

Big Telco, Bigger DW Demands: Moving Towards SQL-on-Hadoop

Keuntae Park



- IT Manager of SK Telecom, South Korea's largest wireless communications provider
- Work on commercial products (~'12)
 - T-FS: Distributed File System
 - Windows compatible layer on TimOS
 - T-MR: on-demand MapReduce service like E-MR
- Open source activity ('13~)
 - Committer of Apache Tajo project



Overview

- Background
 - Telco requirements
- Before Tajo
 - Commercial product
 - Open source (Hadoop) outsourcing
- After Tajo
 - Issues & solutions
 - Performance
- win-win between community and company
- Future Works



Telco data characteristics

- Huge amount of data
 - 40 TB/day (compressed)
 - 15 PB (estimated, end of 2014)
- Report & OLAP ad-hoc query
 - Filtering
 - Summary
 - BI tools

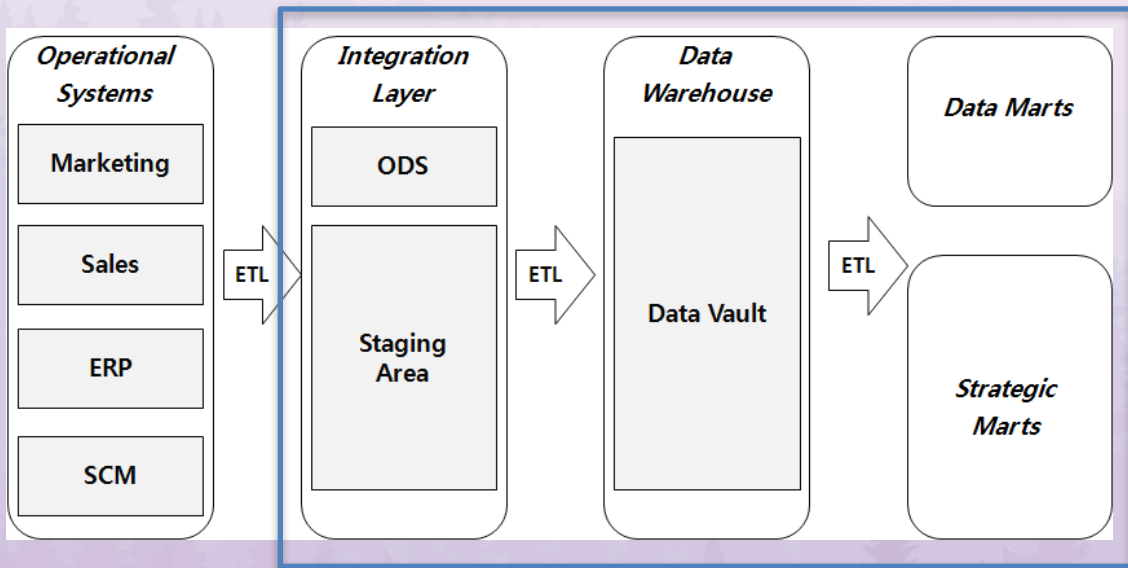


Requirements - different size, different speed

	<i>Filtering & aggregation</i>	<i>Summary</i>	<i>Data re-construction</i>	<i>BI report</i>	<i>Ad-hoc Query</i>
<i>Target</i>	<i>accumulated for 5 minutes</i>	<i>daily sum of filtered data</i>	<i>entire summary data</i>	<i>mart data</i>	<i>summary data</i>
<i>Frequency</i>	<i>every 5 minutes</i>	<i>daily or monthly</i>	<i>non-regularly (rare)</i>	<i>ah-hoc</i>	<i>ah-hoc</i>
<i>Amount of data</i>	<i>terabytes</i>	<i>hundreds of terabytes</i>	<i>petabytes</i>	<i>tens of gigabytes</i>	<i>tens of terabytes</i>
<i>Response time</i>	<i>within a minute</i>	<i>within a hour</i>	<i>no strict deadline</i>	<i>within two seconds</i>	<i>within a hour</i>



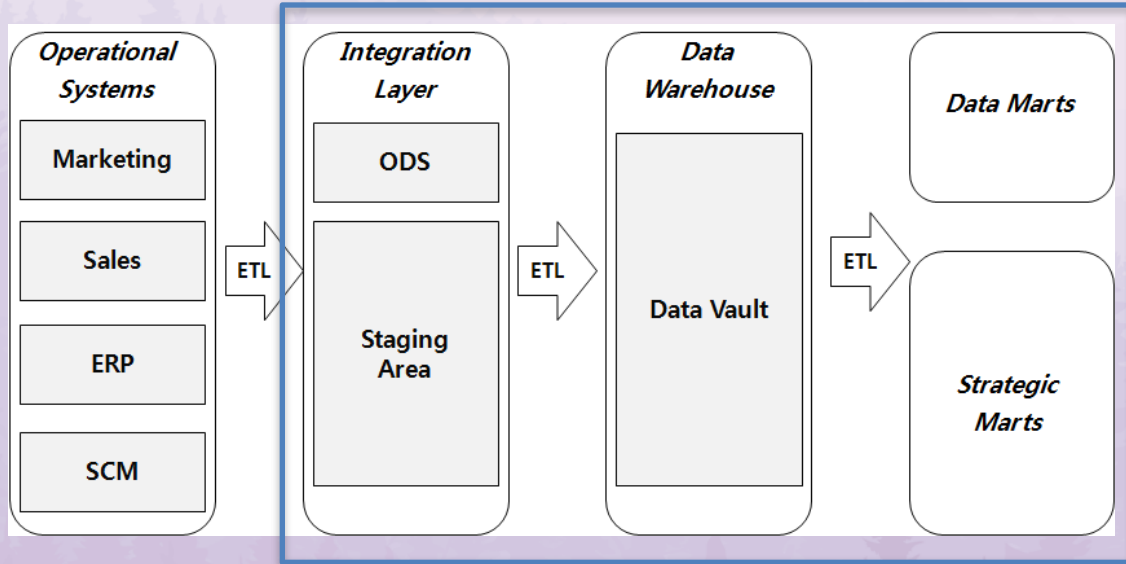
Previous approach - DBMS



based on MPP DBMS



Previous approach - DBMS



based on MPP DBMS

Too Expensive
Not Scalable



Previous approach - DBMS



Too Expensive

Not Scalable

Previous approach - DBMS



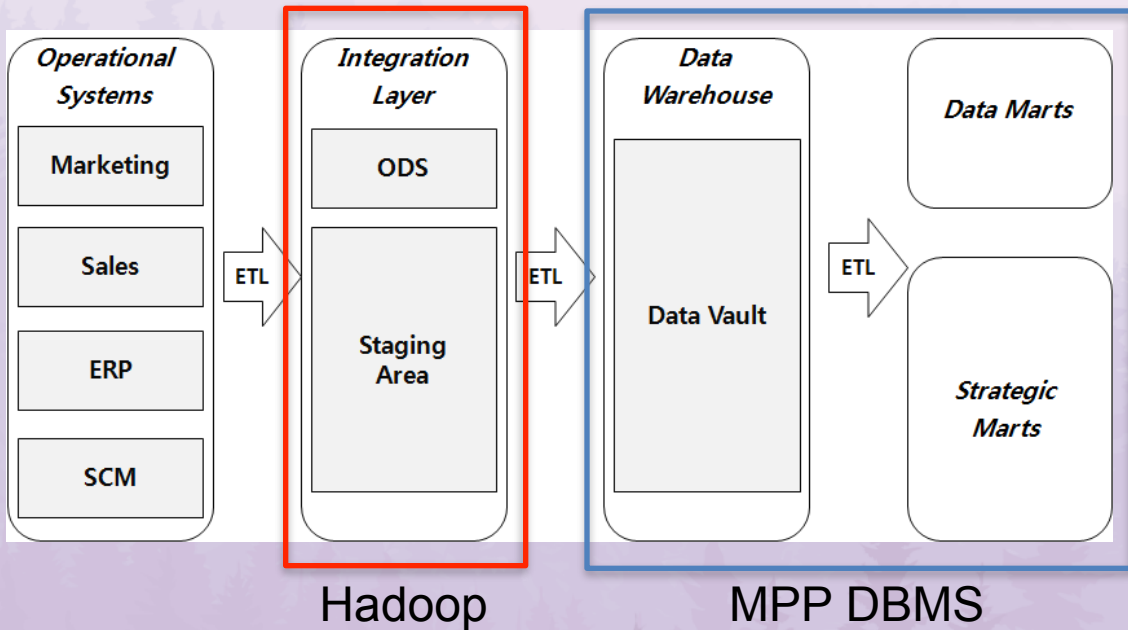
Data Man



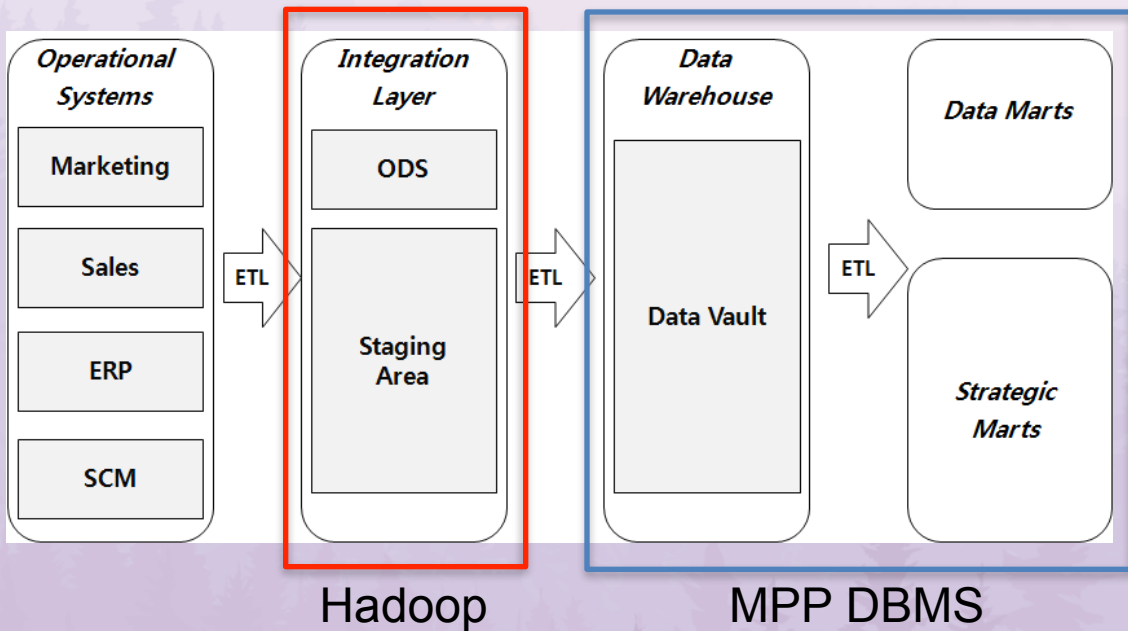
sive
ble

based on DBMS

Previous approach - Hadoop(MapReduce, Hive) + DBMS



Previous approach - Hadoop(MapReduce, Hive) + DBMS



Working
(but...)

Still has Problems

- Hadoop outsourcing
 - quality of outcome is not good (actually bad)
 - communication overhead
 - hard to reflect requirements on open source
- Data Warehouse and Mart becomes bigger



Solution - Tajo!!

- It can replace both DBMS and Hadoop
 - High throughput for batch processing
 - Low latency for ad-hoc queries
 - ANSI SQL compatible
- Can do by myself
 - very open community
 - easily make issues about what I really need
 - fast growing
 - issues solved very fast



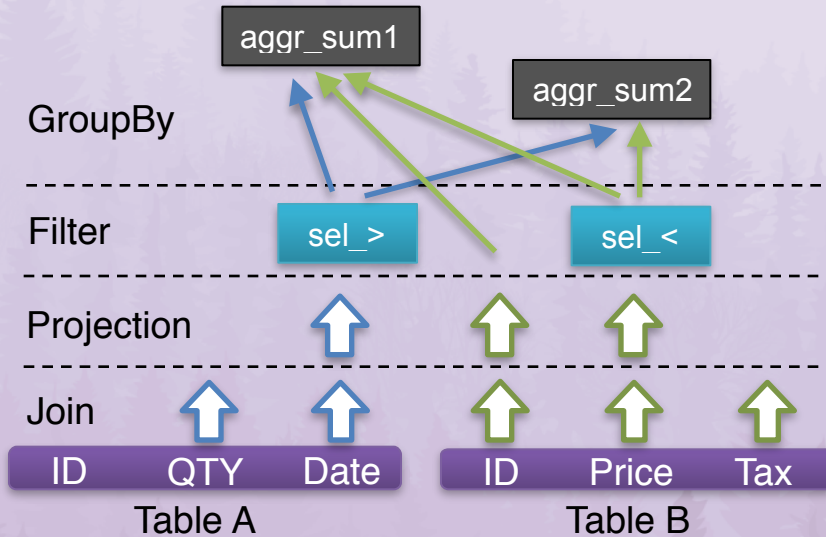
About Tajo

- Tajo (since 2010)
 - Big Data Warehouse System on Hadoop
 - Apache top-level project (entered the ASF in March 2013)
- Features
 - SQL standard compliance
 - Fully distributed SQL query processing
 - HDFS as a primary storage
 - Relational model (will be extended to nested model in the future)
 - ETL as well as low-latency relational query processing (100 ms ~)
- News
 - 0.2-incubating: released November 2013
 - graduation to top-level: April 2014



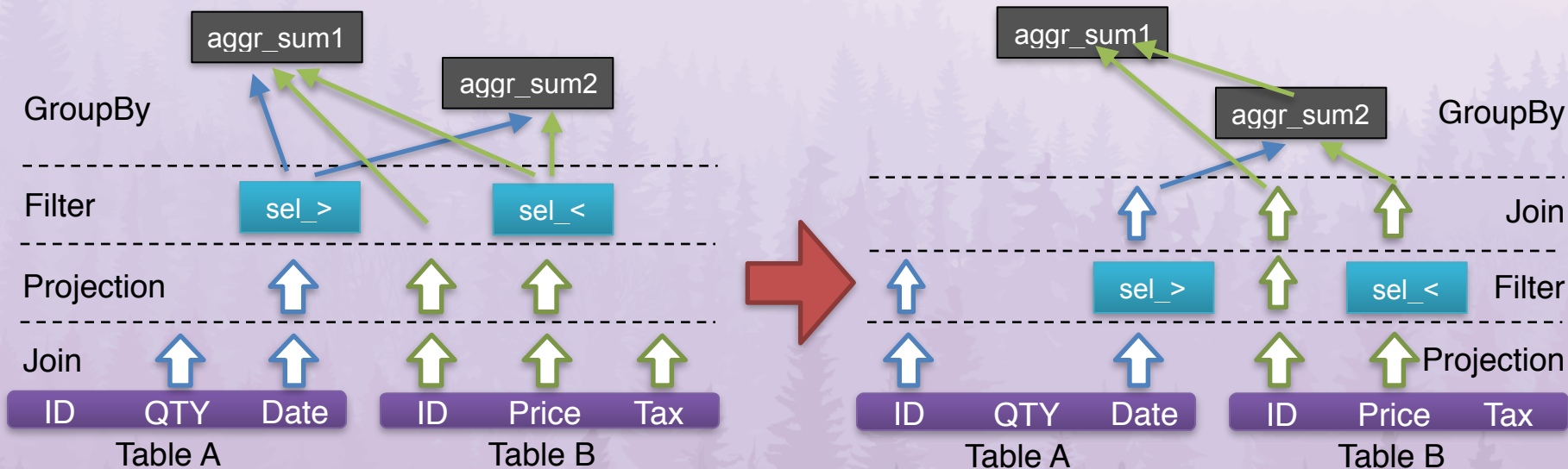
Tajo logical optimizer

- Cost-based join ordering
- Projection/Filter push down & Duplicated expression removal



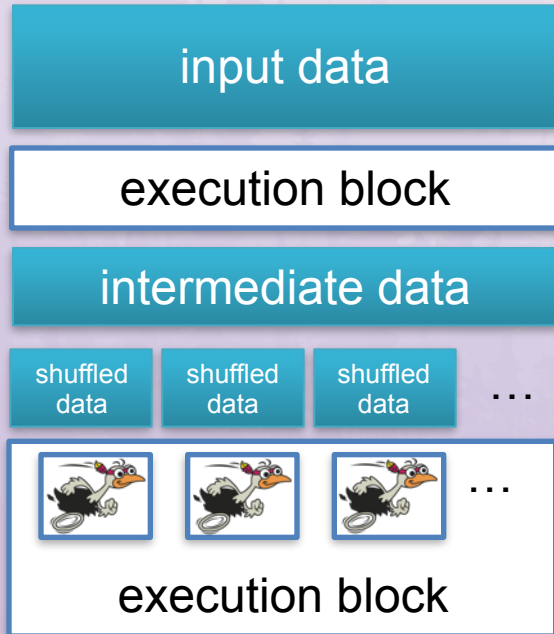
Tajo logical optimizer

- Cost-based join ordering
- Projection/Filter push down & Duplicated expression removal



Tajo progressive optimization

- dynamically adjust number of tasks



unknown priorly

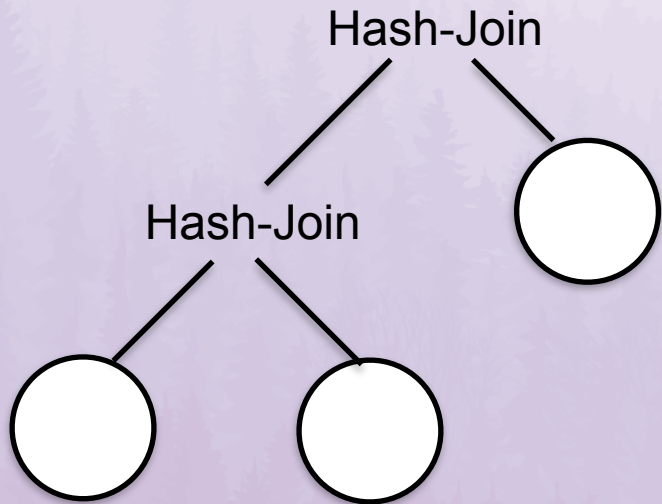
how many tasks
(and workers)?

- estimate data size at planning time
- check size and adjust plan at execution time
- shuffle intermediate data over workers uniformly



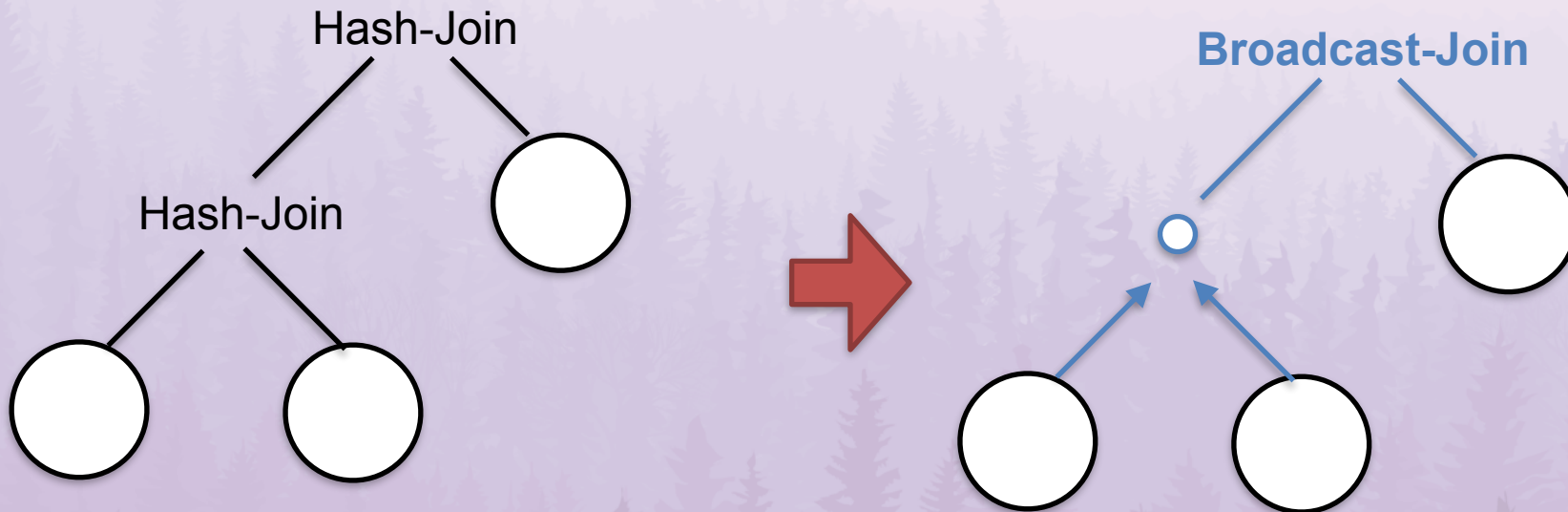
Tajo progressive optimization

- dynamically adjust join order or type



Tajo progressive optimization

- dynamically adjust join order or type



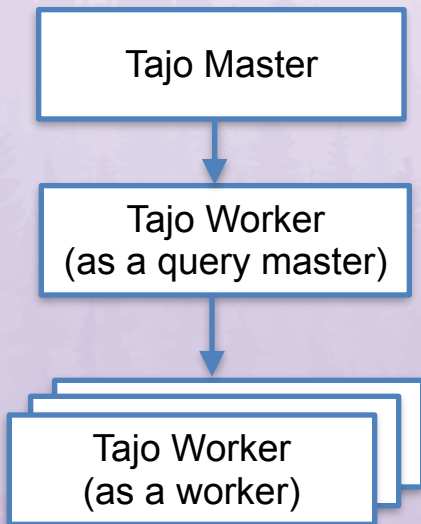
Tajo - what is improved past 9 months ?

- Resource Manager
- Scheduler & Storage Manager
- Data types & Functions
- SQL Interface
- Management



Tajo resource manager

- Fine resource allocation

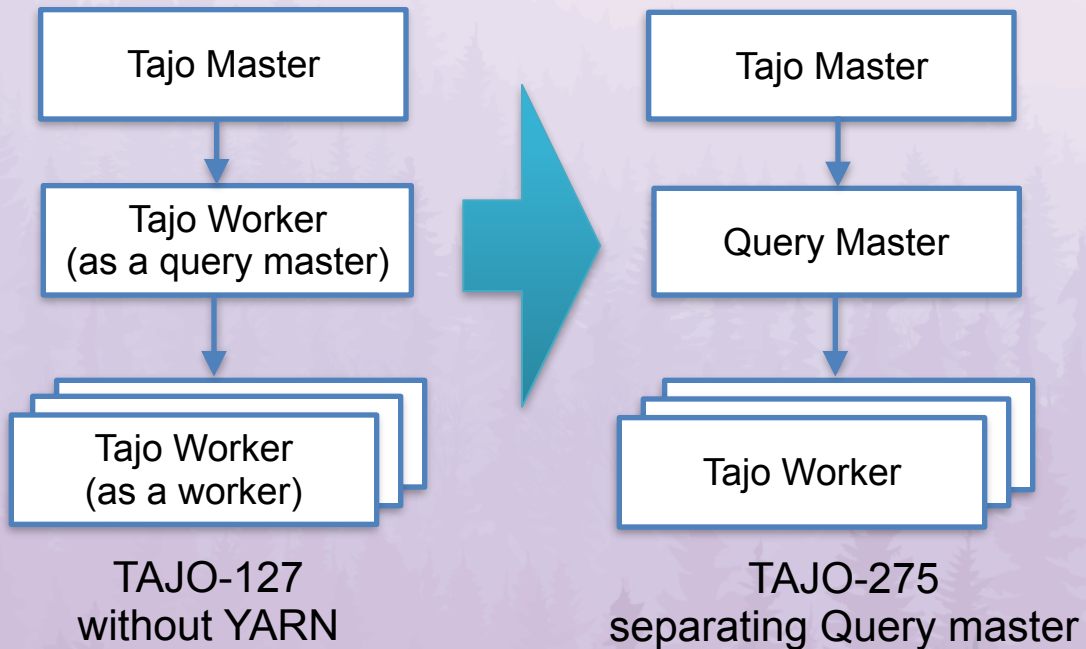


TAJO-127
without YARN



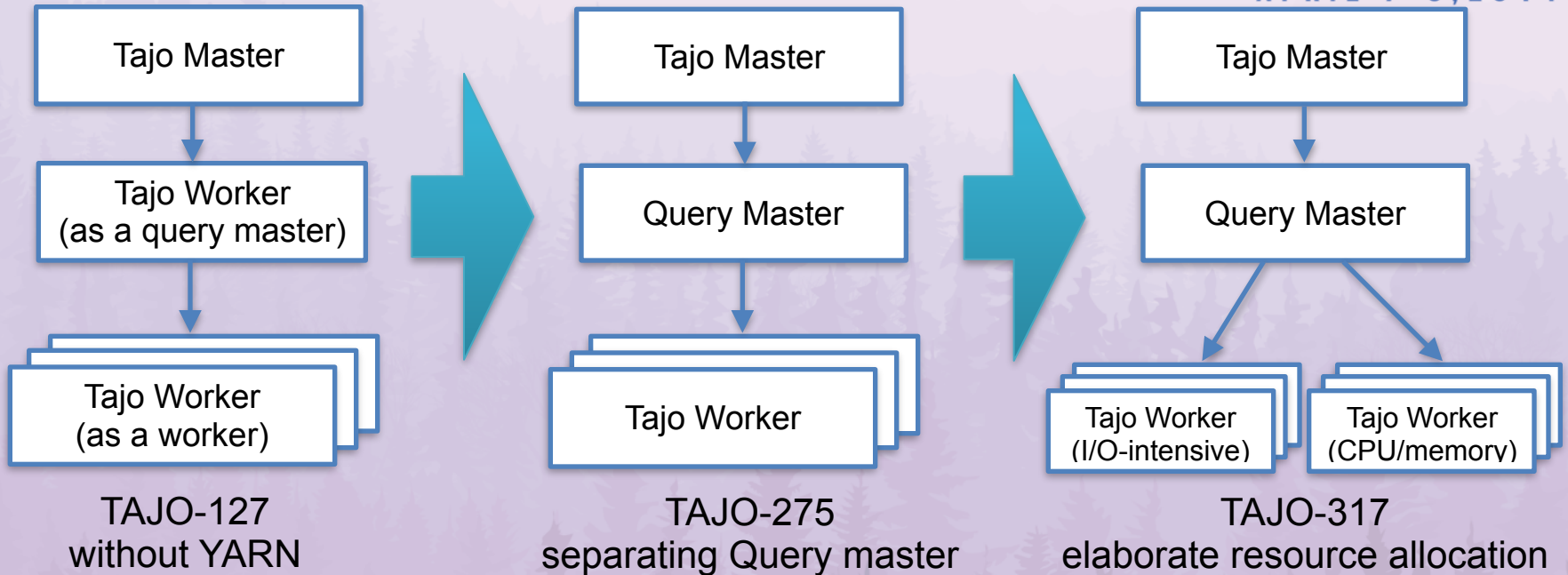
Tajo resource manager

- Fine resource allocation



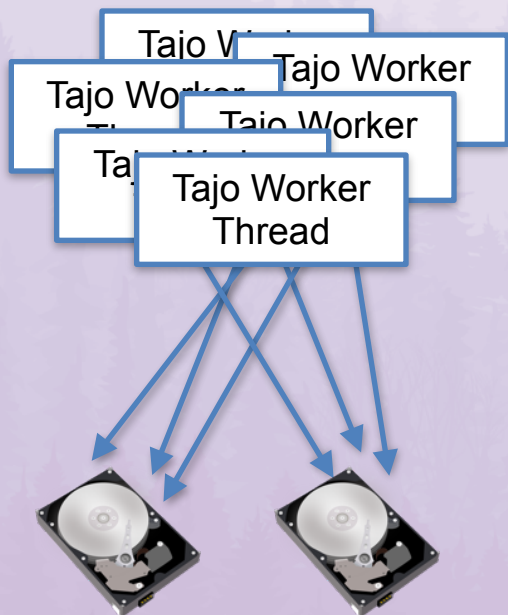
Tajo resource manager

- Fine resource allocation



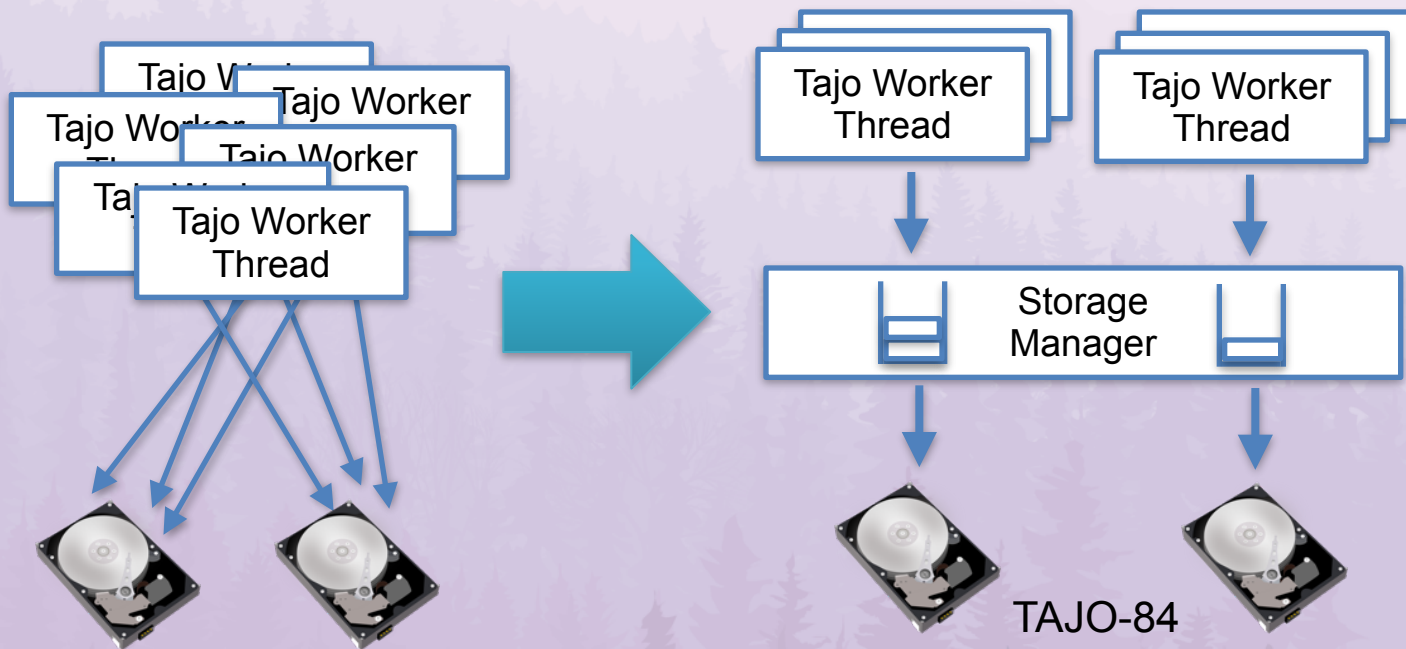
Scheduler & Storage manager

- disk-aware scheduling (volume info from HDFS-3672)



Scheduler & Storage manager

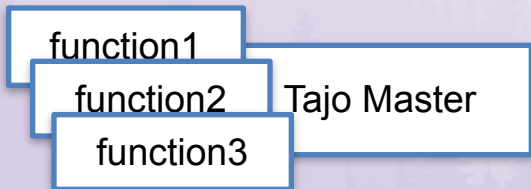
- disk-aware scheduling (volume info from HDFS-3672)



TAJO-84
considering disk load balance
TAJO-178
asynchronous scan

Functions & data types

- supporting more functions and UDFs

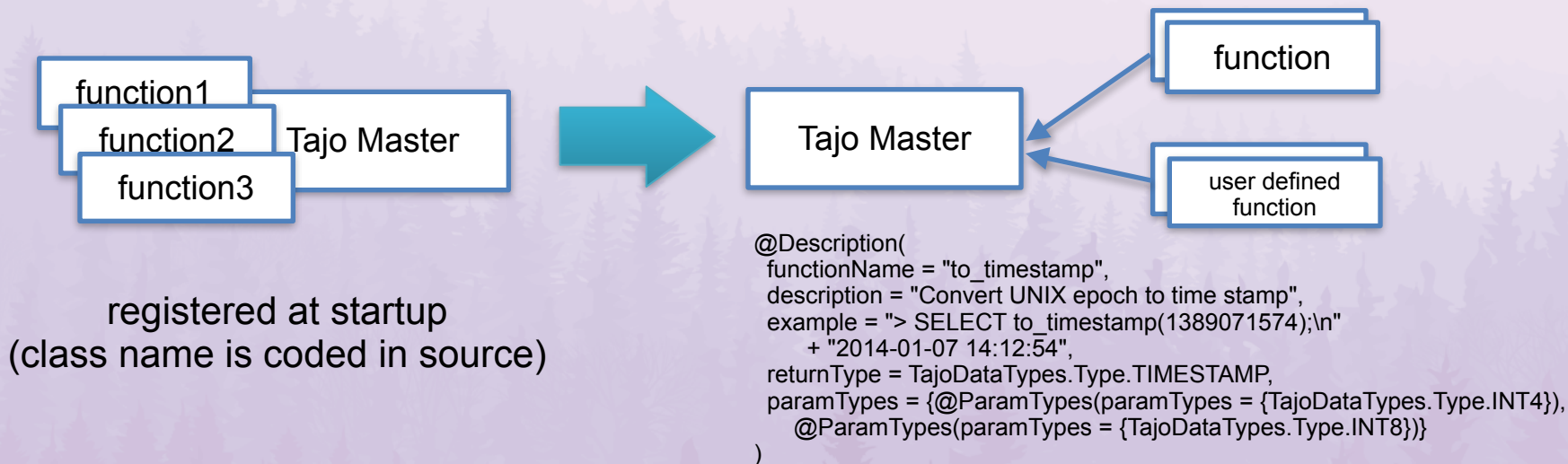


registered at startup
(class name is coded in source)



Functions & data types

- supporting more functions and UDFs

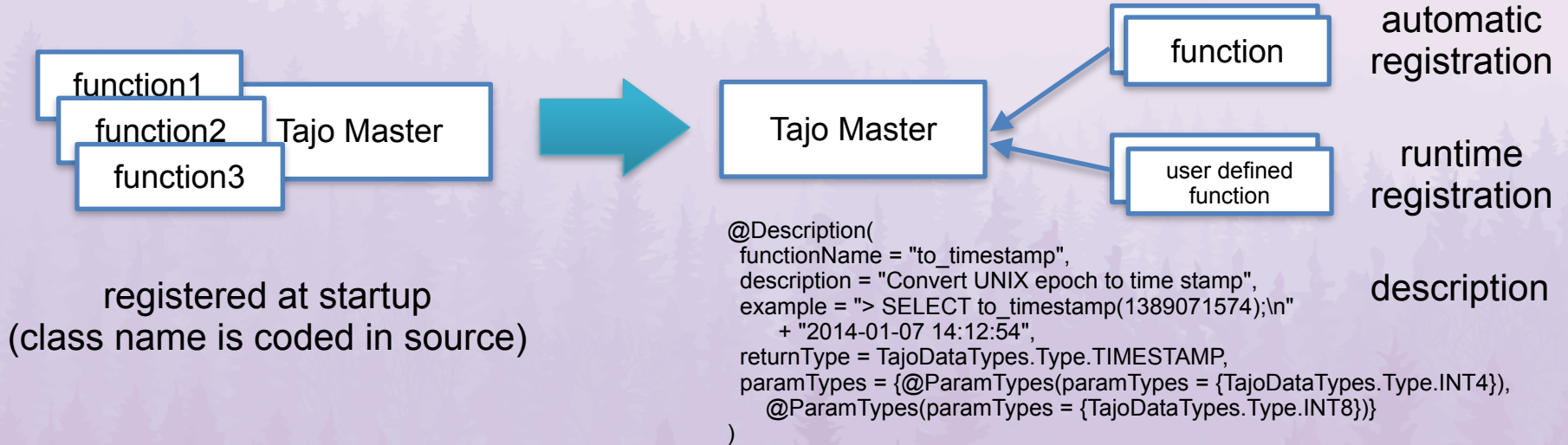


TAJO-408
Improve function system



Functions & data types

- supporting more functions and UDFs

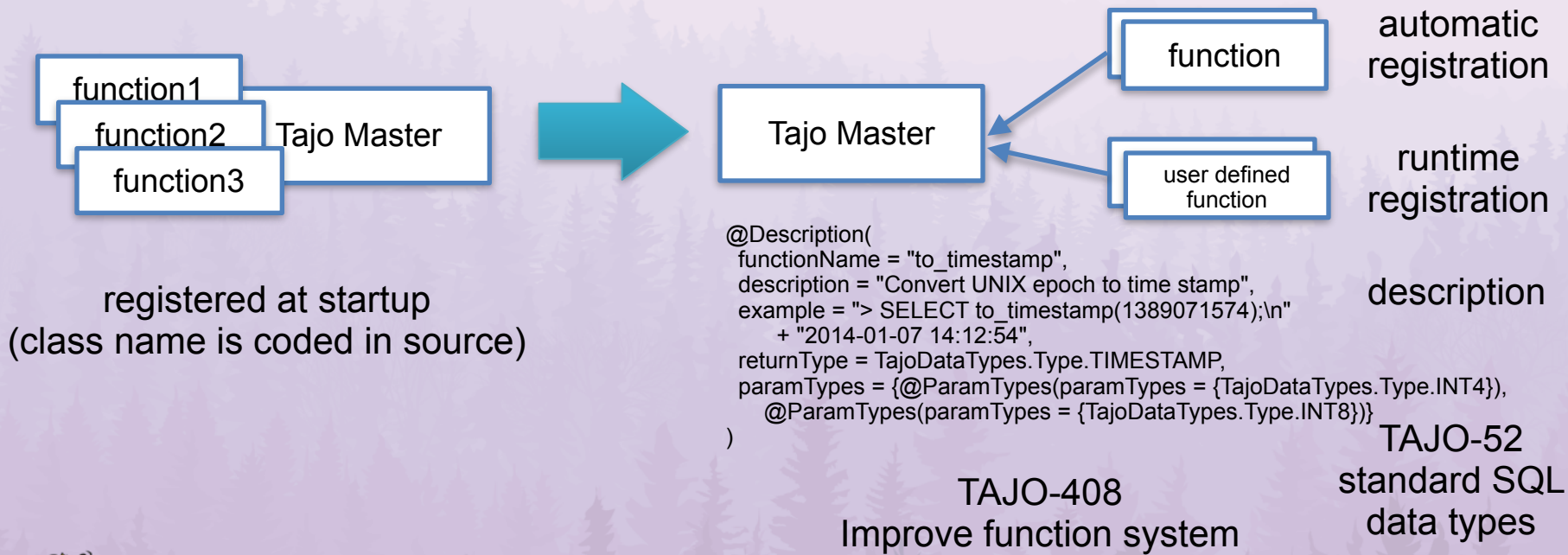


TAJO-408
Improve function system

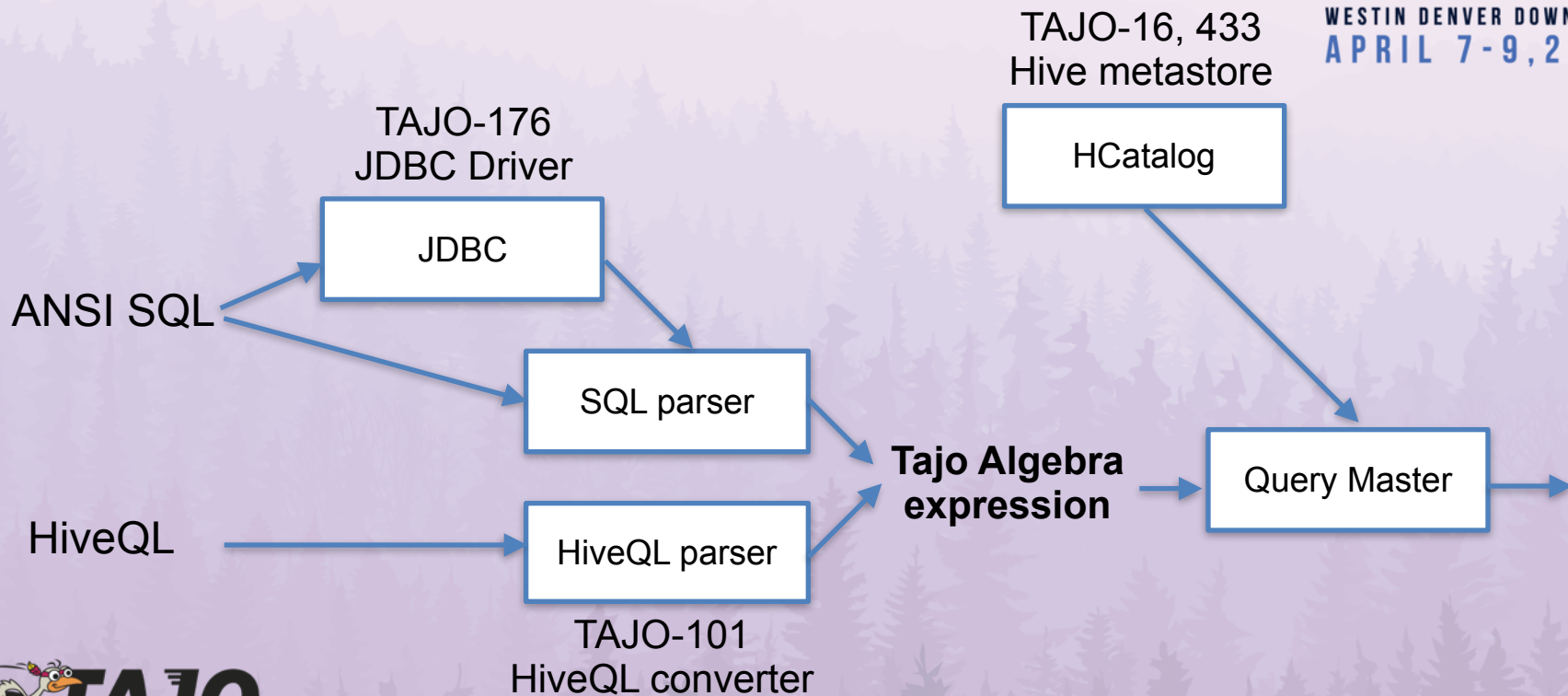


Functions & data types

- supporting more functions and UDFs



JDBC Driver, HCatalog



Management

TAJO Home Cluster Query Catalog Execute Query

Tajo Master: 50.1.102.122:26001

Query Master

Live:1, Dead: 0, QueryMaster Tasks: 0

Live QueryMasters

No	QueryMaster	Client Port	Running Query	Heap(free/total/max)	Heartbeat	Status
1	ceo-tajo02:28093	28092	0	847/921/921 MB	9.0 sec	LIVE

Worker

Live:6, Dead: 0

Live Workers

No	Worker	PullServer Port	Running Tasks	Memory Resource (used/total)	Disk Resource (used/total)	Heap (free/total/max)	Heartbeat	Status
1	ceo-tajo02:28091	45921	0	0/12288	0.0/6.0	47135/59882/59882 MB	.0 sec	LIVE
2	ceo-tajo03:28091	51067	0	0/12288	0.0/6.0	50460/59928/59928 MB	.0 sec	LIVE
3	ceo-tajo04:28091	37758	0	0/12288	0.0/6.0	54665/59877/59877 MB	.0 sec	LIVE
4	ceo-tajo05:28091	52609	0	0/12288	0.0/6.0	55691/59933/59933 MB	.0 sec	LIVE
5	ceo-tajo06:28091	48598	0	0/12288	0.0/6.0	41604/59890/59890 MB	9.0 sec	LIVE
6	ceo-tajo07:28091	51590	0	0/12288	0.0/6.0	54461/59876/59876 MB	.0 sec	LIVE

Dead Workers

No Dead Workers

TAJO-239 Improving Web UI



Management

Tajo Master: 50. Tajo Worker: [ceo-tajo02:28093](#)

Query Master

Live:1, Dead: 0, QueryMa: [q_1395132712581_0015](#) [\[Query Plan\]](#)

Live QueryMasters

No	Queue	ID	State	Started	Finished	Running time	Progress	Tasks
1	ceo-tajo02:28093	eb_1395132712581_0015_000001	SUCCEEDED	2014-03-18 17:55:35	2014-03-18 17:55:35	.0 sec	100.0%	1/1

Logical Plan

Worker

Live:6, Dead: 0

Live Workers

No	Worker
1	ceo-tajo02:28091
2	ceo-tajo03:28091
3	ceo-tajo04:28091
4	ceo-tajo05:28091
5	ceo-tajo06:28091
6	ceo-tajo07:28091

Dead Workers

No Dead Workers

Optimization Log:

```
SCAN(0) on table1
=> target list: table1.id (INT4), table1.name (TEXT), table1.score (FLOAT4), table1.type (TEXT)
=> out schema: {(4) table1.id (INT4),table1.name (TEXT),table1.score (FLOAT4),table1.type (TEXT)}
=> in schema: {(4) table1.id (INT4),table1.name (TEXT),table1.score (FLOAT4),table1.type (TEXT)}
```

Distributed Query Plan

Execution Block Graph (TERMINAL - eb_1395132712581_0015_000002)

```
|-eb_1395132712581_0015_000002
|-eb_1395132712581_0015_000001
```

Block Id: eb_1395132712581_0015_000001 [ROOT]

TAJO-564
Execution block progress



Management

Tajo Master: 50.0.0.1

Tajo Worker: ceo-tajo02:28093

Query Master: q_1395132712581_0015

Tajo Worker: ceo-tajo02:28093

Live QueryMasters

No	Queue	ID
1	ceo-tajo02:28093	eb_1395132712581_0015_000001

Logical Plan

```

SCAN(0) on table1
=> target list: table1.id (INT4), table1.name (TEXT), table1.score (FLOAT4), table1.type (TEXT)
=> out schema: ((4) table1.id (INT4), table1.name (TEXT), table1.score (FLOAT4), table1.type (TEXT))
=> in schema: ((4) table1.id (INT4), table1.name (TEXT), table1.score (FLOAT4), table1.type (TEXT))
    
```

Worker

Live:6, Dead:0

Live Workers

No	Worker
1	ceo-tajo02:28091
2	ceo-tajo03:28091
3	ceo-tajo04:28091
4	ceo-tajo05:28091
5	ceo-tajo06:28091
6	ceo-tajo07:28091

Dead Workers

No Dead Workers

Execution Block Graph

```

|-eb_1395132712581_0015
|-eb_1395132712581_0
    
```

Distributed Query Plan

No	id	Status	Progress	Started	Running Time	Host
1	t_1395132712581_0015_000001_000000	SUCCEEDED	100.0%	2014-03-18 17:55:35	34 ms	50.1.102.123

Block Id: eb_1395132712

TAJO-589
Task progress



Management

The screenshot displays the Tajo web interface with several panels:

- Query Master:** Shows job details for ID `q_1395132712581_0015`.
- Live Query Masters:** A table with columns 'No' and 'Queue'. It lists one entry: `1 ceo-tajo02:28093`.
- Worker:** Shows 'Live:6, Dead:0'.
- Live Workers:** A table with columns 'No' and 'Worker'. It lists six entries: `1 ceo-tajo02:28091`, `2 ceo-tajo03:28091`, `3 ceo-tajo04:28091`, `4 ceo-tajo05:28091`, `5 ceo-tajo06:28091`, and `6 ceo-tajo07:28091`.
- Dead Workers:** Shows 'No Dead Workers'.
- Logical Plan:** Contains the text: `SCAN(0) on table1`, `=> target list: table1`, `=> out schema: ((4) ta`, and `=> in schema: ((4) tab`.
- Distributed Query Plan:** Shows 'Execution Block Graph' with entry `1 -eb_1395132712581_0015`.
- Block Id:** `eb_1395132712581_0015`.
- Job Detail Panel (Right):** Shows details for `Tajo Worker: ceo-tajo02:28093` and `eb_1395132712581_0015_000001`. It includes a table with columns: ID, Progress, State, Launch Time, Finish Time, Running Time, Host, Shuffles, Data Locations, Fragment, Input Statistics, Output Statistics, and Fetches.

TAJO-468 Task detail info



Management

The image shows a composite of Tajo web interface elements and a terminal window. The web interface includes sections for 'Tajo Master: 50...', 'Query Master', 'Live QueryMasters', 'Worker', 'Live Workers', and 'Dead Workers'. The terminal window shows the following commands and output:

```
[hadoop@ceo-tajo01 ~]$ tajo admin -list
QueryId          State      StartTime
-----
q_1395186212310_0040  RUNNING  2014-03-19 11:03:21

[hadoop@ceo-tajo01 ~]$ tajo admin -kill
```

The terminal also shows a 'Logical Plan' section with the following details:

```
SCAN(0) on table1
=> target list: table1
=> out schema: ((4) ta
=> in schema: ((4) tab
```

Below the terminal, there is a slide titled 'Tajo management' with a small screenshot of the Tajo web interface.

TAJO-474
Task admin utility



And lots of Performance enhancement

TAJO-725 Broadcast JOIN should supports multiple tables

TAJO-717 Improve file splitting for large number of splits

TAJO-601 Improve distinct aggregation query processing

TAJO-584 Improve distributed merge sort

TAJO-36 Improve ExternalSortExec with N-merge sort and final pass omission

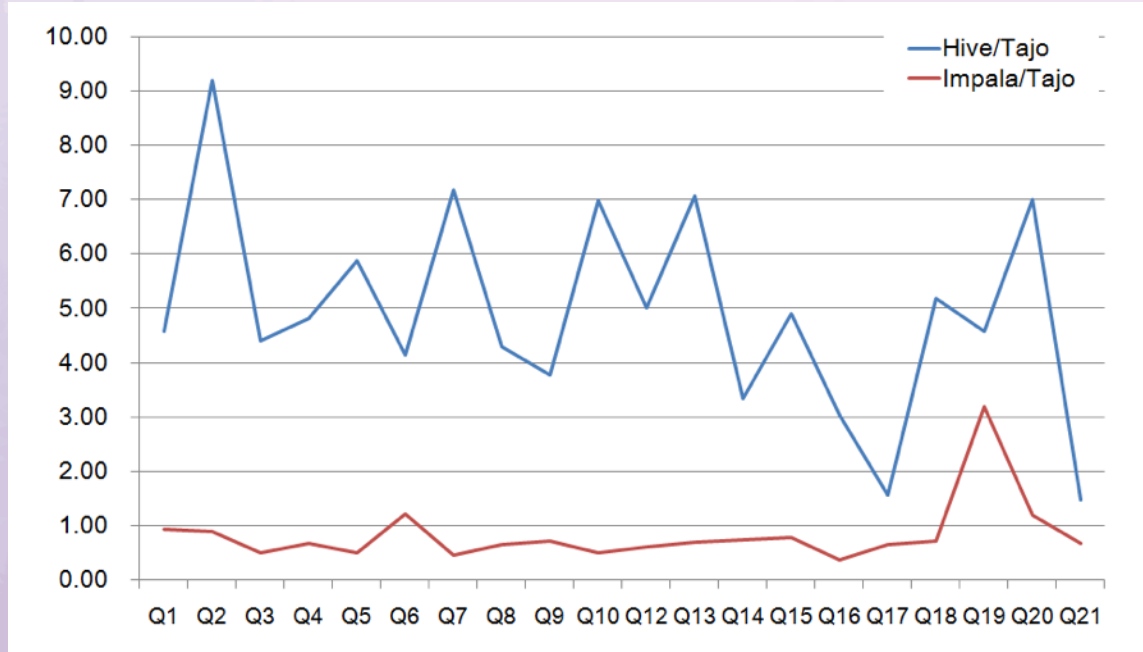
TAJO-345 MergeScanner should support projectable storages

...



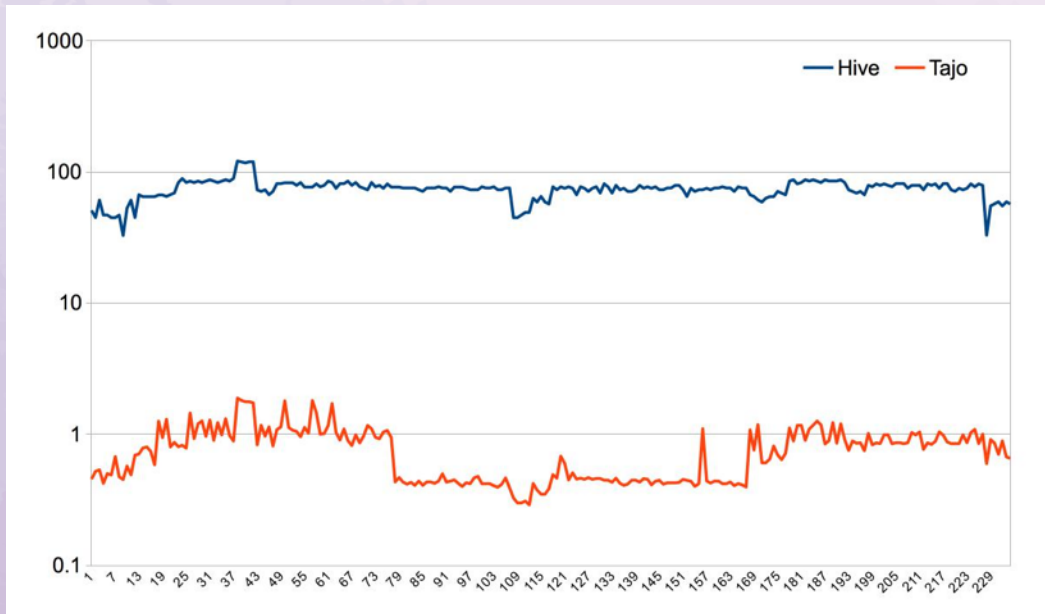
Performance

- TPC-H



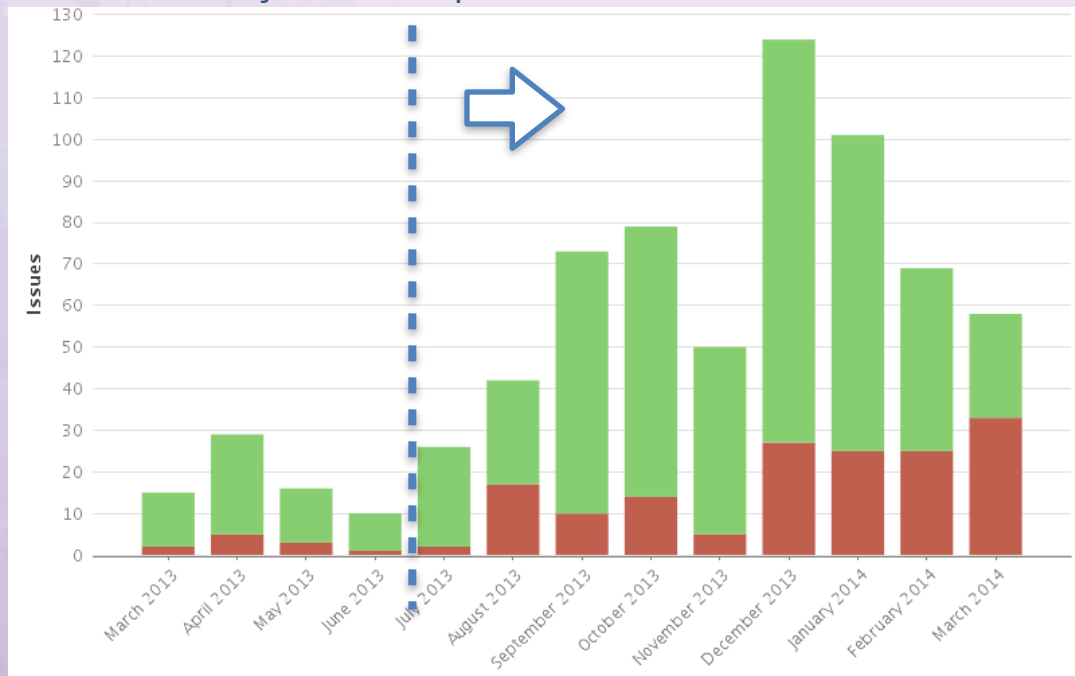
Performance

- OLAP reporting - relatively small data



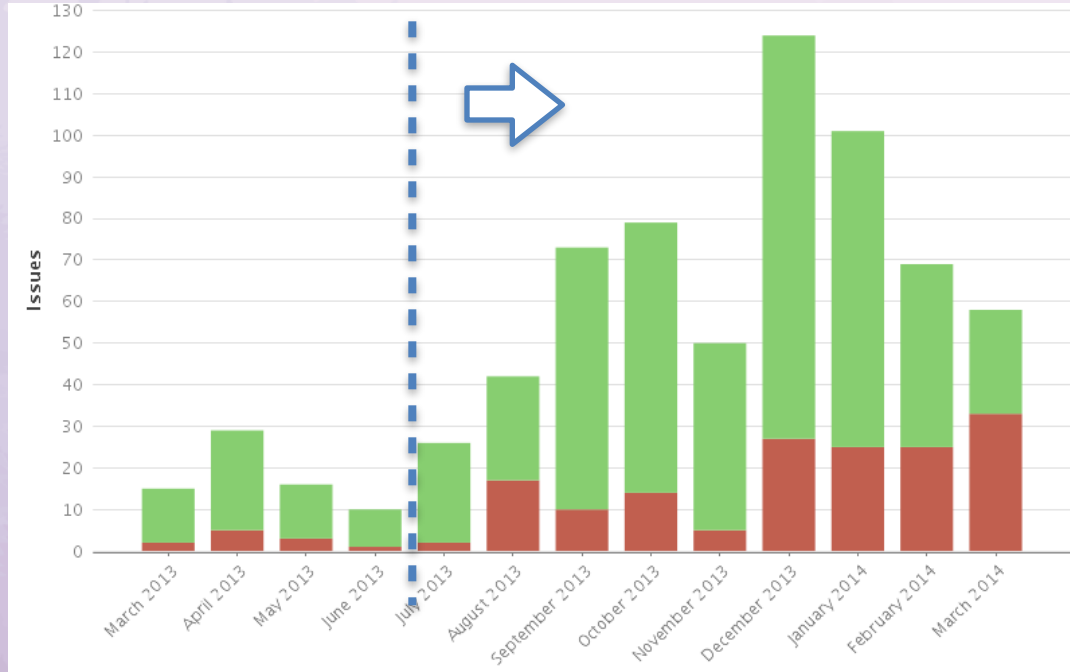
win-win between company and community

- Community boom up



win-win between company and community

- Community boom up



13 → 30



win-win between company and community

- Test in real working cluster
 - Mainly focusing on the scalability test & integration with existing IT systems
 - Finding bugs and function requirements, also



win-win between company and community

- Test in real working cluster
 - Mainly focusing on the scalability test & integration with existing IT systems
 - Finding bugs and function requirements, also

TAJO-691 HashJoin or HashAggregation is too slow if there is many unique keys

TAJO-675 maximum frame size of frameDecoder should be increased

TAJO-673 Assign proper number of tasks when inserting into partitioned table

TAJO-650 Repartitioner::scheduleHashShuffledFetches should adjust the number of tasks

TAJO-647 Work unbalance on disk scheduling of DefaultScheduler

TAJO-292 Too many intermediate partition files

TAJO-283 Add table partitioning

TAJO-592 HCatalogStore should supports RCFile and default hive field delimiter.

...



win-win between company and community



APACHE CON
DENVER
WESTIN DENVER DOWNTOWN
APRIL 7-9, 2014



Presented For The Apache Foundation By
LINUX FOUNDATION

win-win between company and community



- efficient development and operation
- human networking
- brand value up - recruiting

Future Works

- Nested data model (parquet model)
- more SQL compatible
 - window functions, IN, EXIST
- Multi-tenancy
- push shuffle (no materialization)
 - use selectively between push and pull shuffle
 - push shuffle: performance
 - pull shuffle: resilience, schedulability



Q & A

- Getting Started
 - <http://tajo.apache.org/tajo-0.2.0-doc.html#GettingStarted>
- Checkout the development branch
 - <http://tajo.apache.org/downloads.html>
- Jira - Issue Tracker
 - <https://issues.apache.org/jira/browse/TAJO>
- Join the mailing list
 - dev@tajo.apache.org

