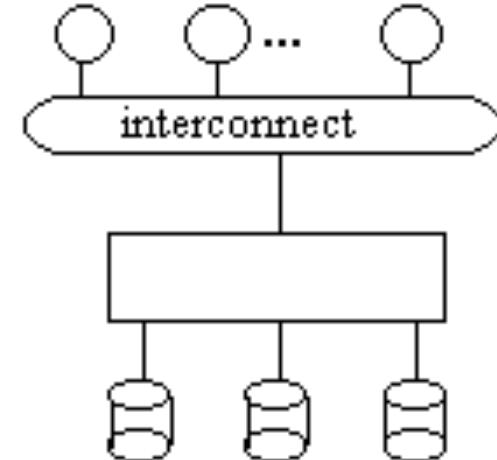
Scales to 1000s of computers



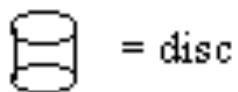
a) shared nothing



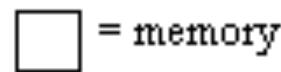
b) shared disc



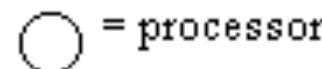
c) shared memory



= disc



= memory



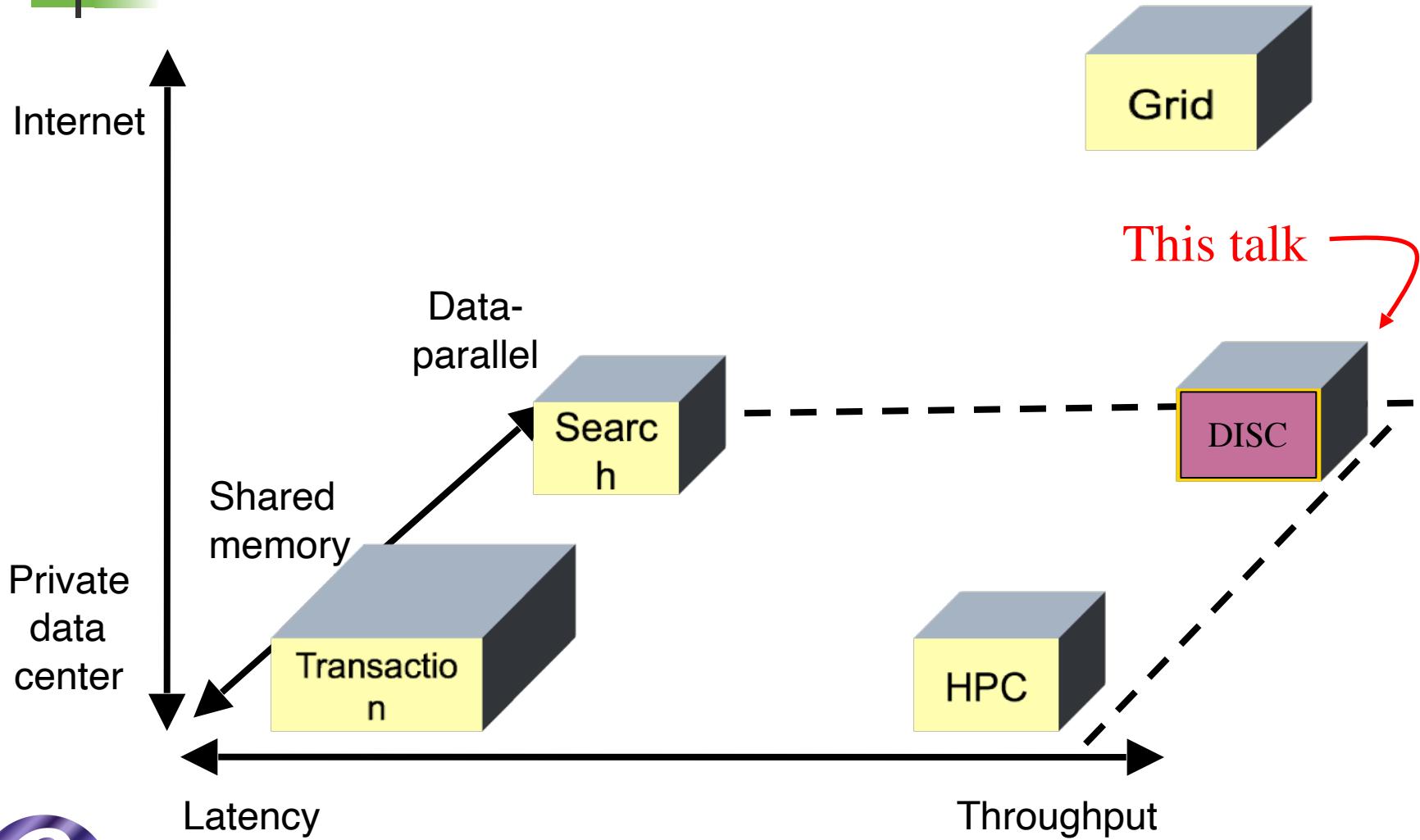
= processor

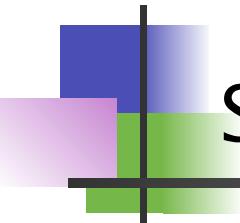
Easiest to program, but

\$\$\$\$

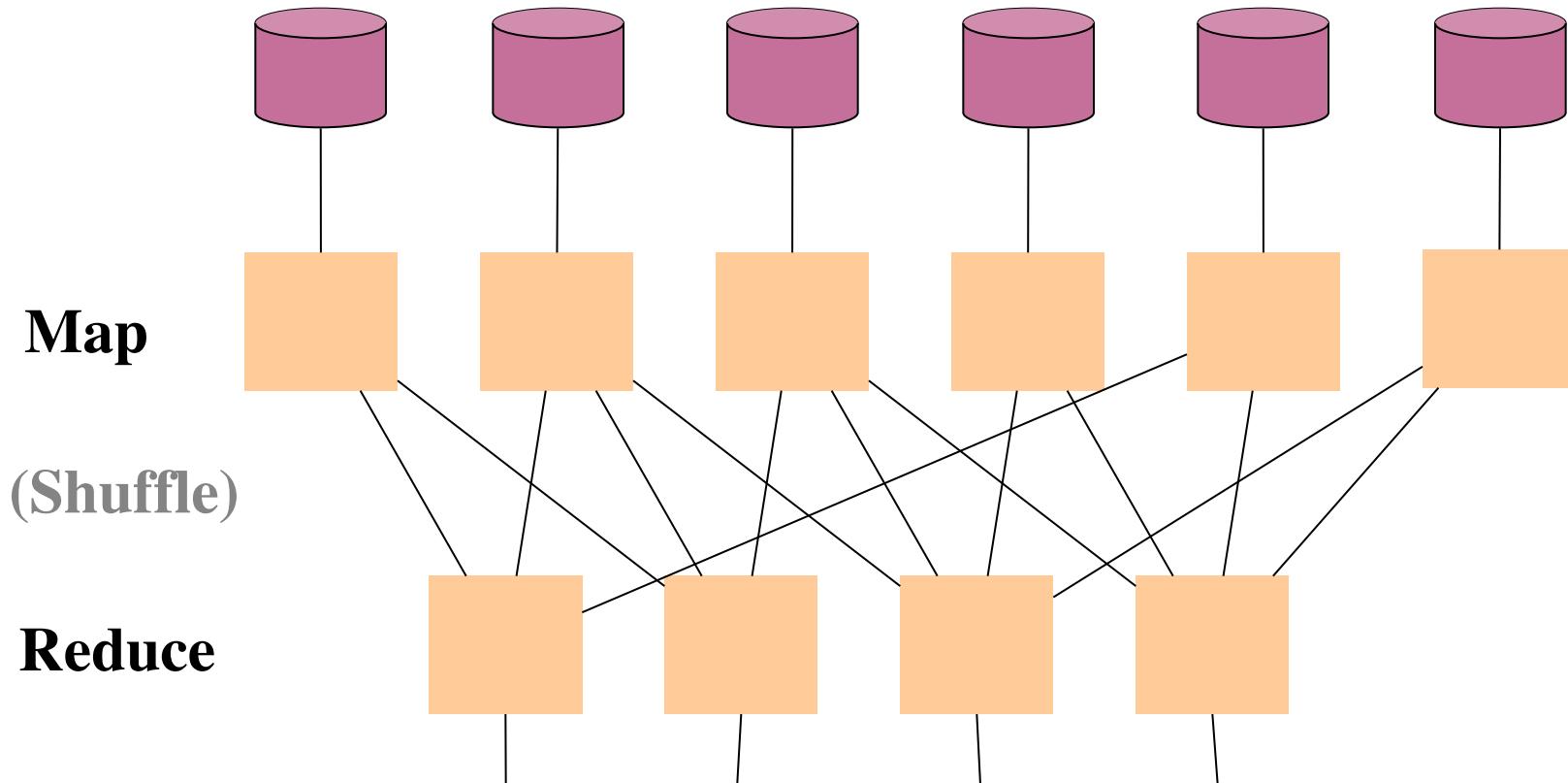
Fig. 3.1 Logical multi-processor database designs (diagram after [DEWI92])

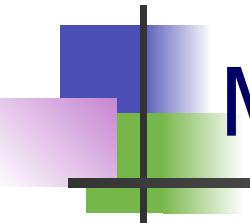
# Design Space





# Some distributed algorithm...





# MapReduce Programming Model

- Input & Output: each a set of key/value pairs
- Programmer specifies two functions:

**map (in\_key, in\_value) -> list(out\_key, intermediate\_value)**

- Processes input key/value pair
- Produces set of intermediate pairs

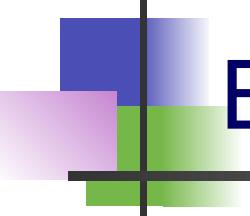
**reduce (out\_key, list(intermediate\_value)) -> list(out\_value)**

- Combines all intermediate values for a particular key
- Produces a set of merged output values (usually just one)

*Inspired by primitives from functional programming  
languages such as Lisp, Scheme, and Haskell*



slide source: Google, Inc.



# Example: What does this do?

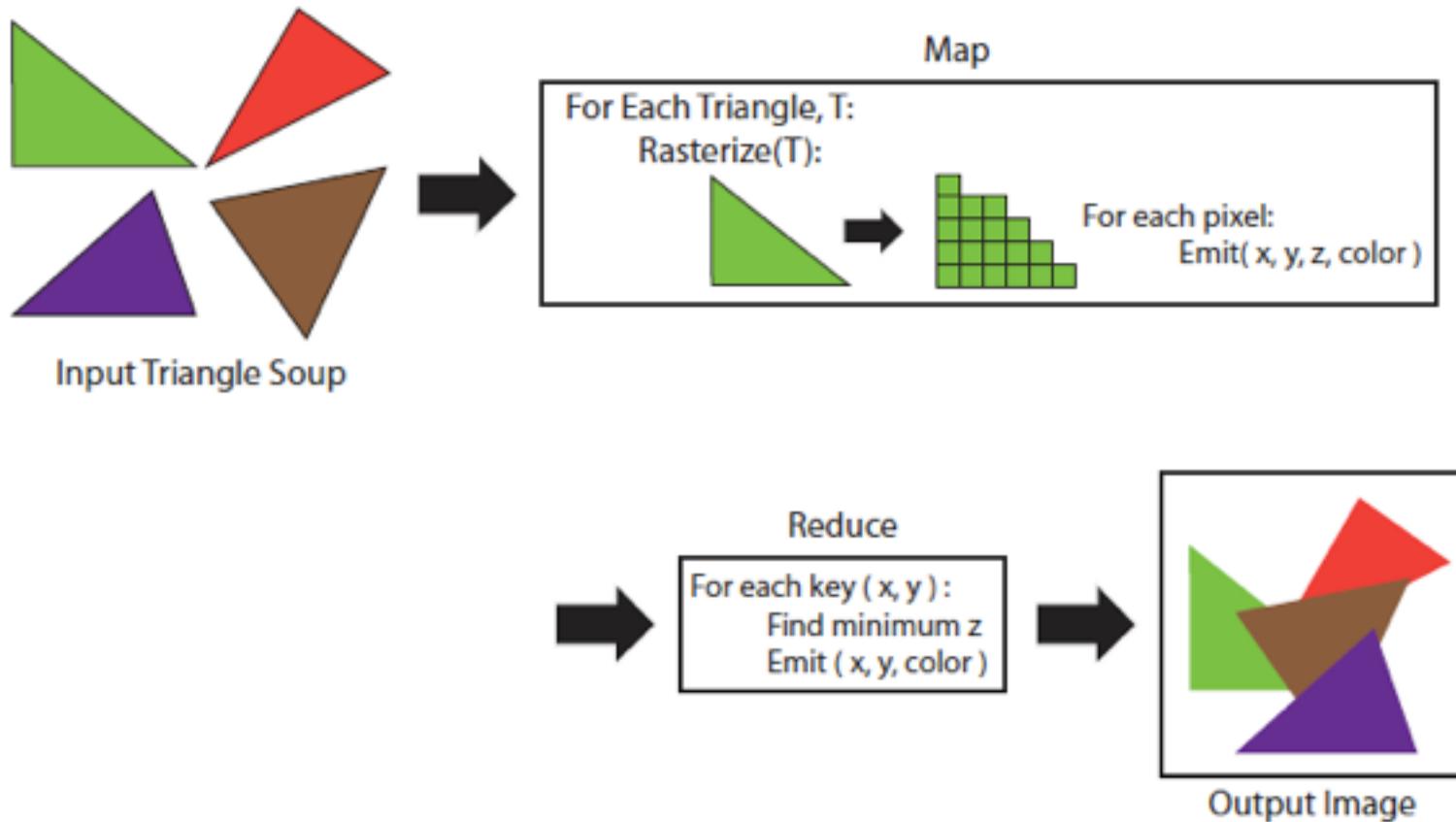
```
map(String input_key, String input_value):  
    // input_key: document name  
    // input_value: document contents  
    for each word w in input_value:  
        EmitIntermediate(w, 1);
```

```
reduce(String output_key, Iterator intermediate_values):  
    // output_key: word  
    // output_values: ????  
    int result = 0;  
    for each v in intermediate_values:  
        result += v;  
    Emit(result);
```



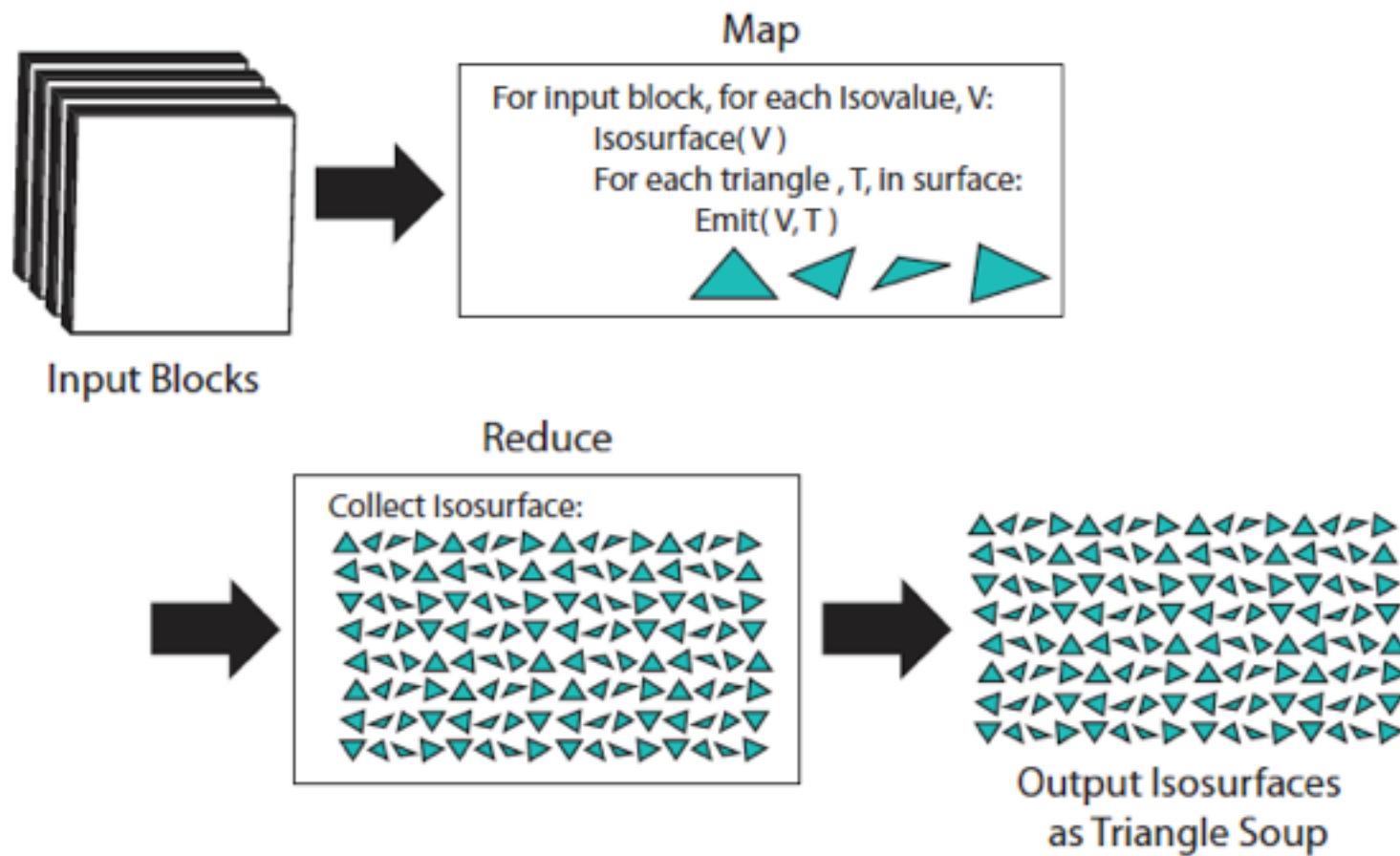
slide source: Google, Inc.

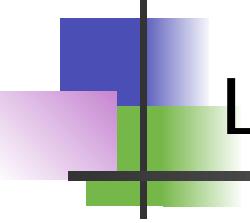
# Example: Rendering



Bronson et al. Vis 2010 (submitted)

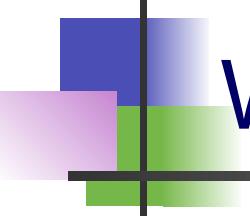
# Example: Isosurface Extraction





# Large-Scale Data Processing

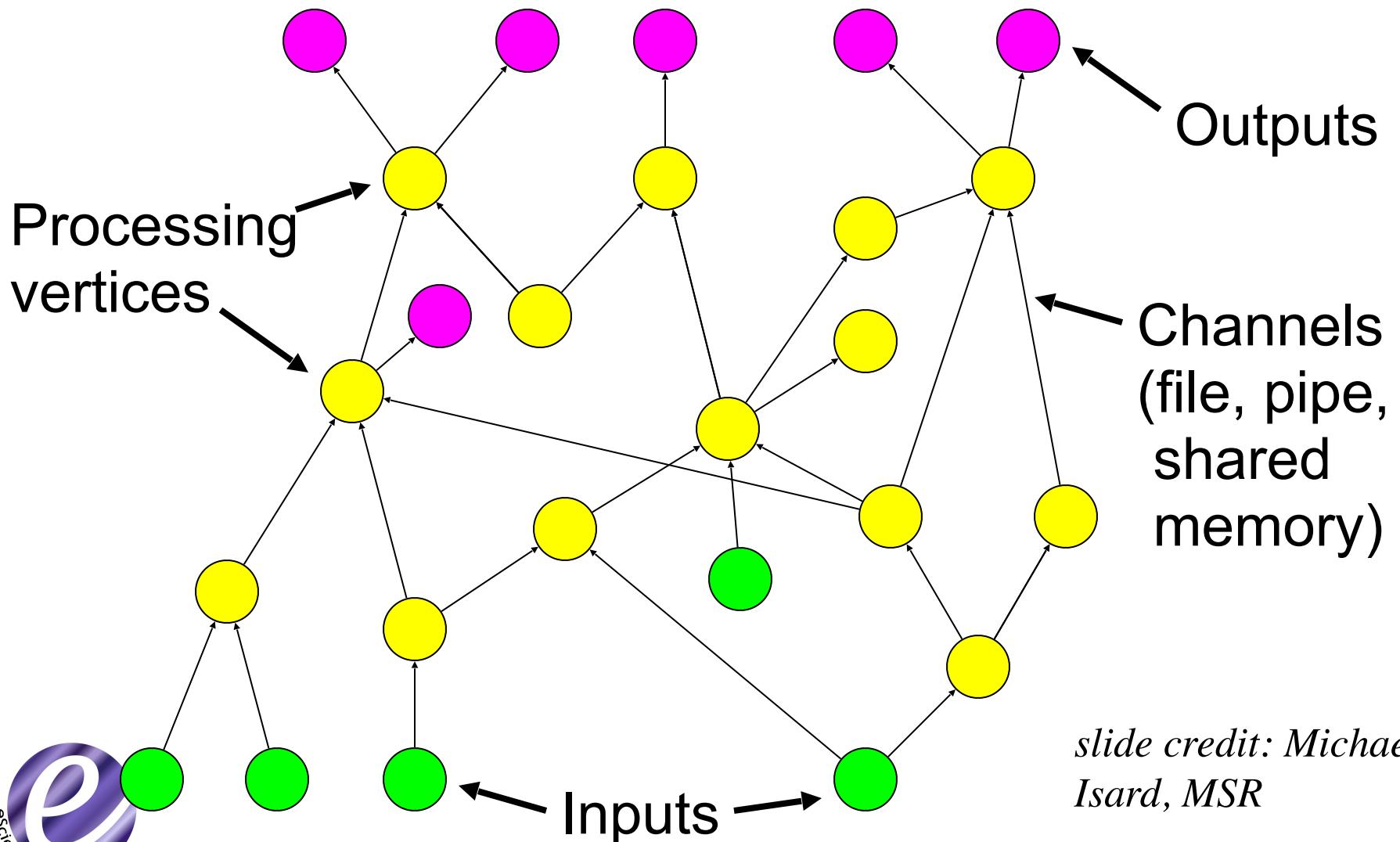
- Many tasks process big data, produce big data
- Want to use hundreds or thousands of CPUs
  - ... but this needs to be easy
  - **Parallel databases** exist, but they are expensive, difficult to set up, and do not necessarily scale to hundreds of nodes.
- MapReduce is a lightweight framework, providing:
  - **Automatic parallelization and distribution**
  - **Fault-tolerance**
  - **I/O scheduling**
  - **Status and monitoring**

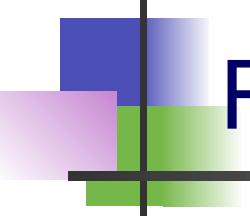


# What's wrong with MapReduce?

- Literally Map then Reduce and that's it...
  - Reducers write to replicated storage
- Complex jobs pipeline multiple stages
  - No fault tolerance between stages
    - Map assumes its data is always available: simple!
- What else?

# Realistic Job = Directed Acyclic Graph





# Relational Database History

Pre-Relational: if your data changed, your application broke.

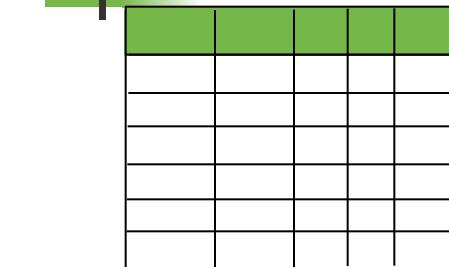
Early RDBMS were buggy and slow (and often reviled), but required only 5% of the application code.

*“Activities of users at terminals and most application programs should remain unaffected when the internal representation of data is changed and even when some aspects of the external representation are changed.”* -- Codd 1979

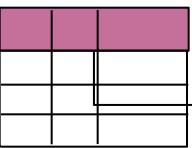
**Key Ideas:** Programs that manipulate tabular data exhibit an algebraic structure allowing reasoning and manipulation independently of physical data representation



# Key Idea: Data Independence



*logical data independence*



*physical data independence*

views

relations

files and  
pointers

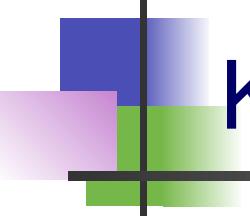
```
SELECT *
FROM my_sequences
```

```
SELECT seq
FROM ncbi_sequences
WHERE seq =
'GATTACGATATTA' ;
```

```
f = fopen('table_file');
fseek(10030440);
while (True) {
    fread(&buf, 1, 8192, f);
    if (buf == GATTACGATATTA) {
```

...  
...

Bill Howe, UW



# Key Idea: Indexes

- Databases are especially, but exclusively, effective at “Needle in Haystack” problems:
  - Extracting small results from big datasets
  - Transparently provide “old style” scalability
  - Your query will **always\*** finish, regardless of dataset size.
- Indexes are easily built and automatically used when appropriate

```
CREATE INDEX seq_idx ON sequence(seq);
```

```
SELECT seq  
      FROM sequence  
     WHERE seq = 'GATTACGATATTA';
```



\*almost

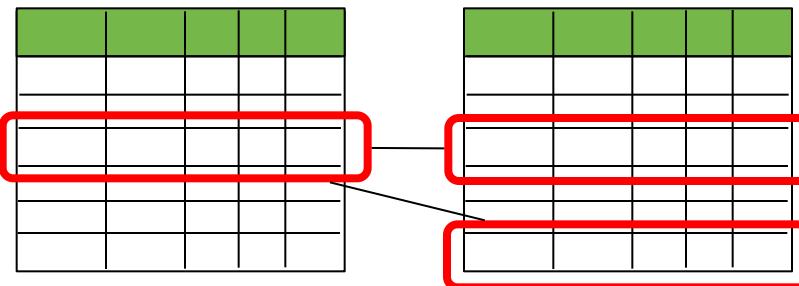
# Key Idea: An *Algebra of Tables*



select

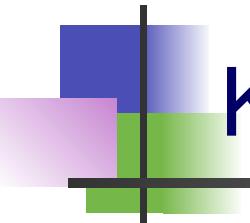


project



join

*Other operators: aggregate, union, difference, cross product*



# Key Idea: Algebraic Optimization

$$N = ((z^*2)+((z^*3)+0))/1$$

Algebraic Laws:

1. (+) identity:  $x+0 = x$
2. (/) identity:  $x/1 = x$
3. (\*) distributes:  $(n*x+n*y) = n*(x+y)$
4. (\*) commutes:  $x*y = y*x$

Apply rules 1, 3, 4, 2:

$$N = (2+3)*z$$

two operations instead of five, no division operator

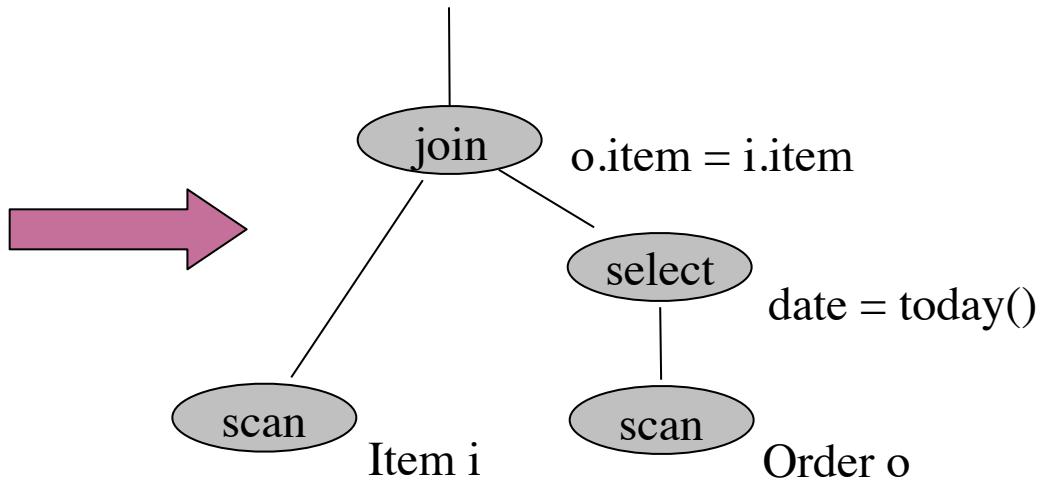
*Same idea works with the Relational Algebra!*

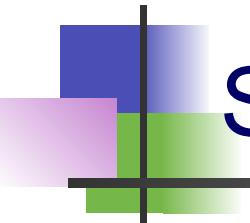


# Key Idea: Declarative Languages

*Find all orders from today, along with the items ordered*

```
SELECT *
  FROM Order o, Item i
 WHERE o.item = i.item
   AND o.date = today()
```

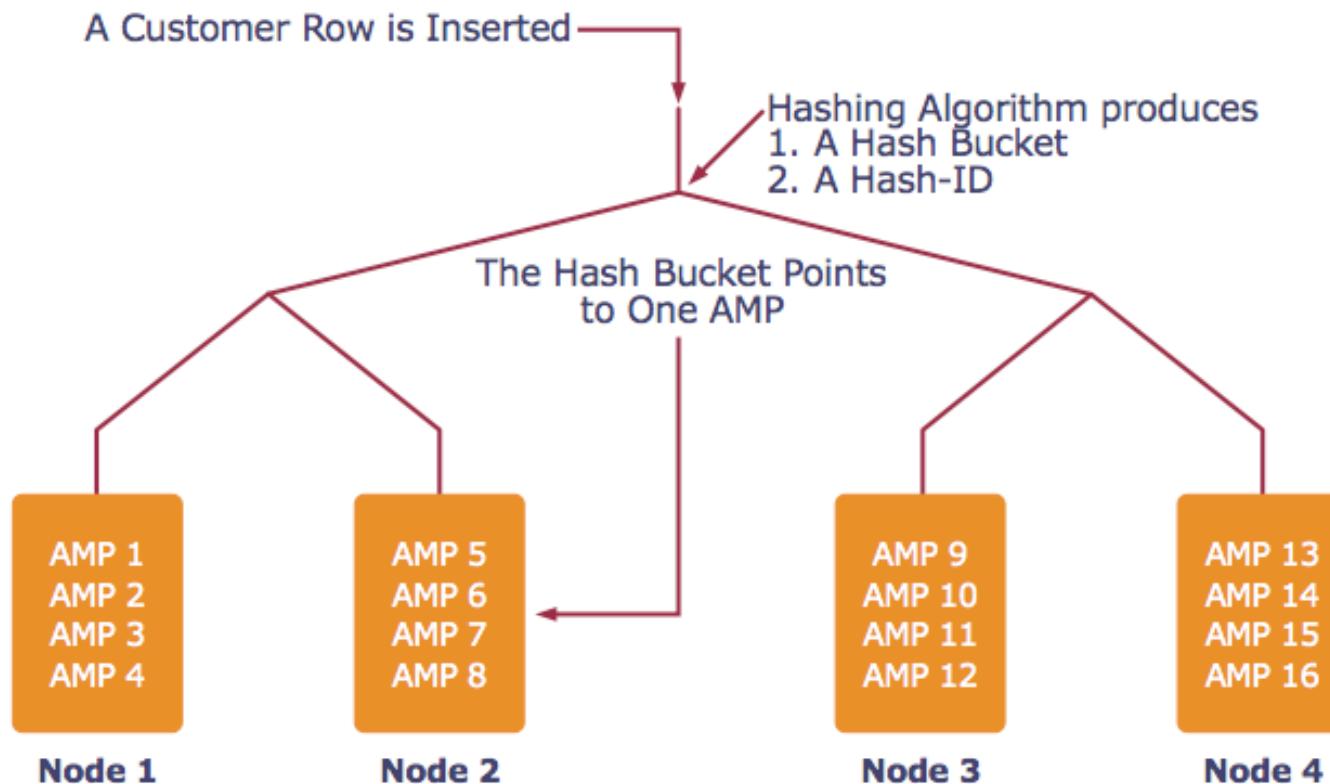




# Shared Nothing Parallel Databases

- Teradata
- Greenplum
- Netezza
- Aster Data Systems
- ~~Datallegro~~ Microsoft
- Vertica
- MonetDB **Recently commercialized as “Vectorwise”**

# Example System: Teradata

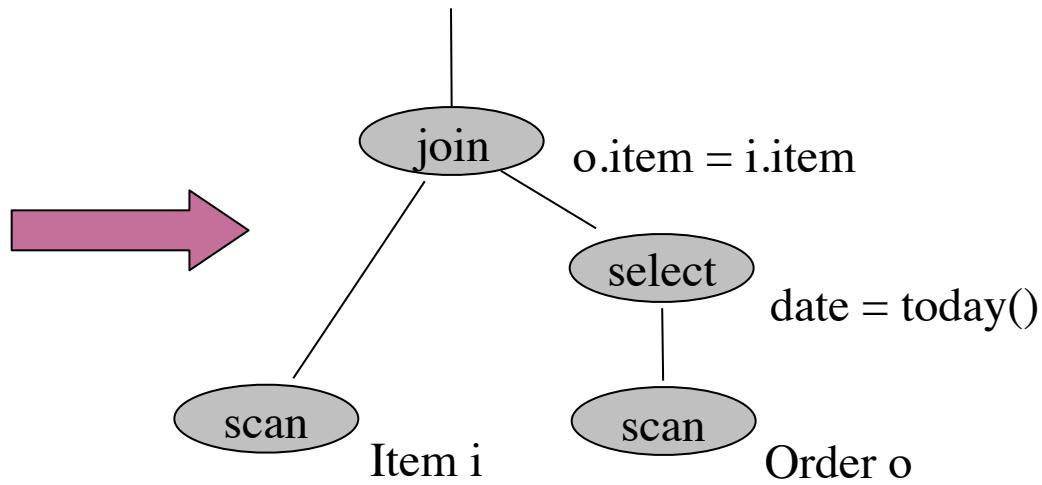


*AMP = unit of parallelism*

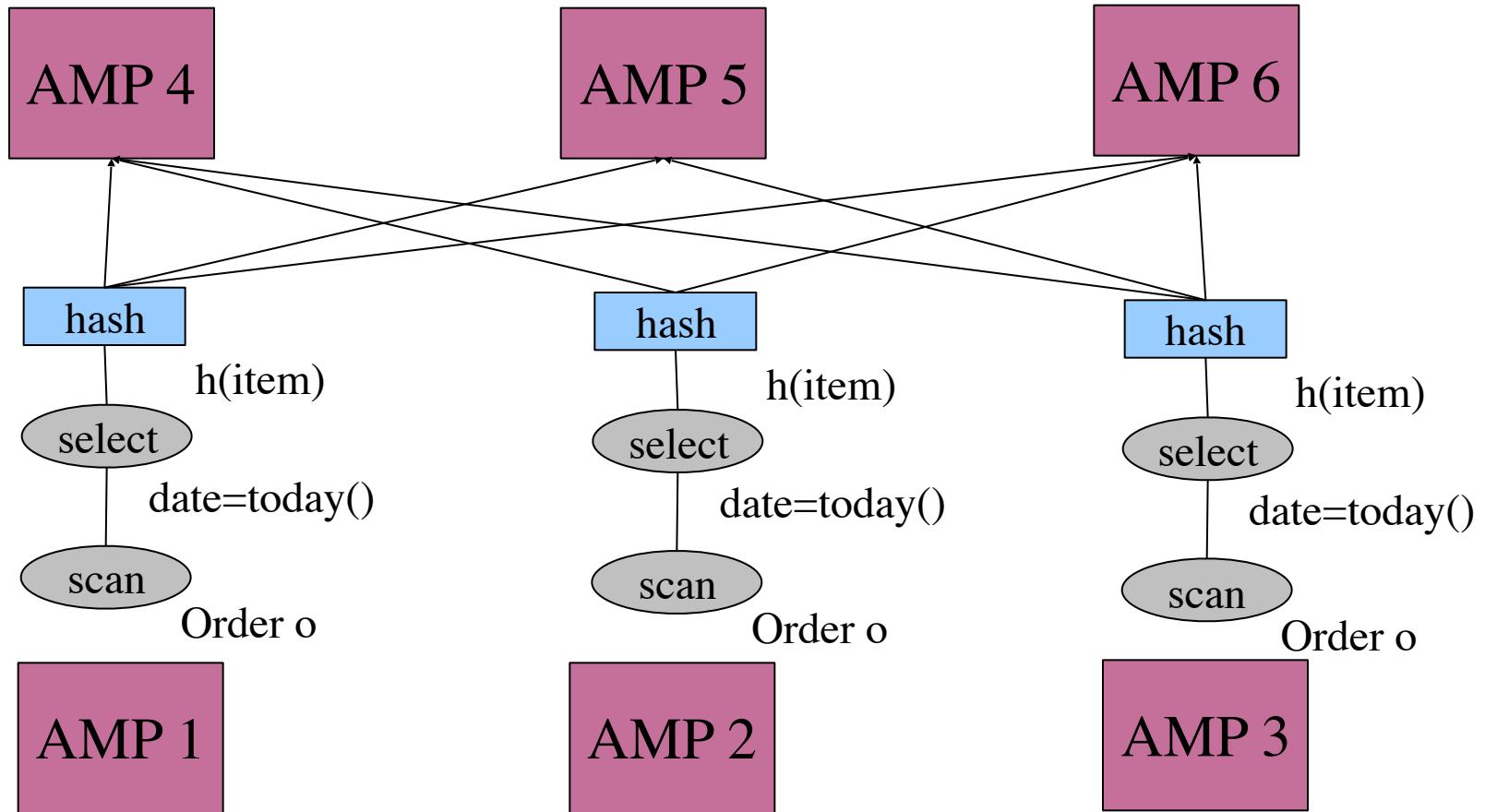
# Example System: Teradata

*Find all orders from today, along with the items ordered*

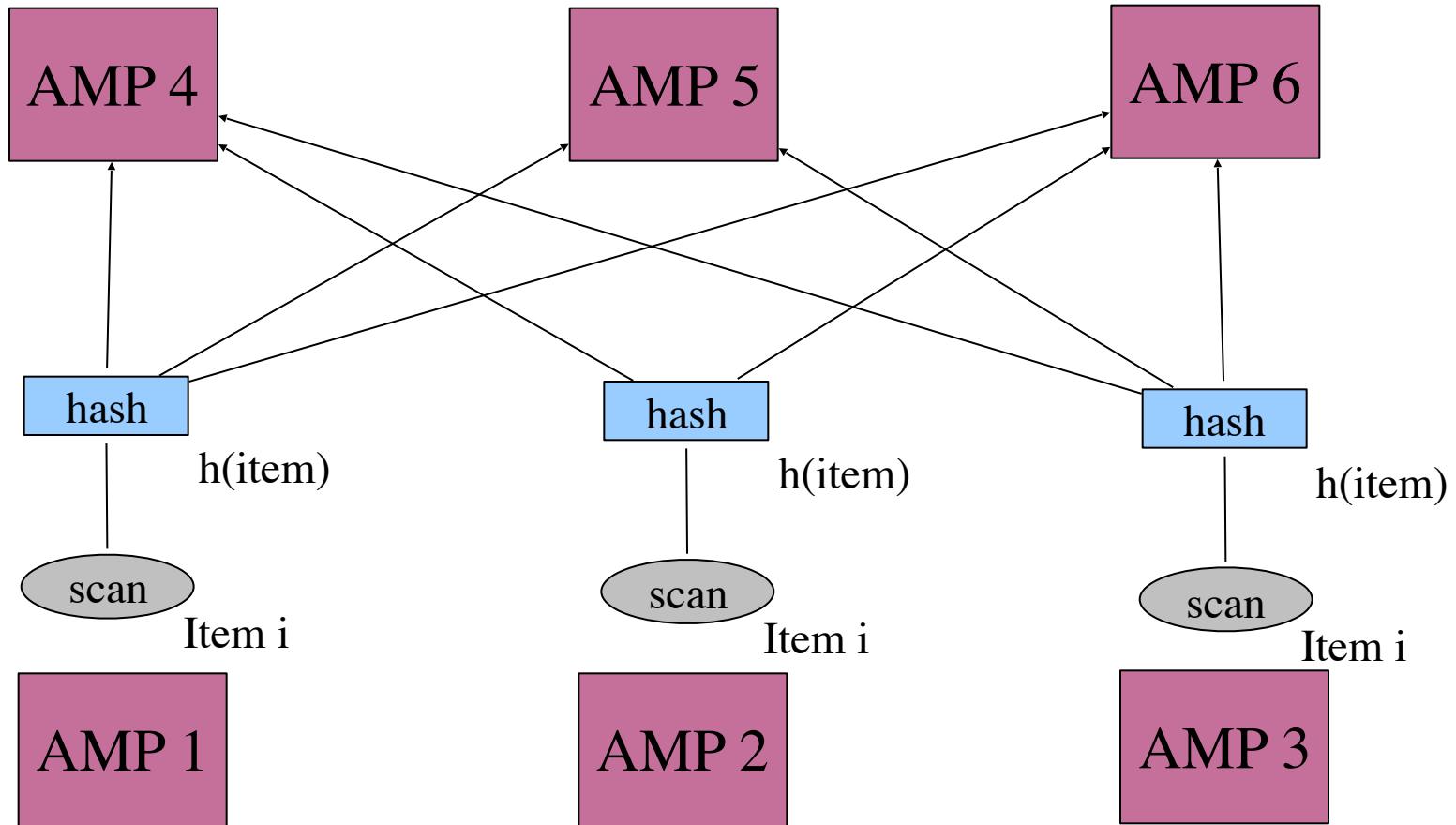
```
SELECT *
  FROM Orders o, Lines i
 WHERE o.item = i.item
   AND o.date = today()
```



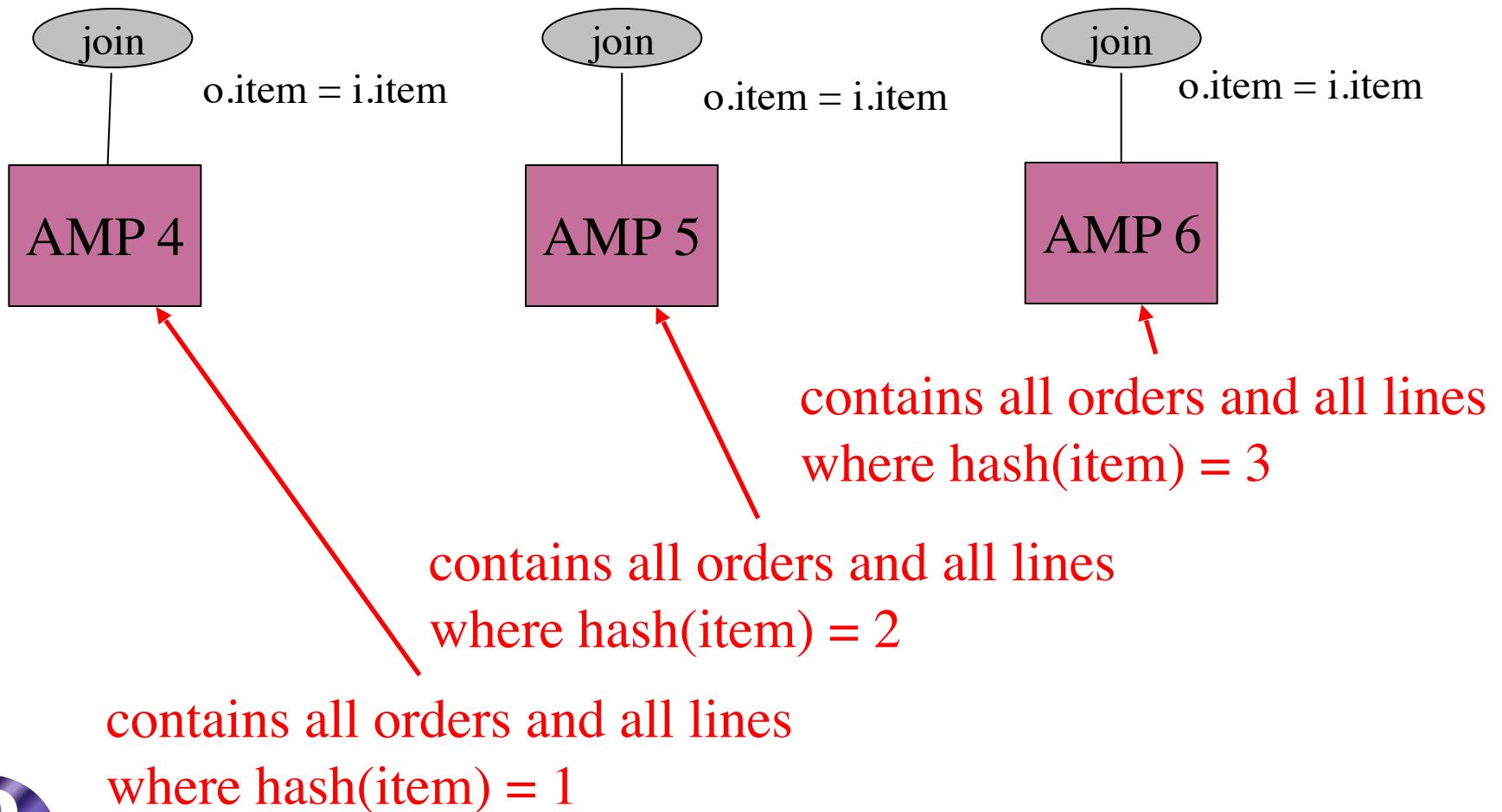
# Example System: Teradata

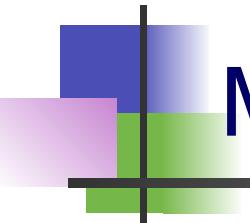


# Example System: Teradata



# Example System: Teradata

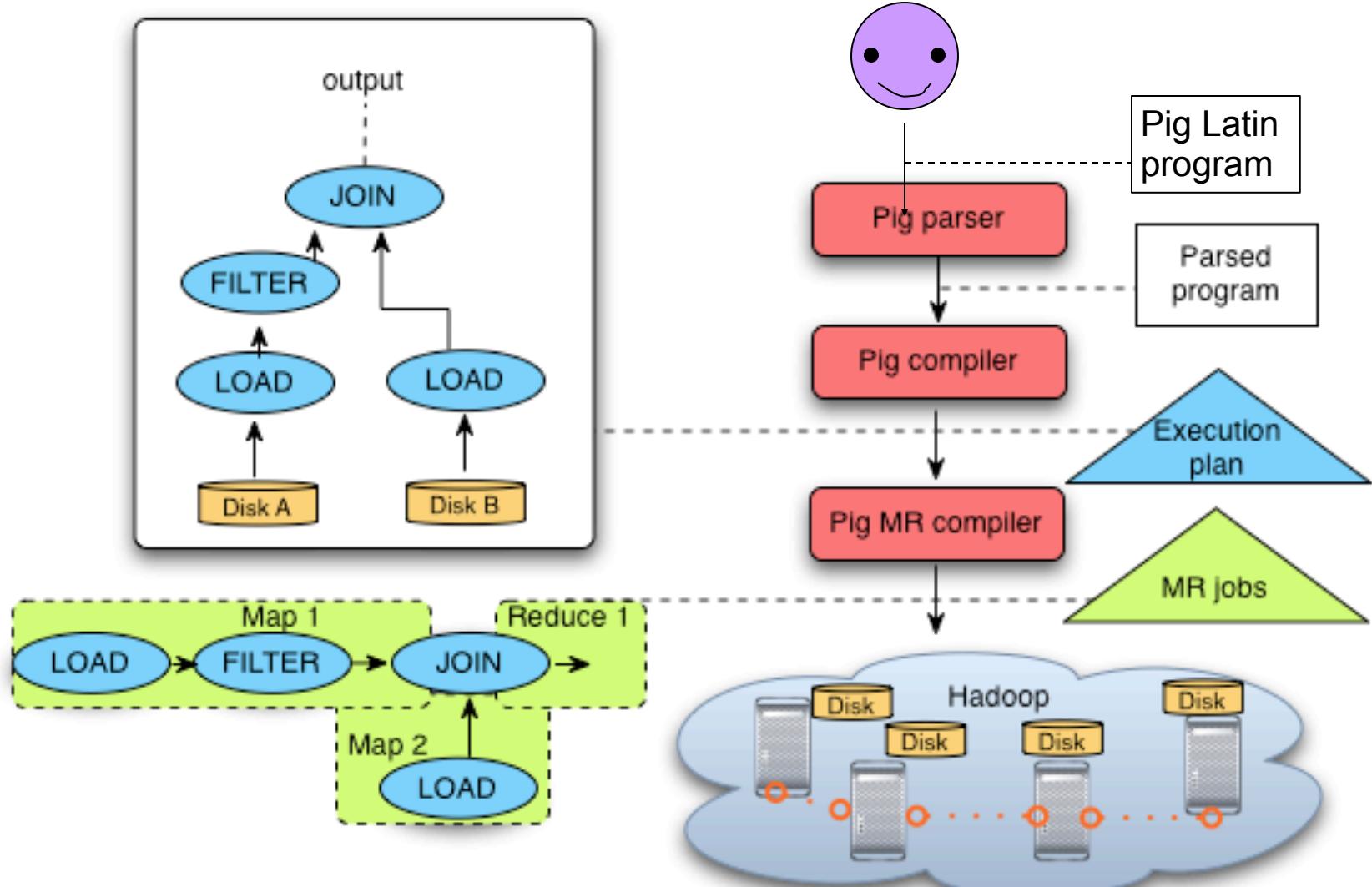


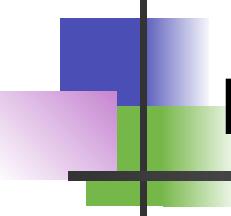


# MapReduce Contemporaries

- Dryad (Microsoft)
  - Relational Algebra
- Pig (Yahoo)
  - Near Relational Algebra over MapReduce
- HIVE (Facebook)
  - SQL over MapReduce
- Cascading
  - Relational Algebra
- Clustera
  - U of Wisconsin
- Hbase
  - Indexing on HDFS

# Example System: Yahoo Pig

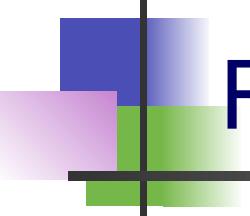




# MapReduce vs RDBMS

- RDBMS
    - Declarative query languages
    - Schemas
    - Logical Data Independence
    - Indexing
    - Algebraic Optimization
    - Caching/Materialized Views
    - *ACID/Tra*n*sactions*
  - MapReduce
    - High Scalability
    - Fault-tolerance
    - “One-person deployment”
- Dryad, Pig, HIVE  
HIVE, Pig, Dryad
- Hbase  
Pig, (Dryad, HIVE)

	<b>Data Model</b>	<b>Prog. Model</b>	<b>Services</b>
<b>GPL</b>	*	*	Typing (maybe)
<b>Workflow</b>	*	dataflow	typing, provenance, scheduling, caching, task parallelism, reuse
<b>Relational Algebra</b>	Relations	Select, Project, Join, Aggregate, ...	optimization, physical data independence, data parallelism
<b>MapReduce</b>	[(key,value)]	Map, Reduce	massive data parallelism, fault tolerance
<b>MS Dryad</b>	IQueryable, IEnumerable	RA + Apply + Partitioning	typing, massive data parallelism, fault tolerance
<b>MPI</b>	Arrays/ Matrices	70+ ops	data parallelism, full control



# Roadmap

- Introduction
- Context: RDBMS, MapReduce, etc.
- New Extensions for Science
  - Recursive MapReduce
  - Skew Handling



# PageRank

Rank Table  $R_0$ 

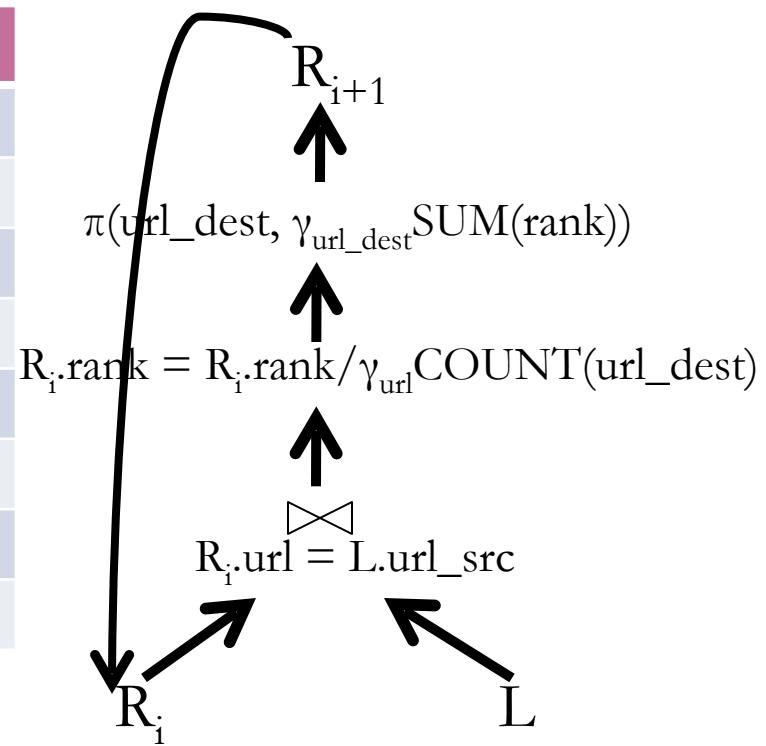
url	rank
<a href="http://www.a.com">www.a.com</a>	1.0
<a href="http://www.b.com">www.b.com</a>	1.0
<a href="http://www.c.com">www.c.com</a>	1.0
<a href="http://www.d.com">www.d.com</a>	1.0
<a href="http://www.e.com">www.e.com</a>	1.0

Rank Table  $R_3$ 

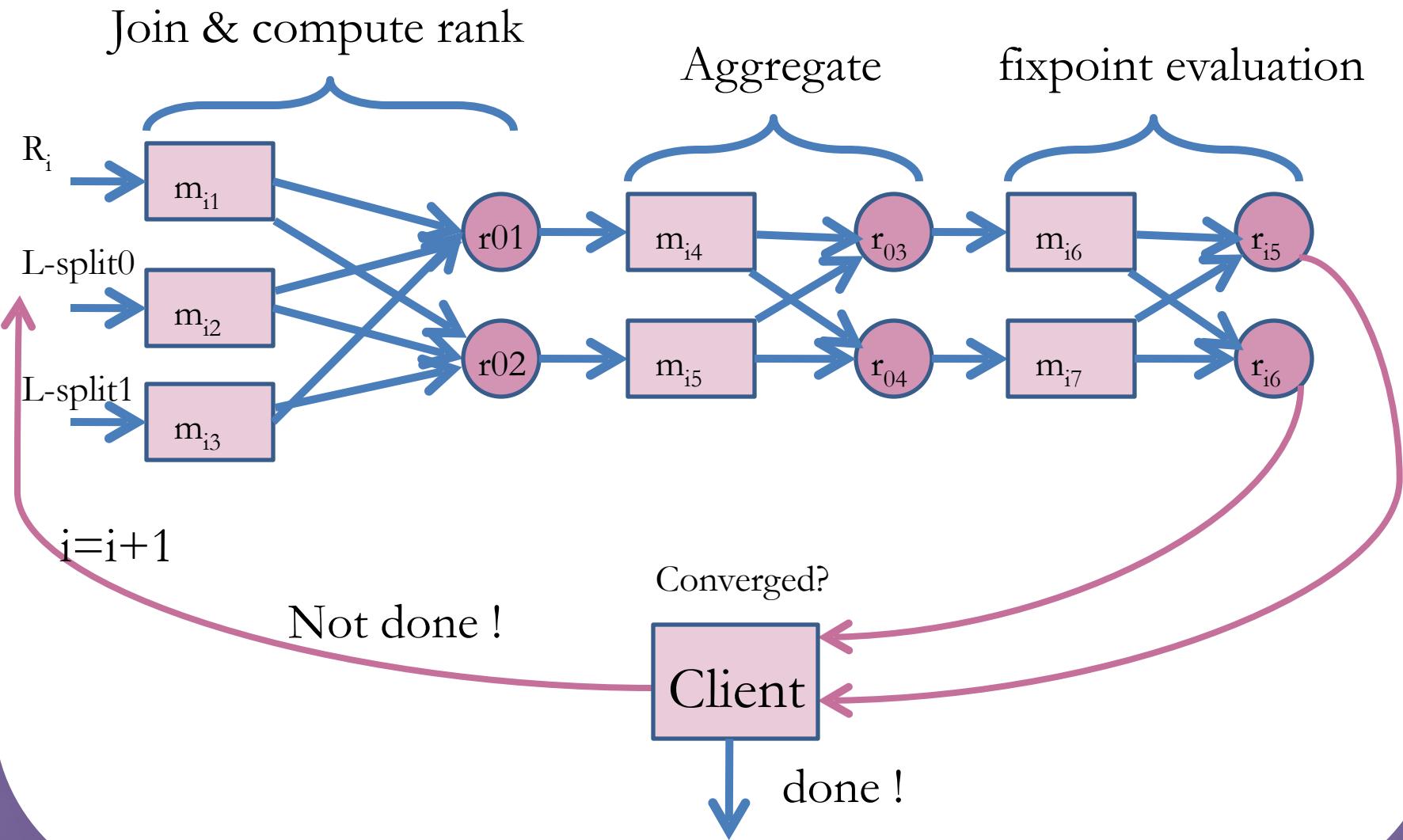
url	rank
<a href="http://www.a.com">www.a.com</a>	2.13
<a href="http://www.b.com">www.b.com</a>	3.89
<a href="http://www.c.com">www.c.com</a>	2.60
<a href="http://www.d.com">www.d.com</a>	2.60
<a href="http://www.e.com">www.e.com</a>	2.13

Linkage Table  $L$ 

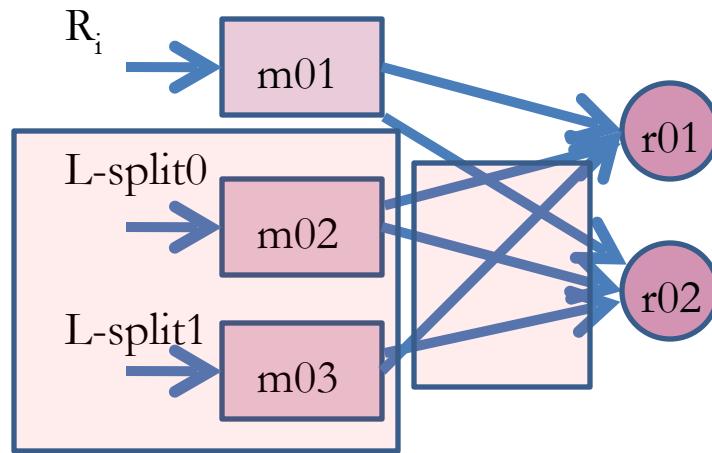
url_src	url_dest
<a href="http://www.a.com">www.a.com</a>	<a href="http://www.b.com">www.b.com</a>
<a href="http://www.a.com">www.a.com</a>	<a href="http://www.c.com">www.c.com</a>
<a href="http://www.c.com">www.c.com</a>	<a href="http://www.a.com">www.a.com</a>
<a href="http://www.e.com">www.e.com</a>	<a href="http://www.c.com">www.c.com</a>
<a href="http://www.d.com">www.d.com</a>	<a href="http://www.b.com">www.b.com</a>
<a href="http://www.c.com">www.c.com</a>	<a href="http://www.e.com">www.e.com</a>
<a href="http://www.e.com">www.e.com</a>	<a href="http://www.c.com">www.c.com</a>
<a href="http://www.a.com">www.a.com</a>	<a href="http://www.d.com">www.d.com</a>



# MapReduce Implementation



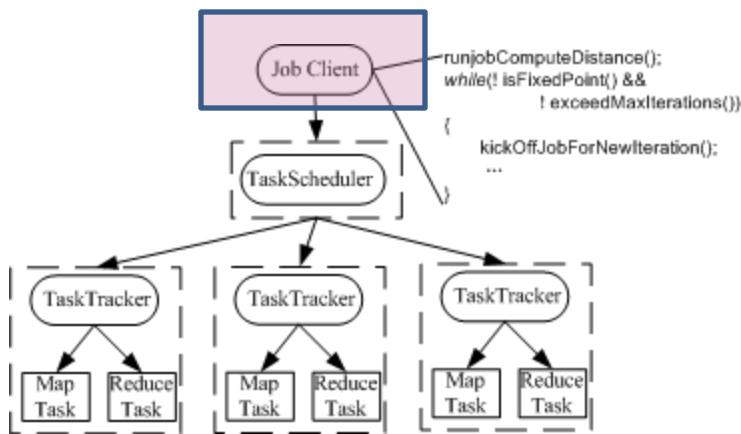
# What's the problem?



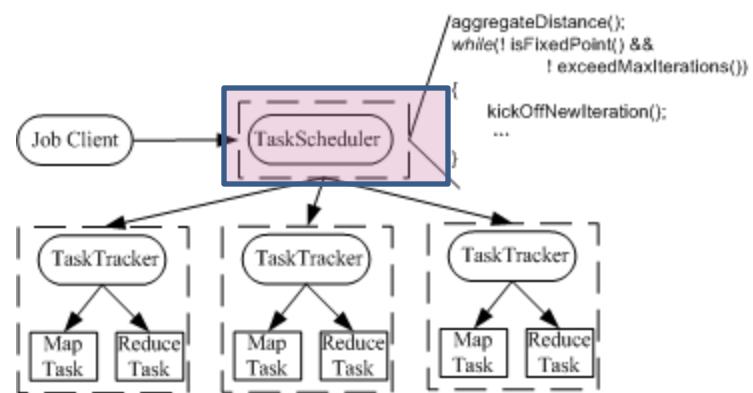
- $L$  is loaded and shuffled in each iteration
- $L$  never changes

# HaLoop: Loop-aware Hadoop

Hadoop



HaLoop

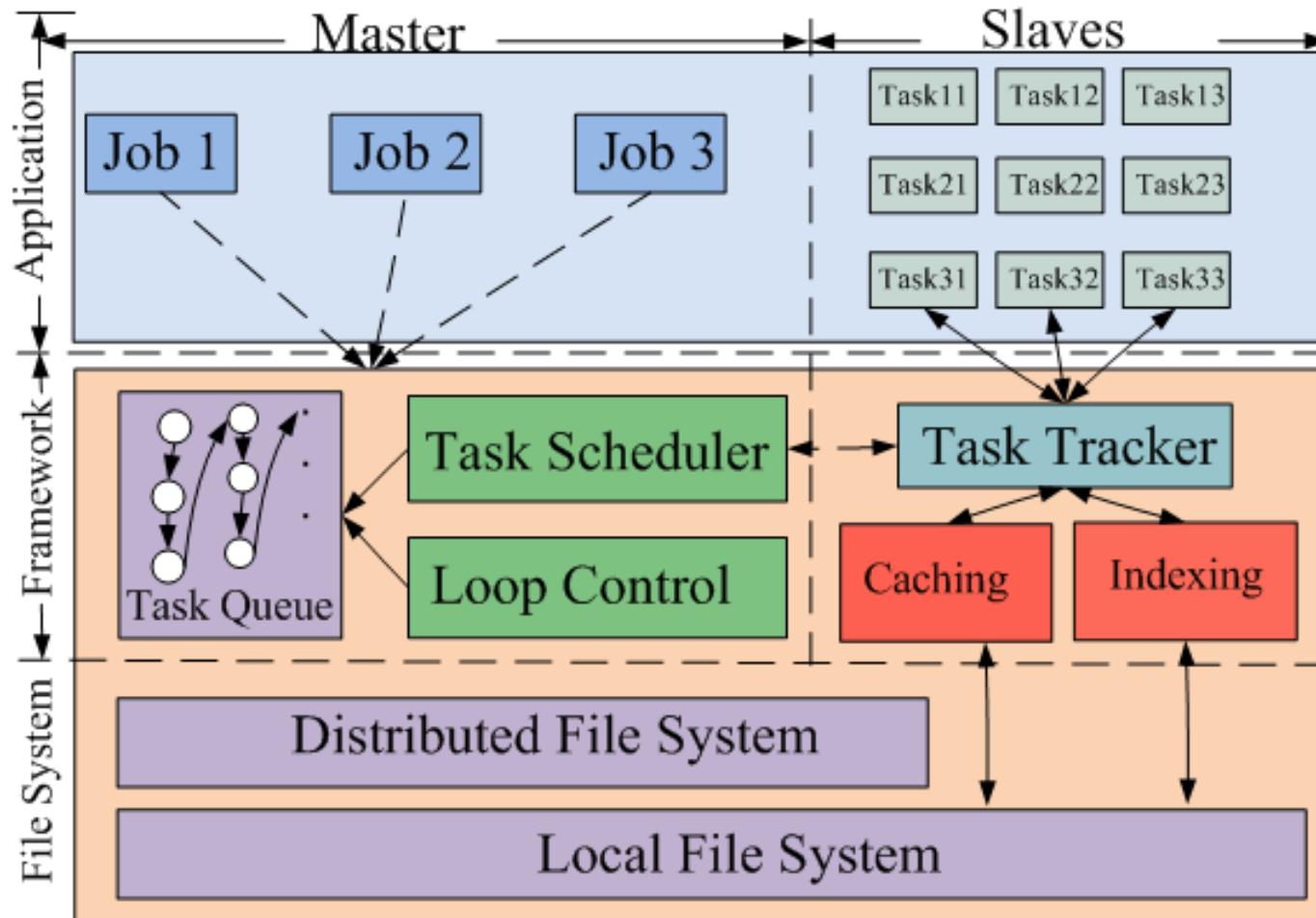


- Hadoop: Loop control in client program
- HaLoop: Loop control in master node

# Feature: Inter-iteration Locality

- Mapper Output Cache
  - K-means
  - Neural network analysis
- Reducer Input Cache
  - Recursive join
  - PageRank
  - HITs
  - Social network analysis
- Reducer Output Cache
  - Fixpiont evaluation

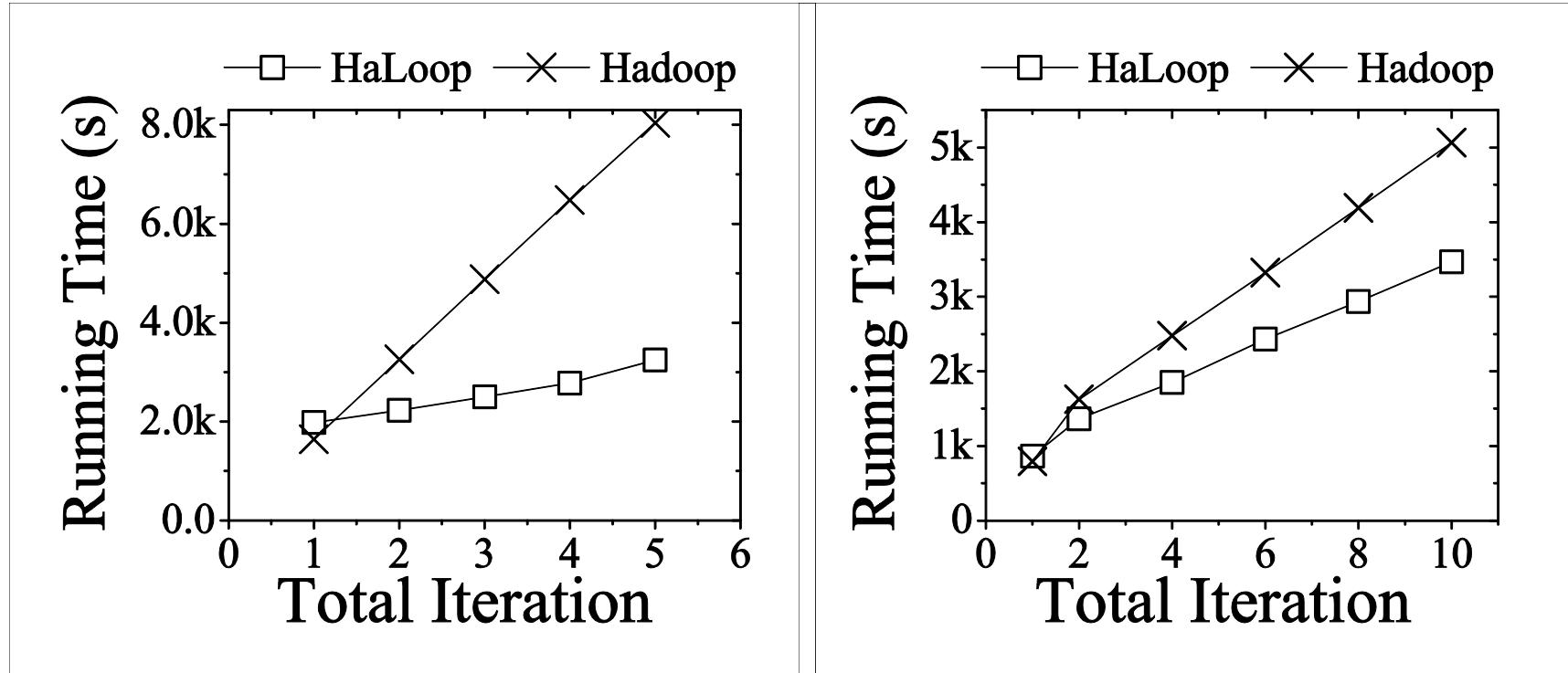
# HaLoop Architecture



# Experiments

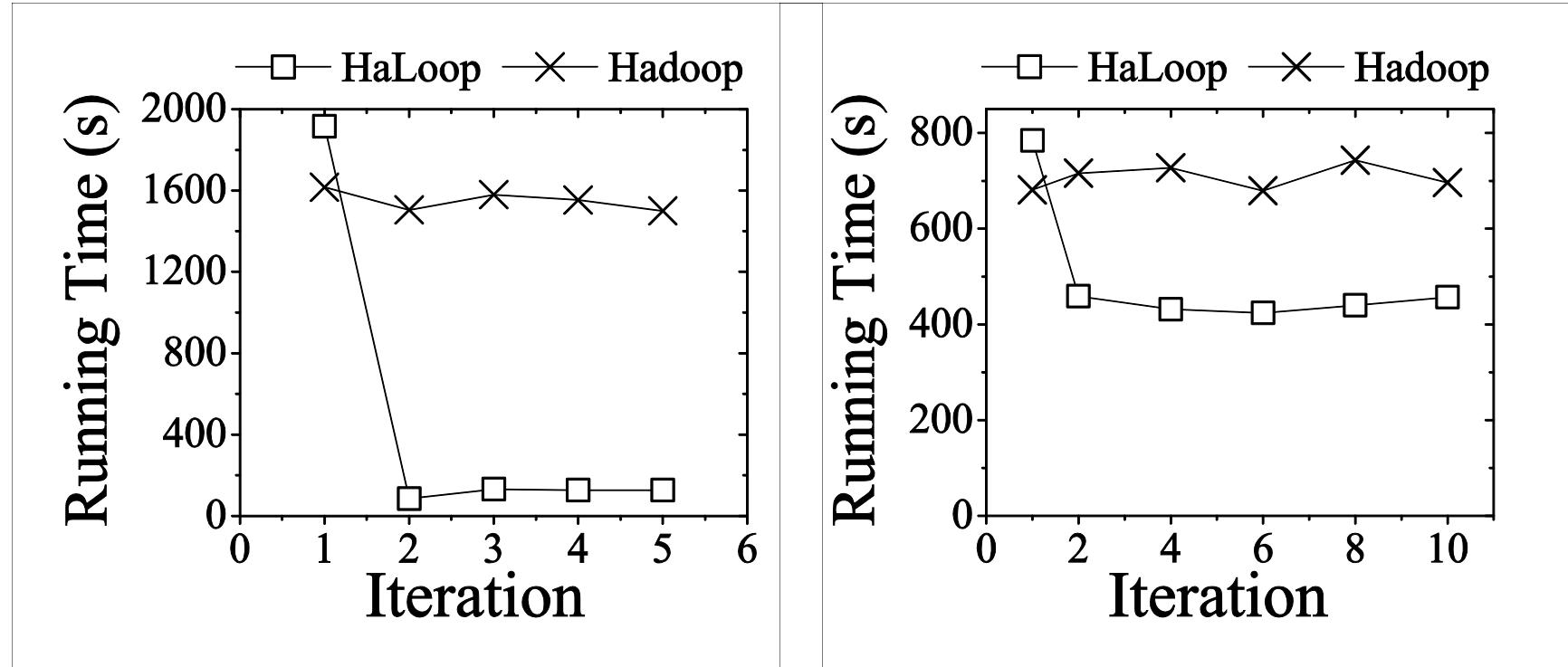
- Amazon EC2
  - 20, 50, 90 default small instances
- Datasets
  - Billions of Triples (120GB)
  - Freebase (12GB)
  - Livejournal social network (18GB)
- Queries
  - Transitive Closure
  - PageRank
  - k-means

# Application Run Time



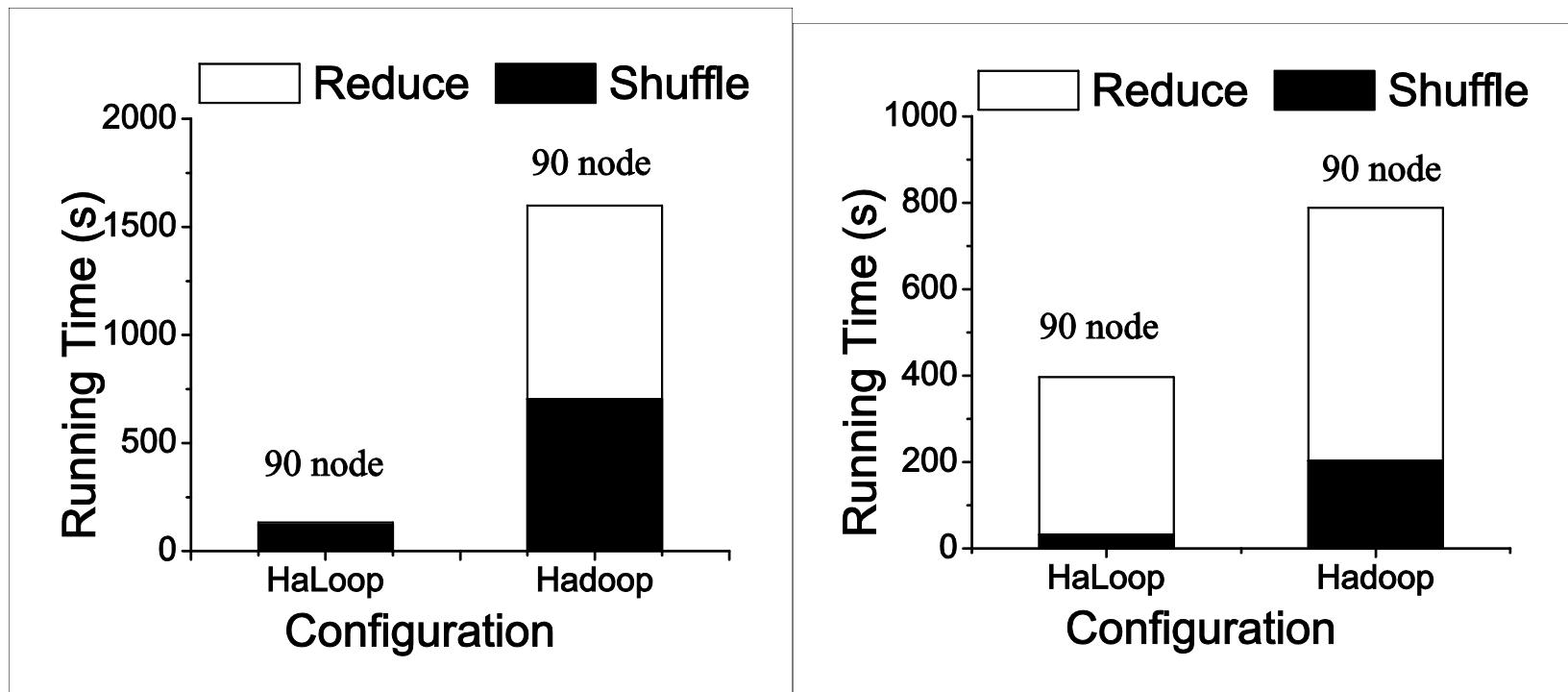
- Transitive Closure
- PageRank

# Join Time



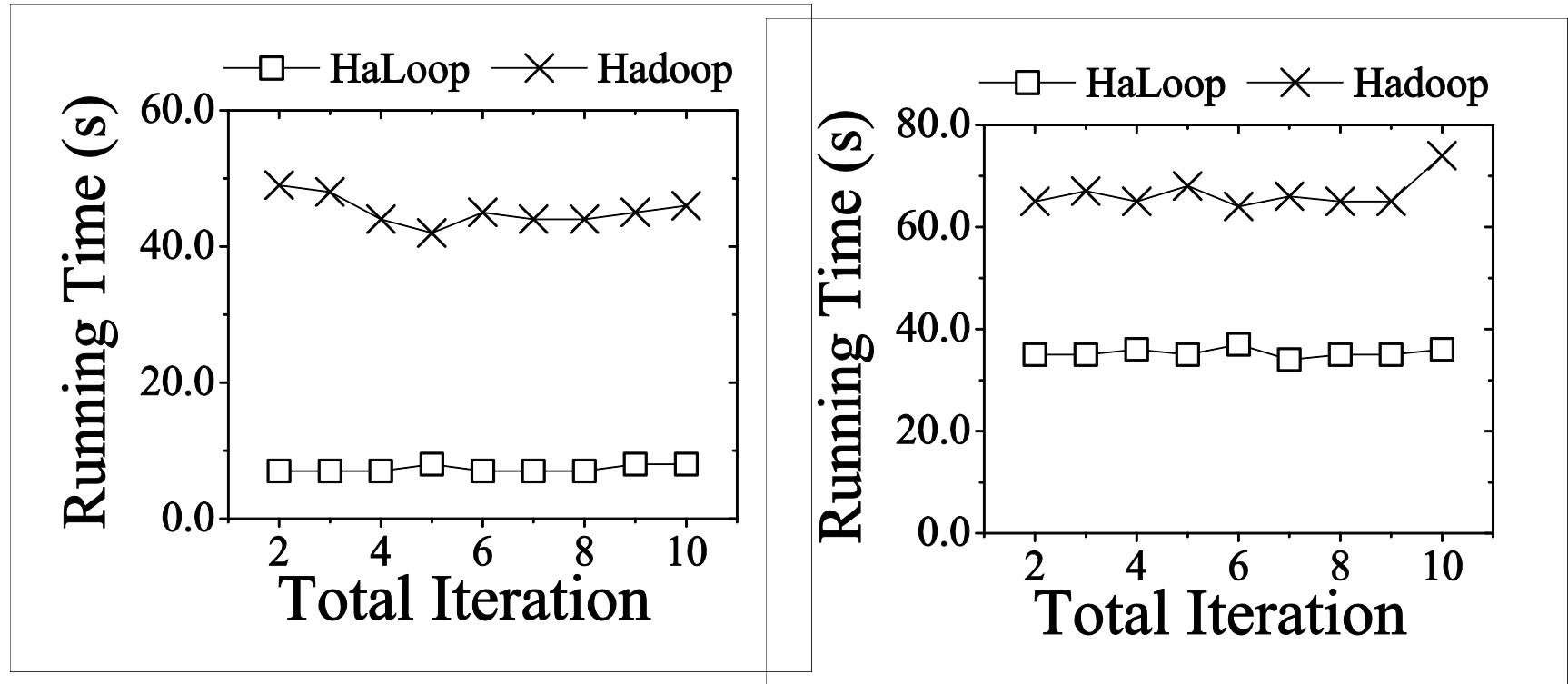
- Transitive Closure
- PageRank

# Run Time Distribution

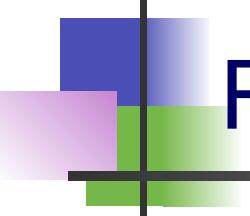


- Transitive Closure
- PageRank

# Fixpoint Evaluation



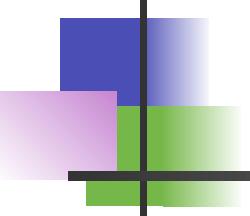
- PageRank



# Roadmap

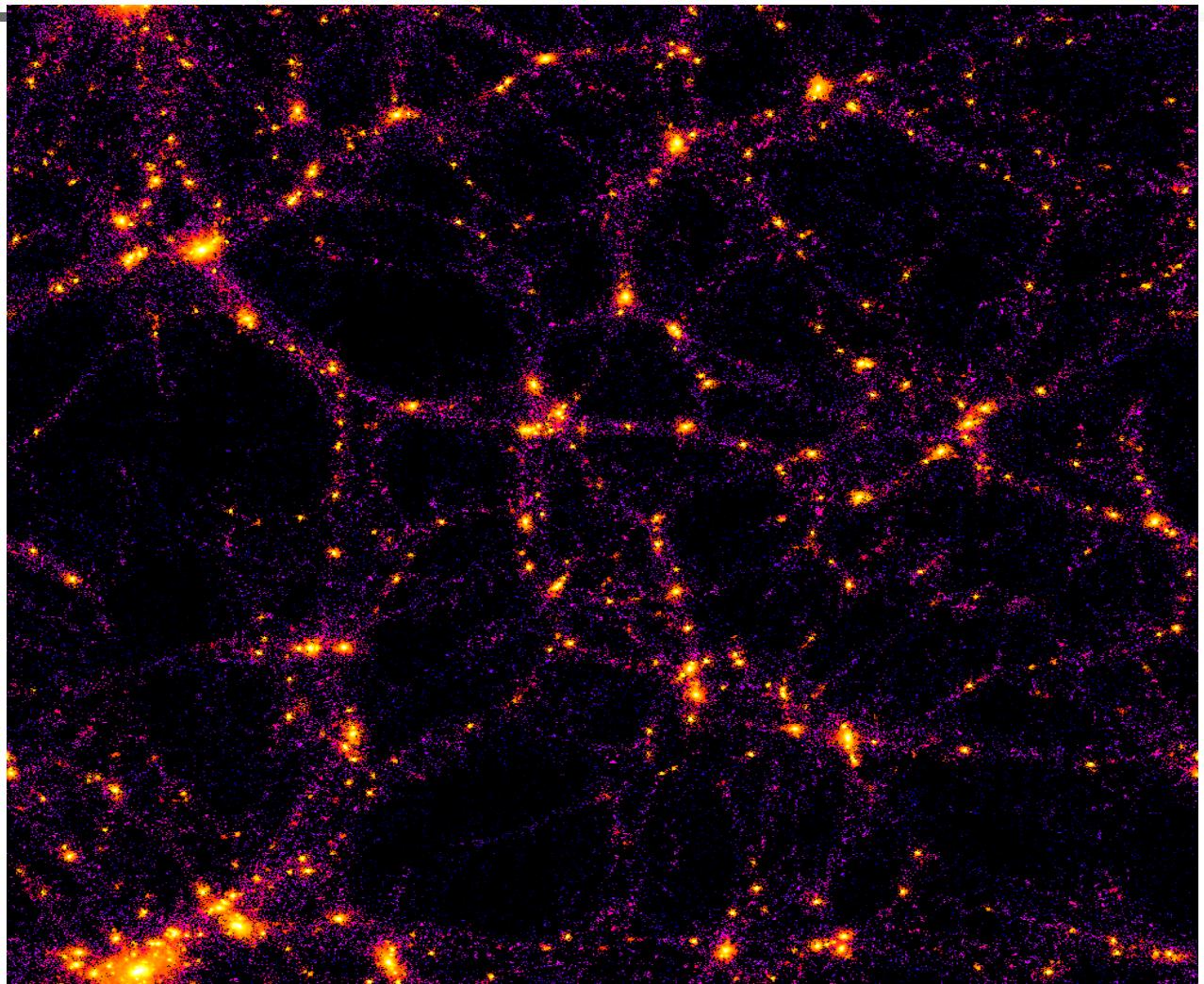
- Introduction
- Context: RDBMS, MapReduce, etc.
- New Extensions for Science
  - Recursive MapReduce
  - Skew Handling



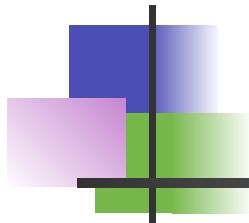


# N-body Astrophysics Simulation

- 15 years in dev
- $10^9$  particles
- Months to run
- 7.5 million CPU hours
- 500 timesteps
- Big Bang to now



Simulations from Tom Quinn's Lab, work by Sarah Loebman, YongChul Kwon, Bill Howe, Jeff Gardner, Magda Balazinska  
Bill Howe, UW

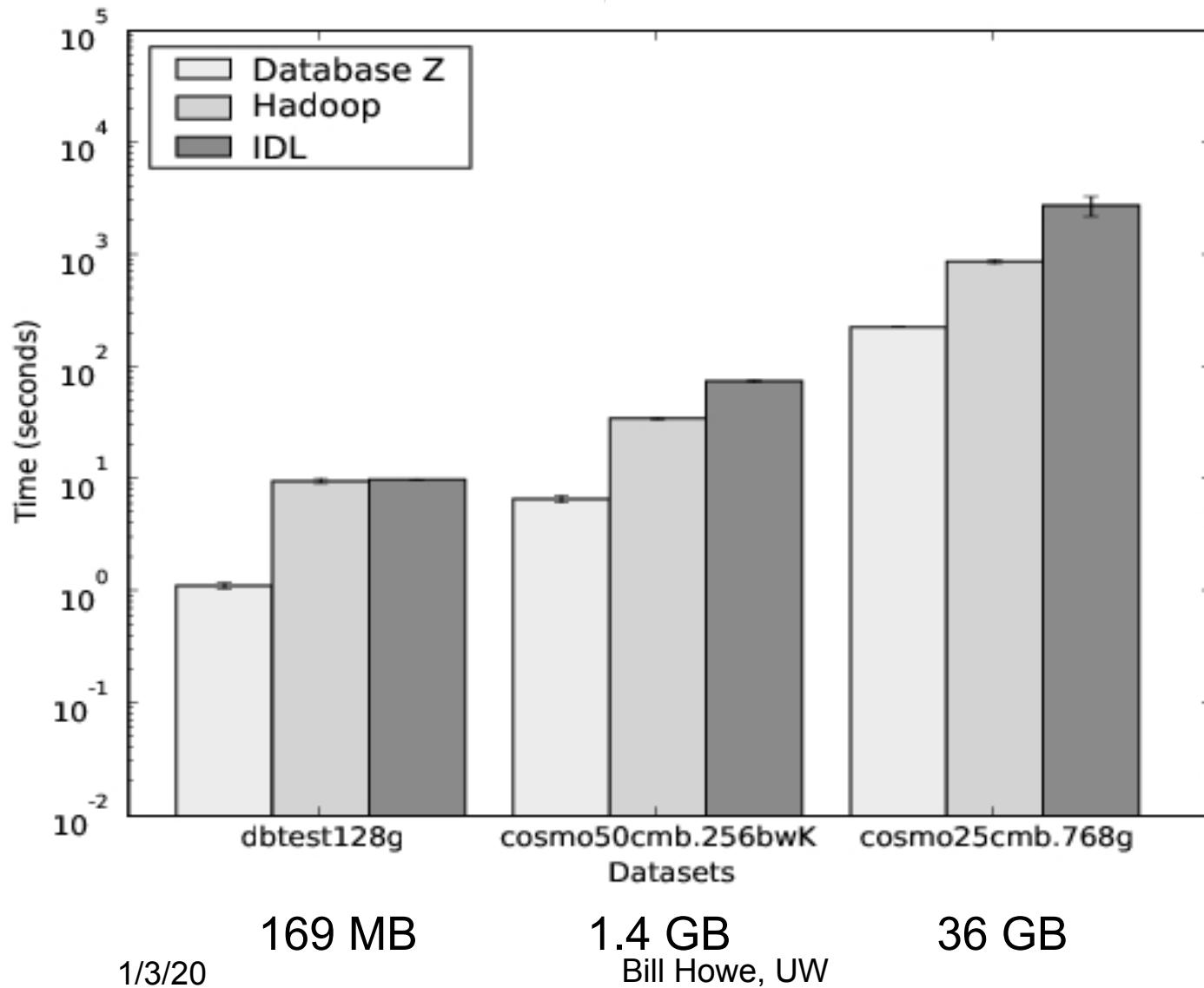


# Q1: Find Hot Gas

```
SELECT id  
      FROM gas  
     WHERE temp > 150000
```

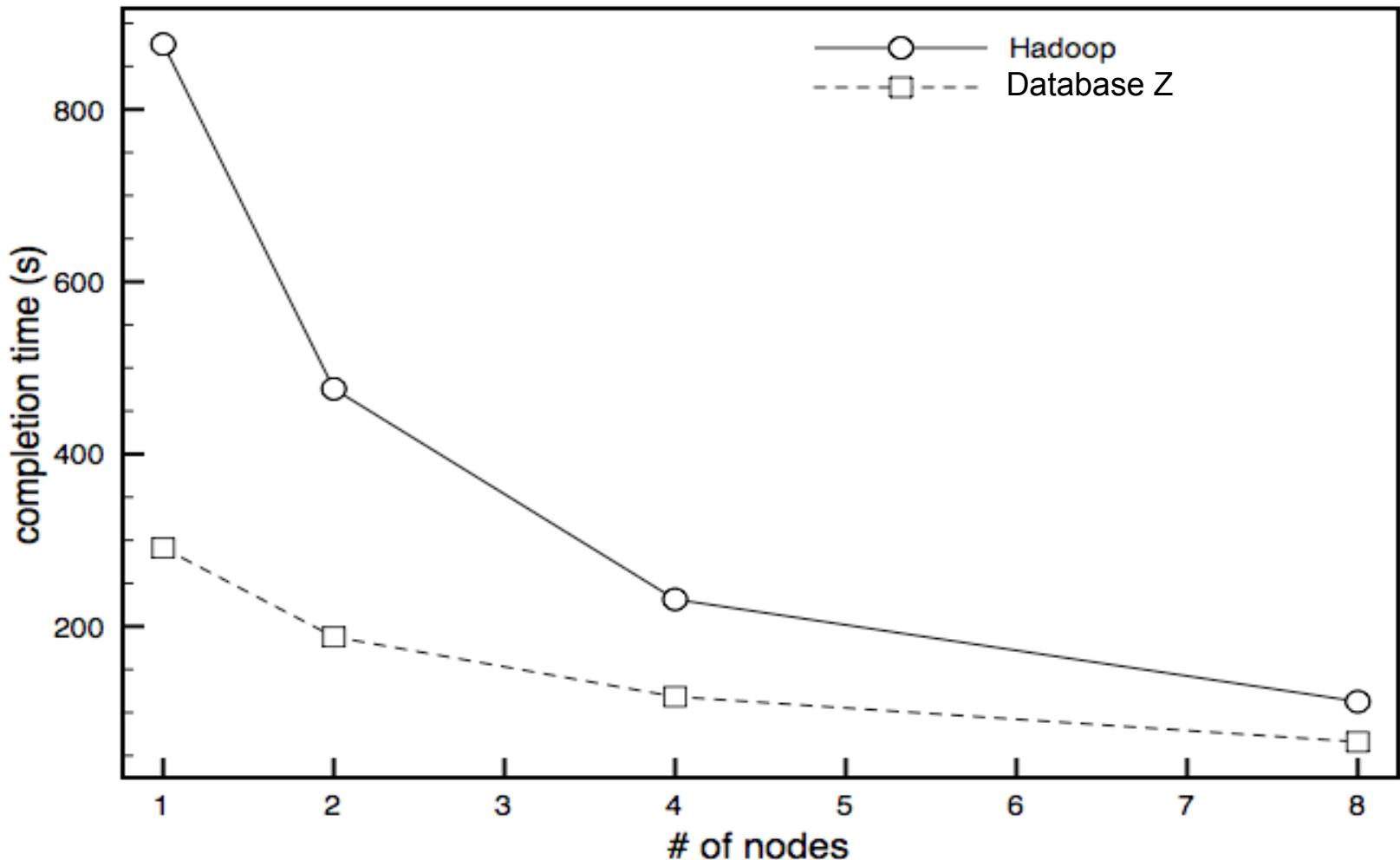
# Single Node: Query 1

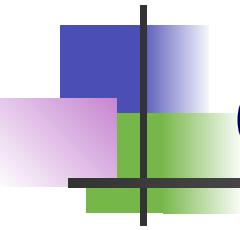
[IASDS 09]



# Multiple Nodes: Query 1

[IASDS 09]





# Q4: Gas Deletion

[IASDS 09]

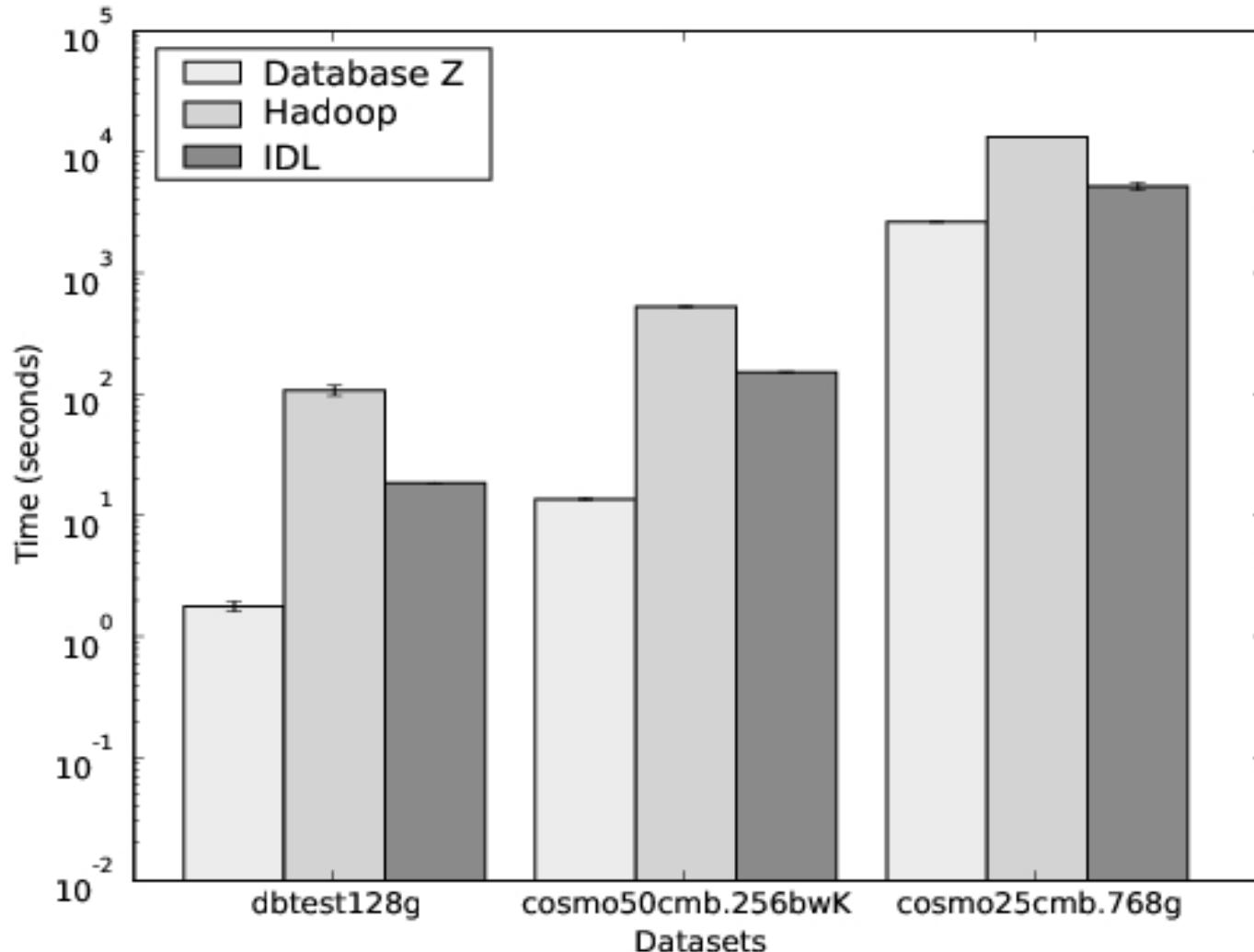
```
SELECT gas1.id  
FROM gas1  
FULL OUTER JOIN gas2  
ON gas1.id=gas2.id  
WHERE gas2.id=NULL
```

Particles removed  
between two timesteps



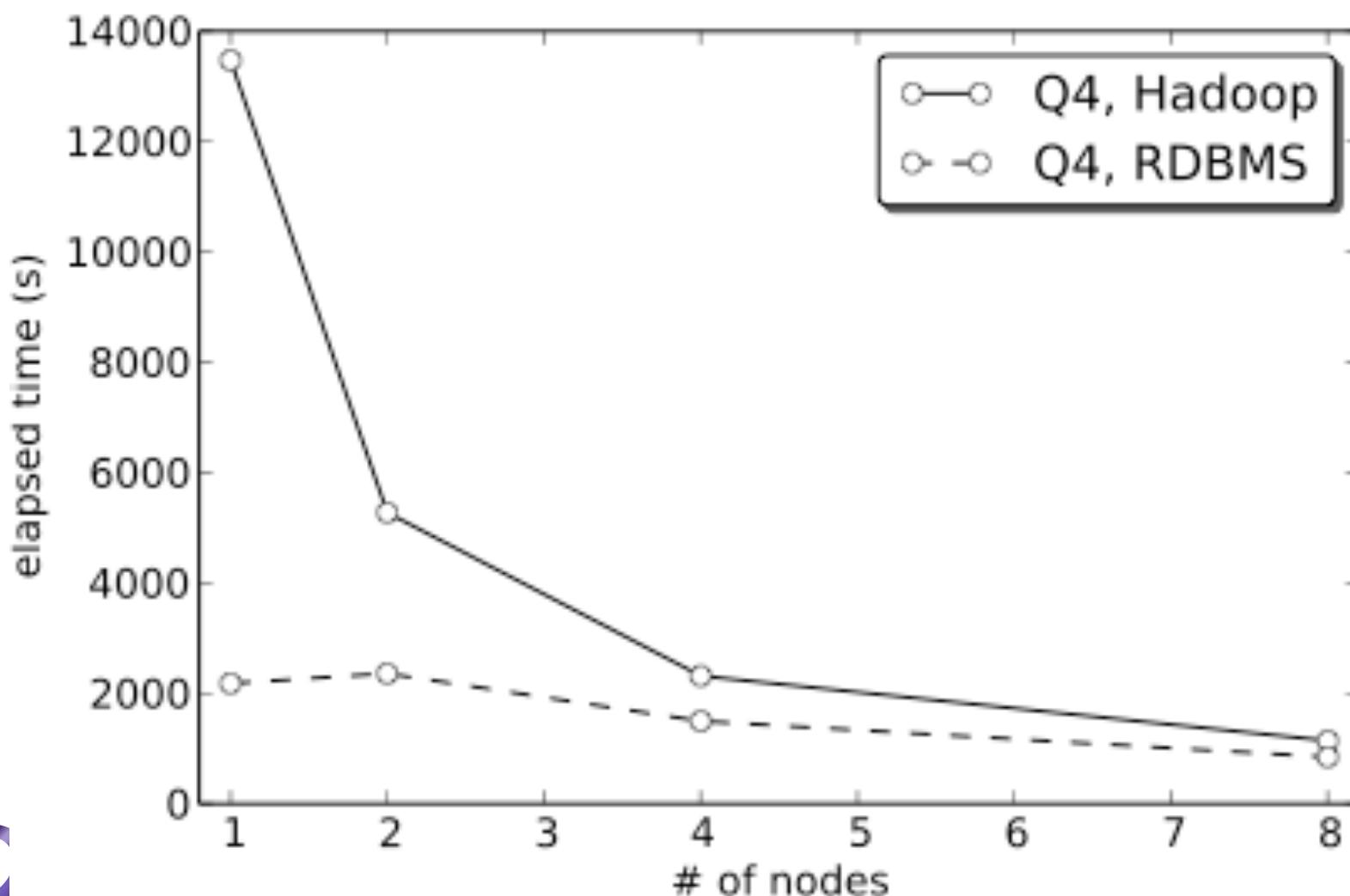
# Single Node: Query 4

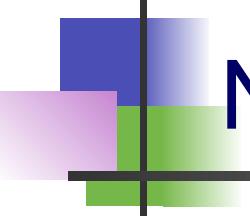
[IASDS 09]



# Multiple Nodes: Query 4

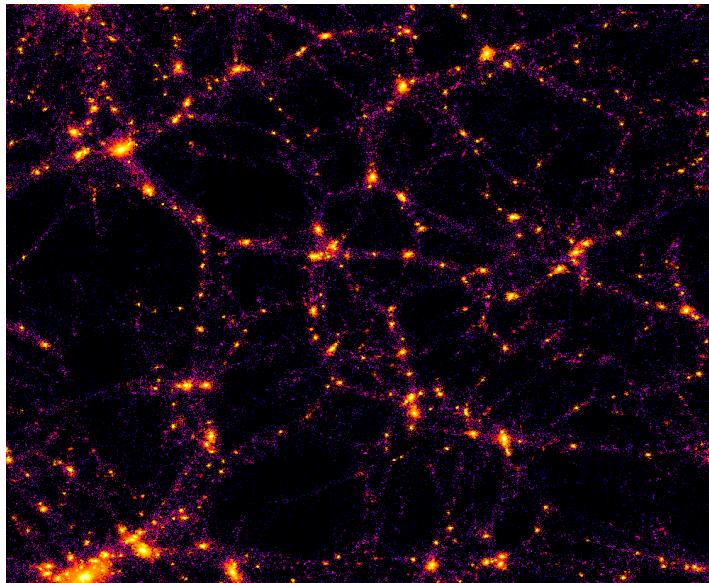
[IASDS 09]

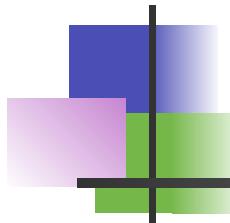




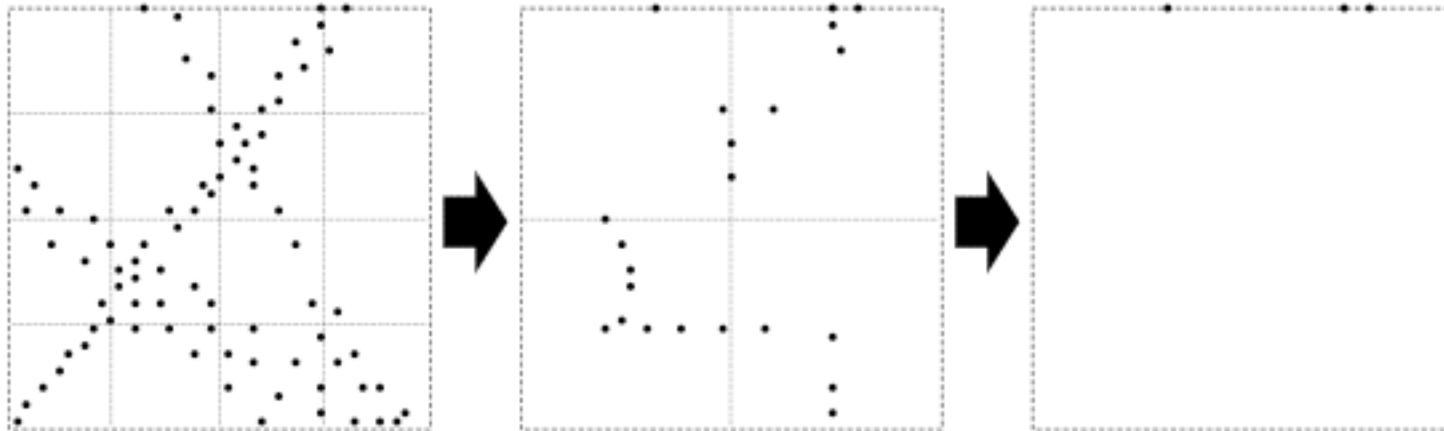
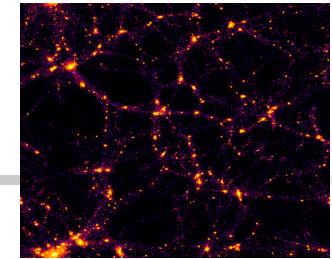
# New Task: Scalable Clustering

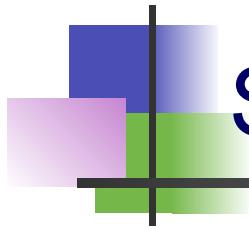
- Group particles into spatial clusters



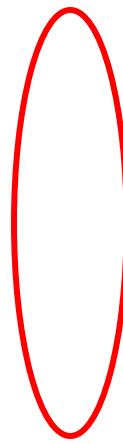


# Scalable Clustering

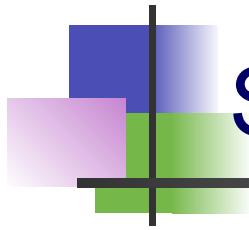




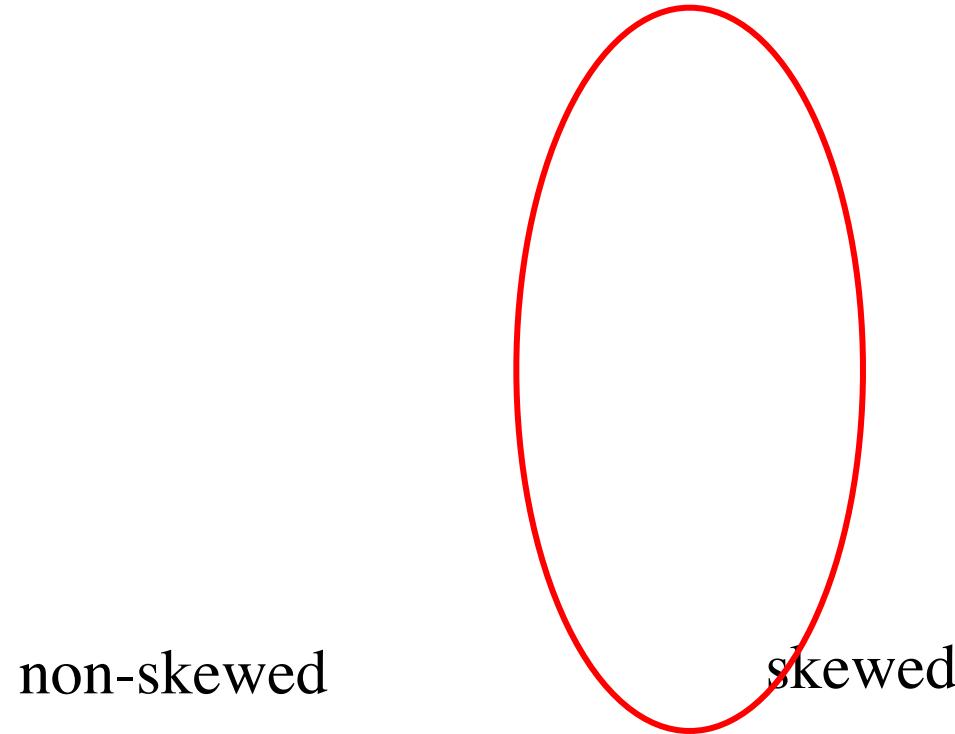
# Scalable Clustering in Dryad



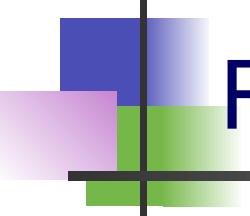
[Kwon SSDBM 2010]



# Scalable Clustering in Dryad



YongChul Kwon, Dylan Nunlee, Jeff Gardner, Sarah Loebman, Magda Balazinska, Bill Howe



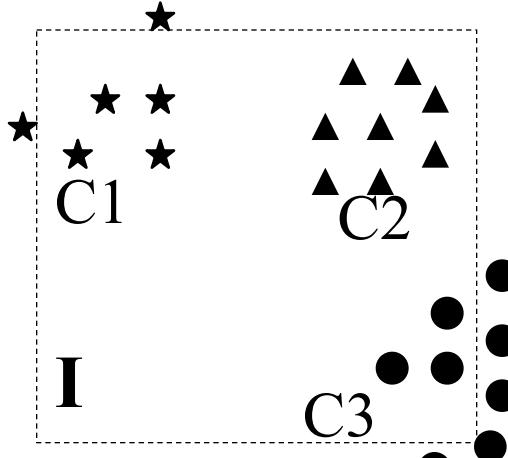
# Roadmap

- Introduction
- Context: RDBMS, MapReduce, etc.
- New Extensions for Science
  - Recursive MapReduce
  - Skew Handling

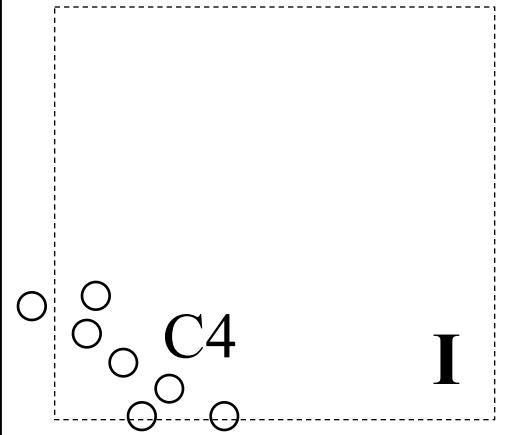


# Example: Friends of Friends

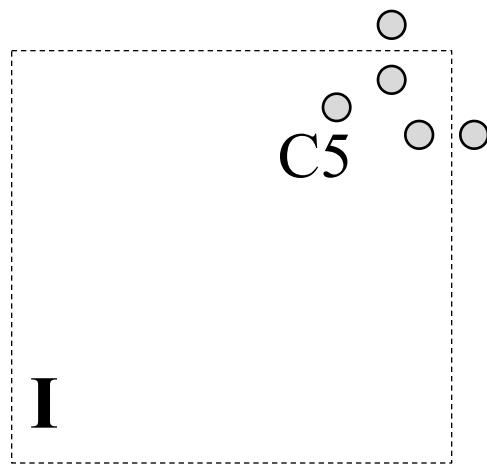
P1



P2

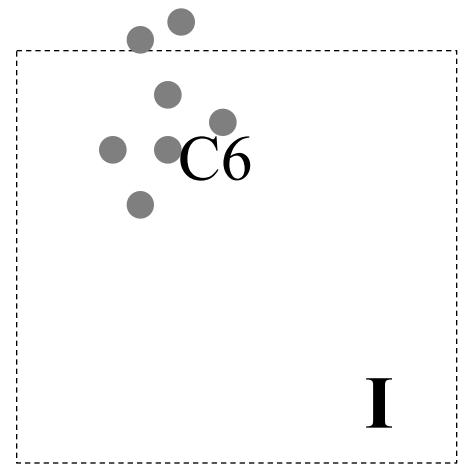


P3

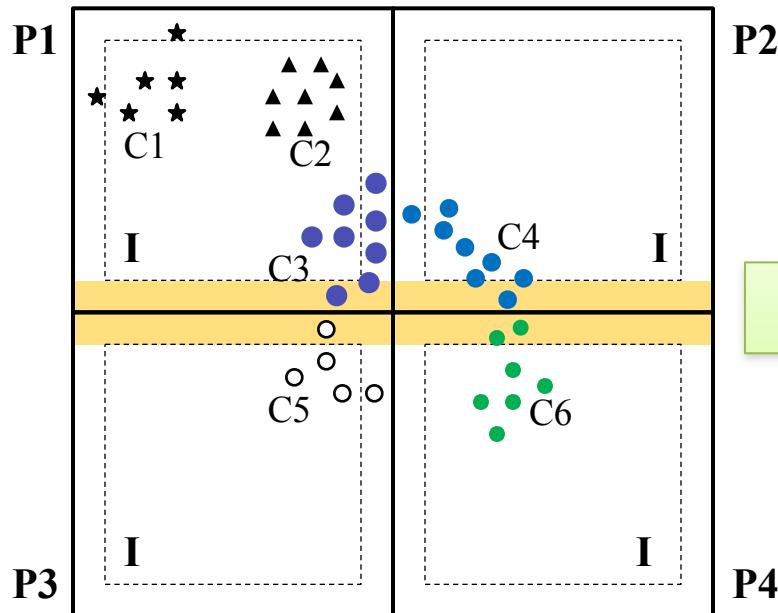


I

P4

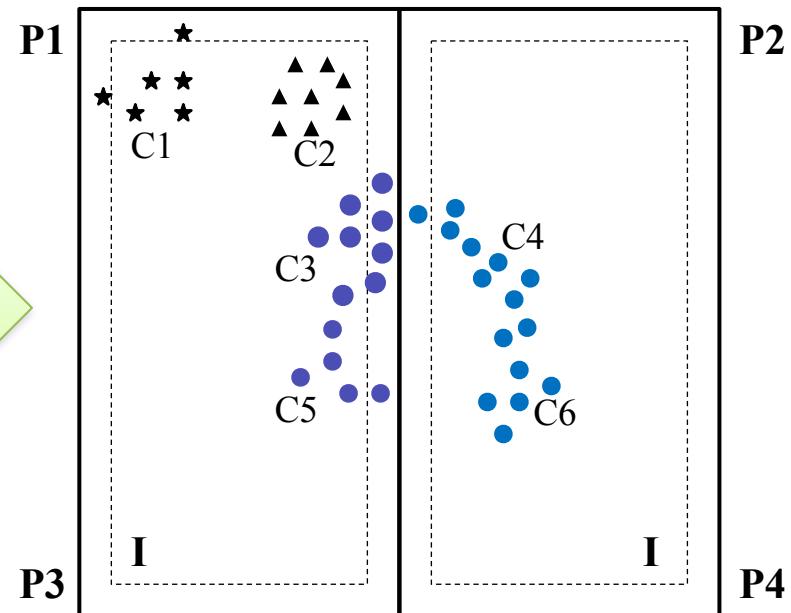


# Example: Friends of Friends



P2

P4



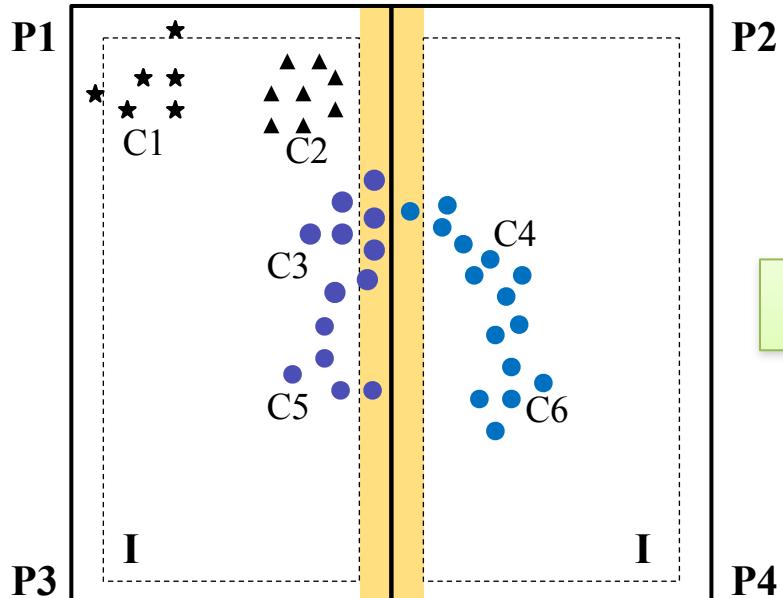
P2

P4

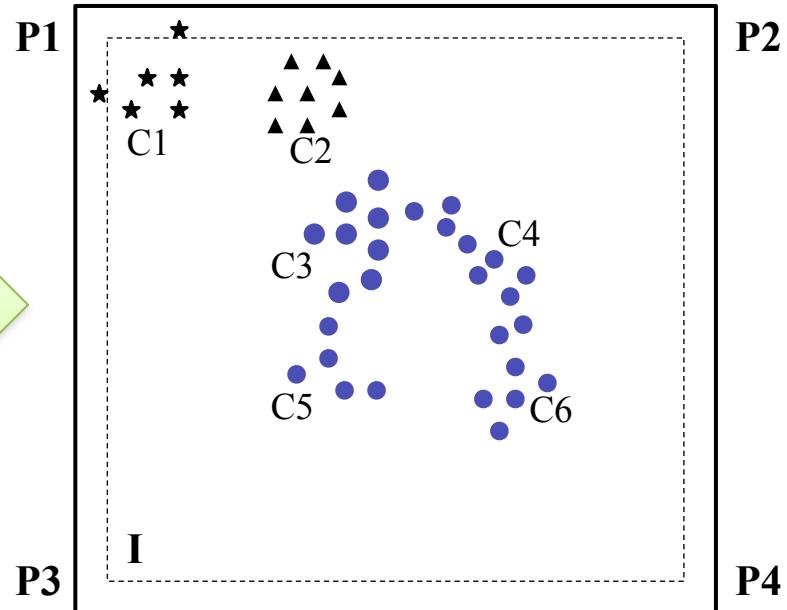
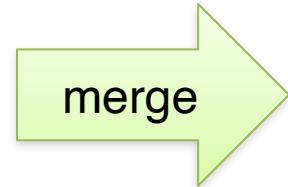
Merge P1, P3  
Merge P2, P4

C5 → C3  
C6 → C4

# Example: Friends of Friends



P2



P3

P2

P4

Merge P1-P3, P2-P4

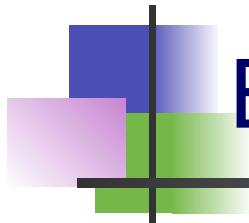
C5 → C3

C6 → C4

C4 → C3

C5 → C3

C6 → C3



# Example: Unbalanced Computation

Local FoF



5 minutes



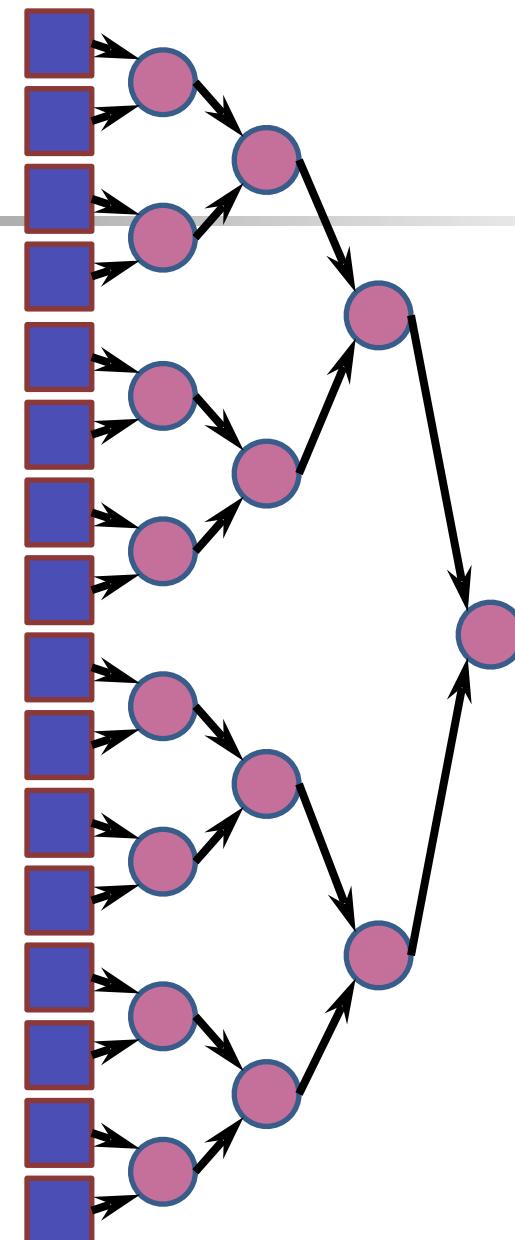
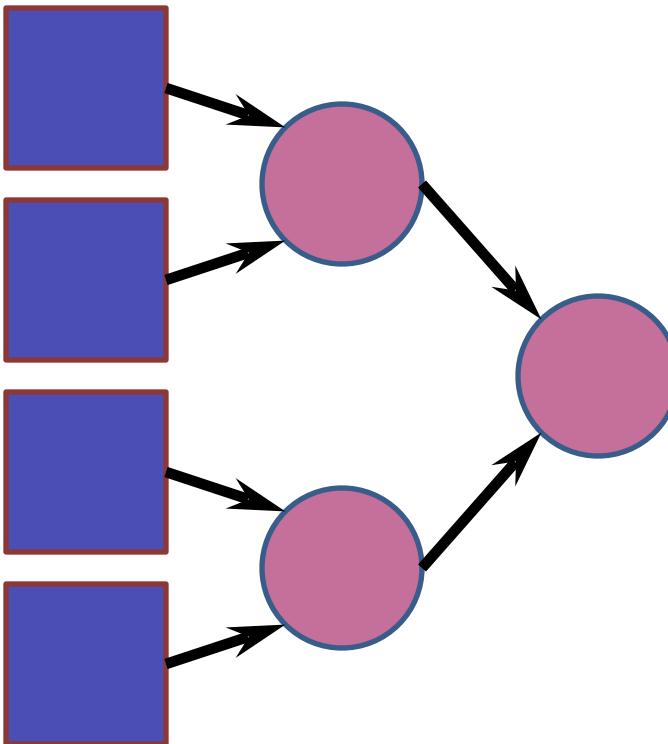
What's going on?!

Merge

The top red line runs for 1.5 hours

1/3/20

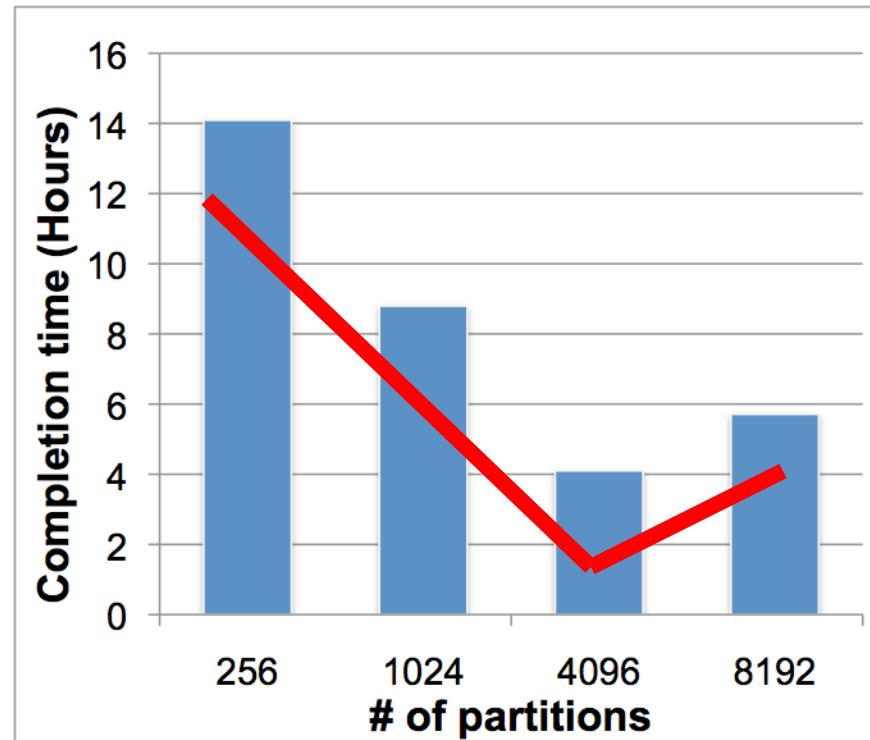
# Which one is better?



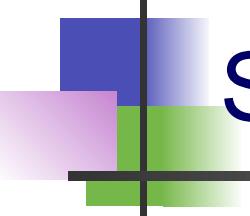
- How to decompose space?
- How to schedule?
- How to avoid memory overrun?

# Optimal Partitioning Plan Non-Trivial

- Fine grained partitions
  - Less data = Less skew
  - Framework overhead dominates
- Finding optimal point is time consuming
- No guarantee of successful merge phase



***Can we find a good partitioning plan without trial and error?***



# Skew Reduce Framework

- User provides three functions

`process :: <Seq. of  $T$ > → < $F$ , Seq. of  $S$ >`

`merge :: < $F$ ,  $F$ > → < $F$ , Seq. of  $S$ >`

`finalize :: < $F$ ,  $S$ > → <Seq. of  $Z$ >`

- Plus (optionally) two cost functions

$C_p :: (S, \alpha, B) \rightarrow \mathbb{R}$

$C_m :: (S, \alpha, B) \times (S, \alpha, B) \rightarrow \mathbb{R}$

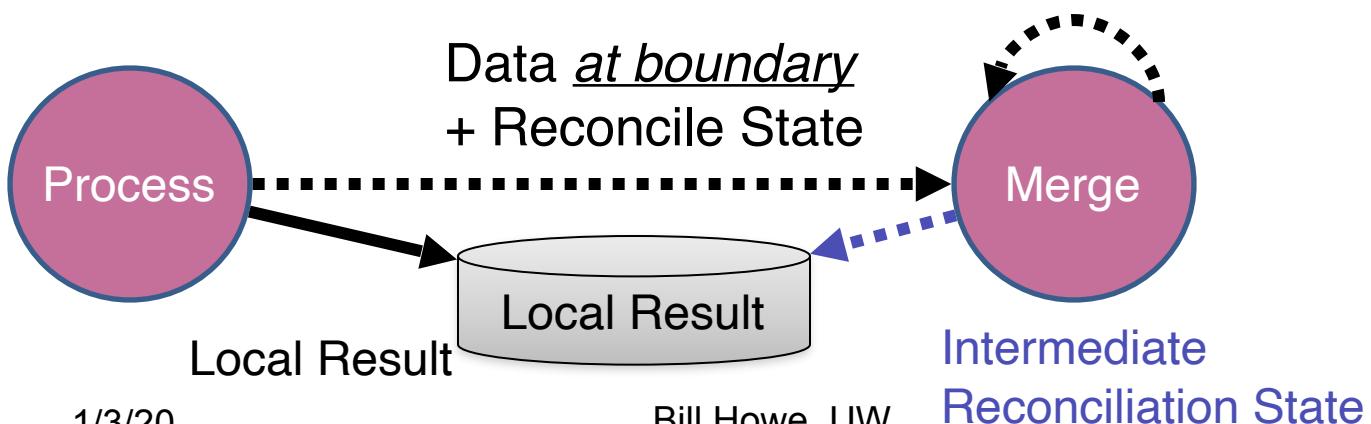
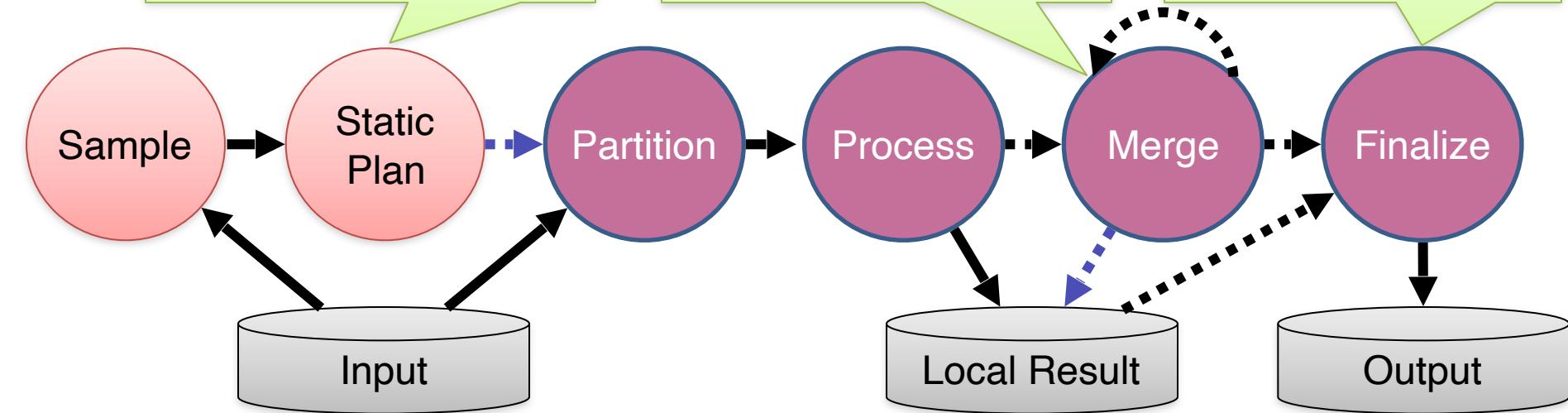
$S$  = sample of the input block;  $\alpha$  and  $B$  are metadata about the block

# Skew Reduce Framework

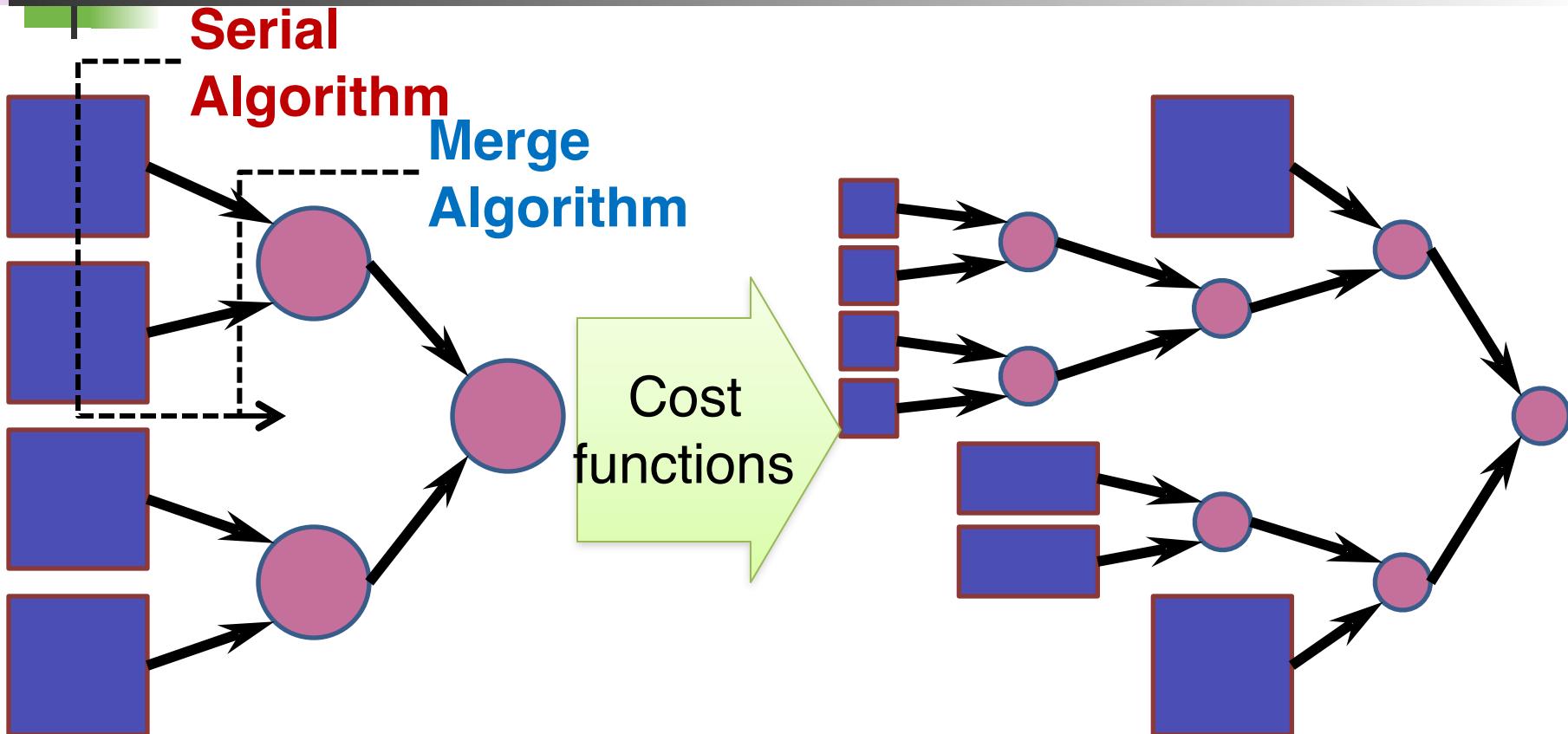
- User supplied cost function
- Could run in offline

- Hierarchically reconcile local result

- Update local result and produce final result

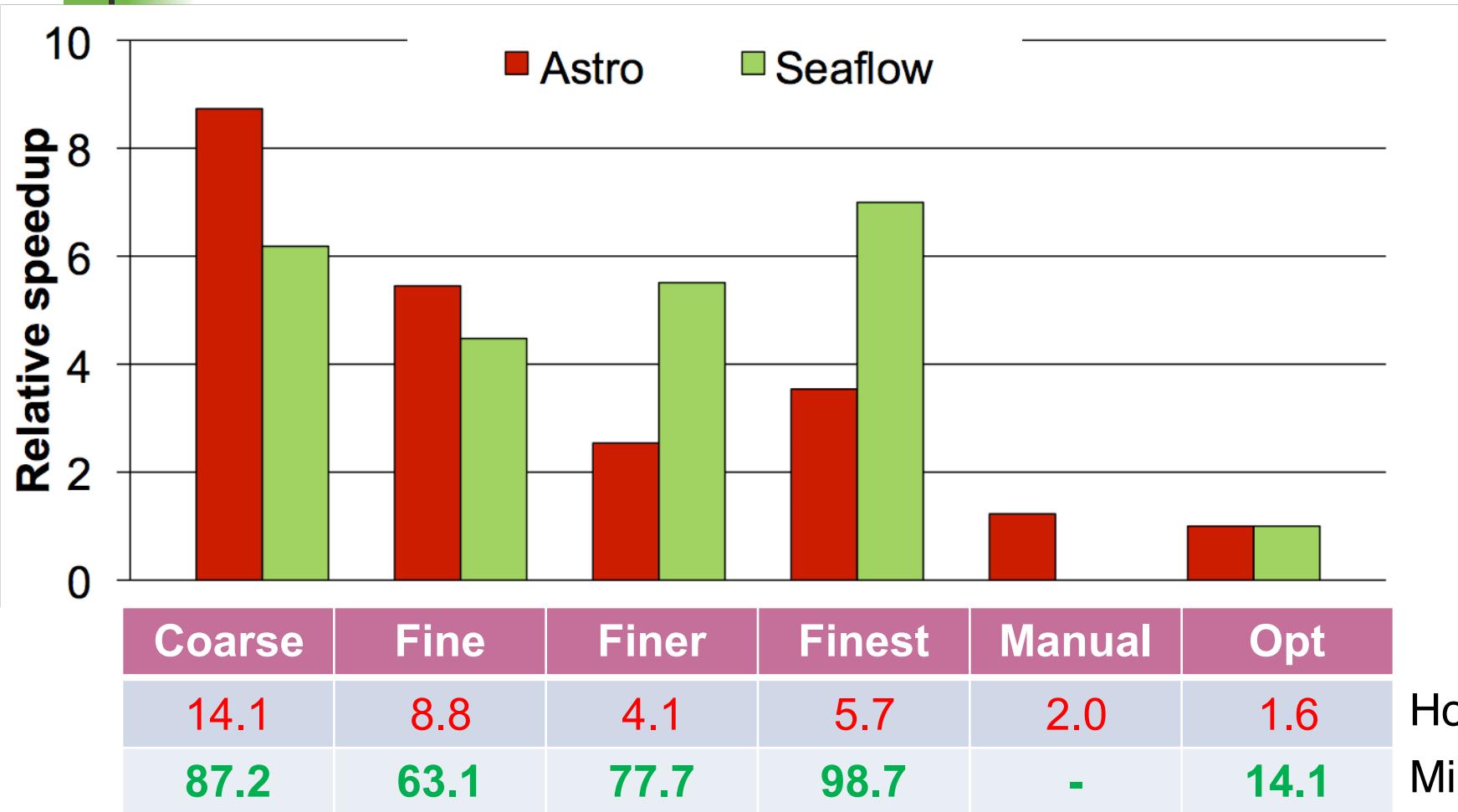


# Contiribution: SkewReduce



- Two algorithms: Serial/Merge algorithm
- Two cost functions for each algorithm
- Find a good partition plan and schedule

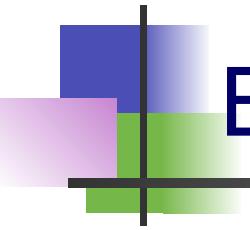
# Does SkewReduce work?



- Static plan yields 2 ~ 8 times faster running time

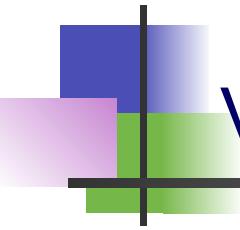


*Data-Intensive  
Scalable Science*



# BACKUP SLIDES





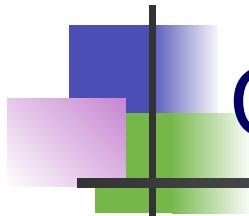
# Visualization + Data Management

*We can no longer afford two separate systems*

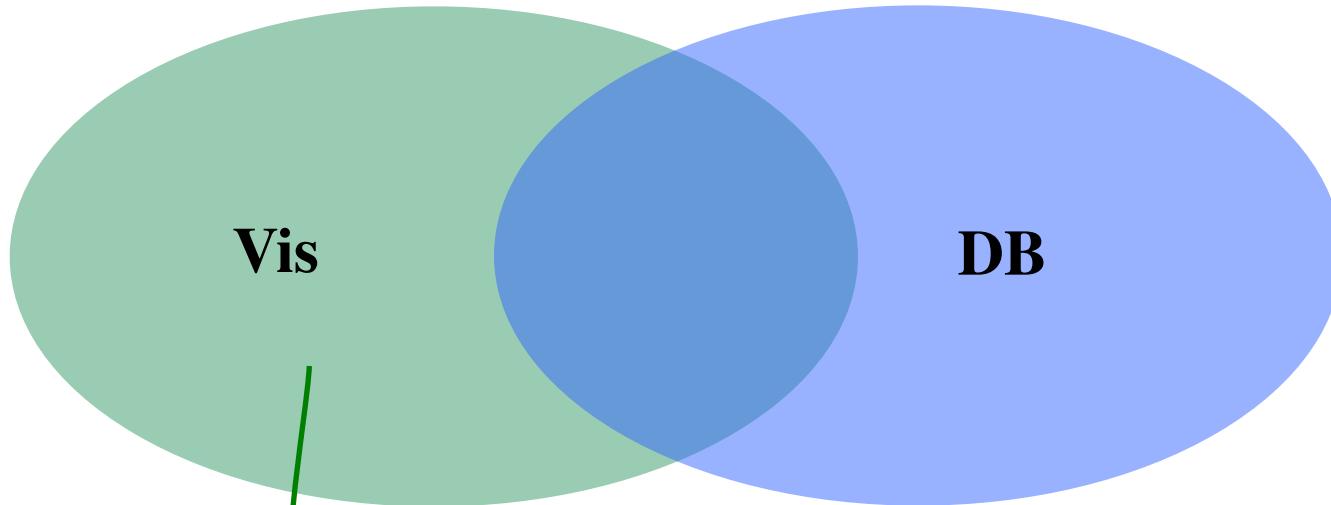
“Transferring the whole data generated ... to a storage device or a visualization machine could become a serious bottleneck, because I/O would take most of the ... time. A more feasible approach is to **reduce and prepare the data in situ** for subsequent visualization and data analysis tasks.”

-- SciDAC Review

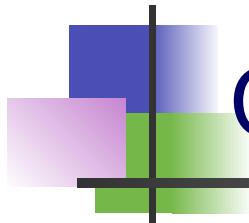




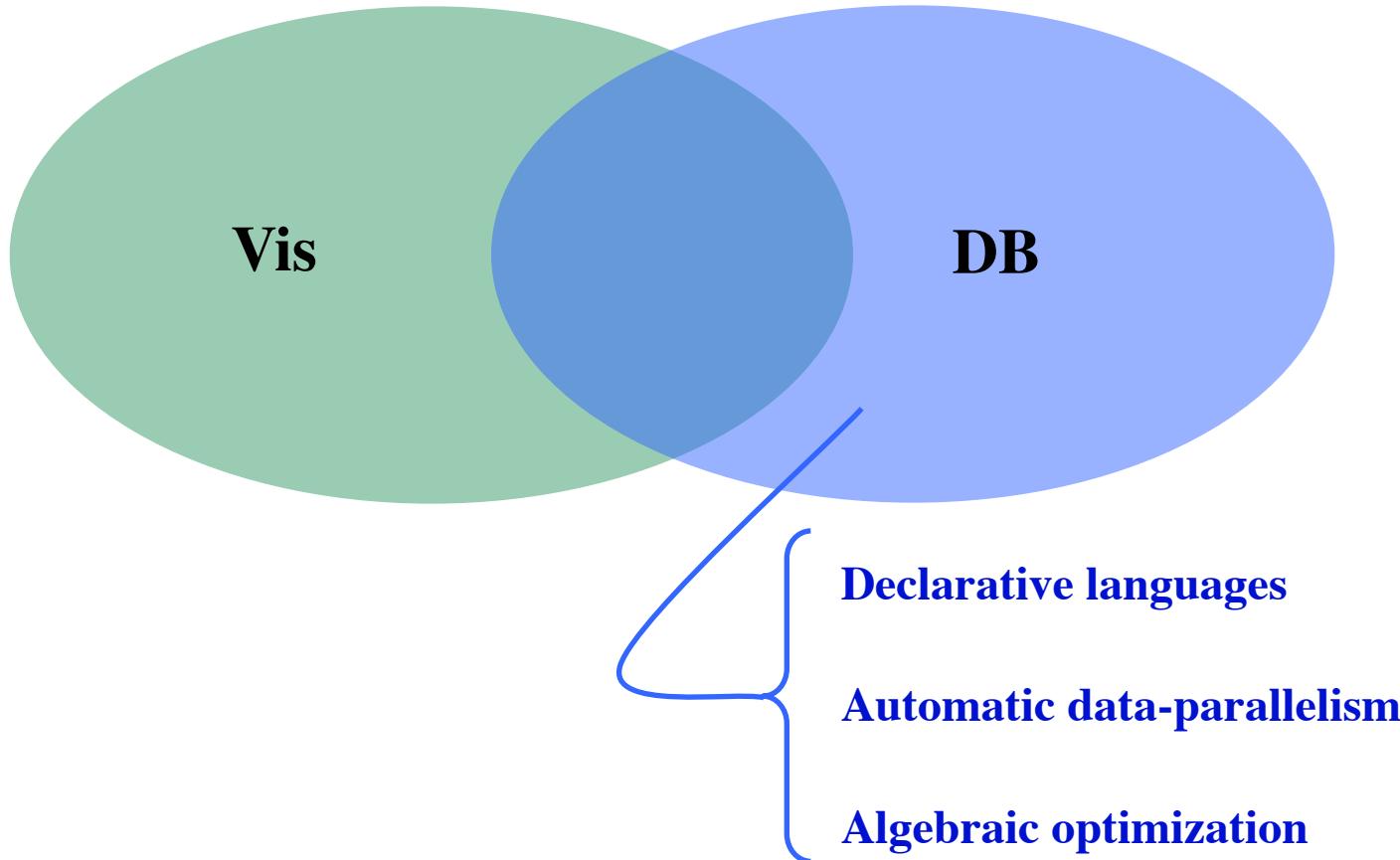
# Converging Requirements

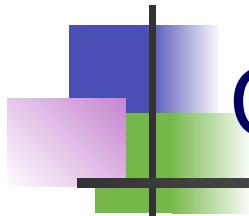


- Core vis techniques (isosurfaces, volume rendering, ...)
- Emphasis on interactive performance
- Mesh data as a first-class citizen

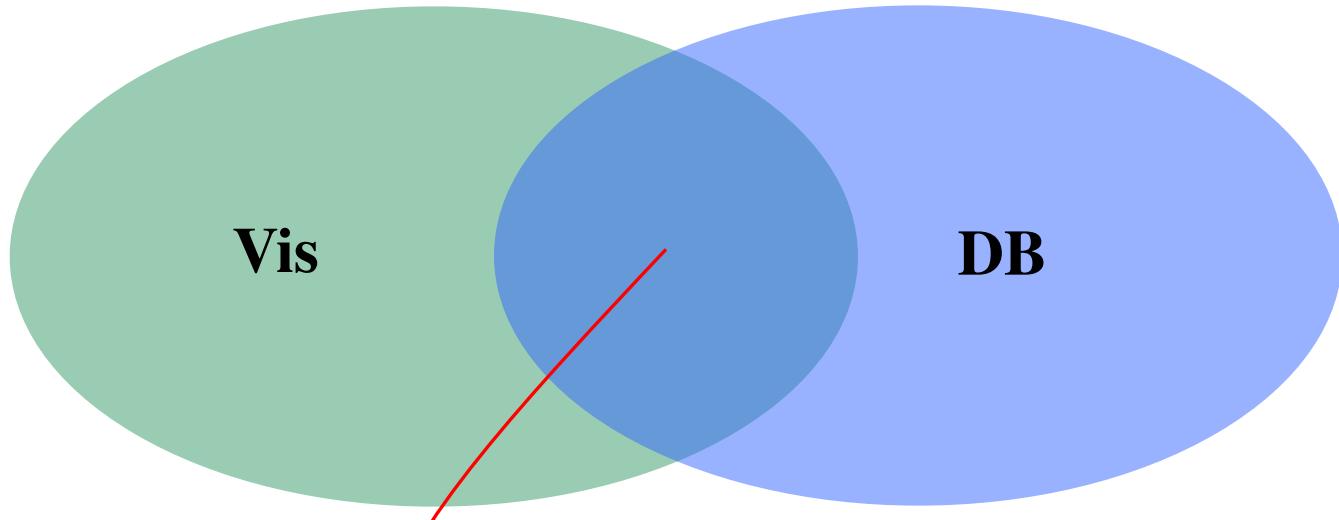


# Converging Requirements





# Converging Requirements



Vis: “Query-driven Visualization”

Vis: “In Situ Visualization”

Vis: “Remote Visualization”

DB: “Push the computation to the data”

# Desiderata for a “VisDB”

- New Data Model
  - Structured and Unstructured Grids
- New Query Language
  - Scalable grid-aware operators
  - Native visualization primitives
- New indexing and optimization techniques
- “Smart” Query Results
  - Interactive Apps/Dashboards

