# Serving Large-scale Batch Computed Data with Project Voldemort
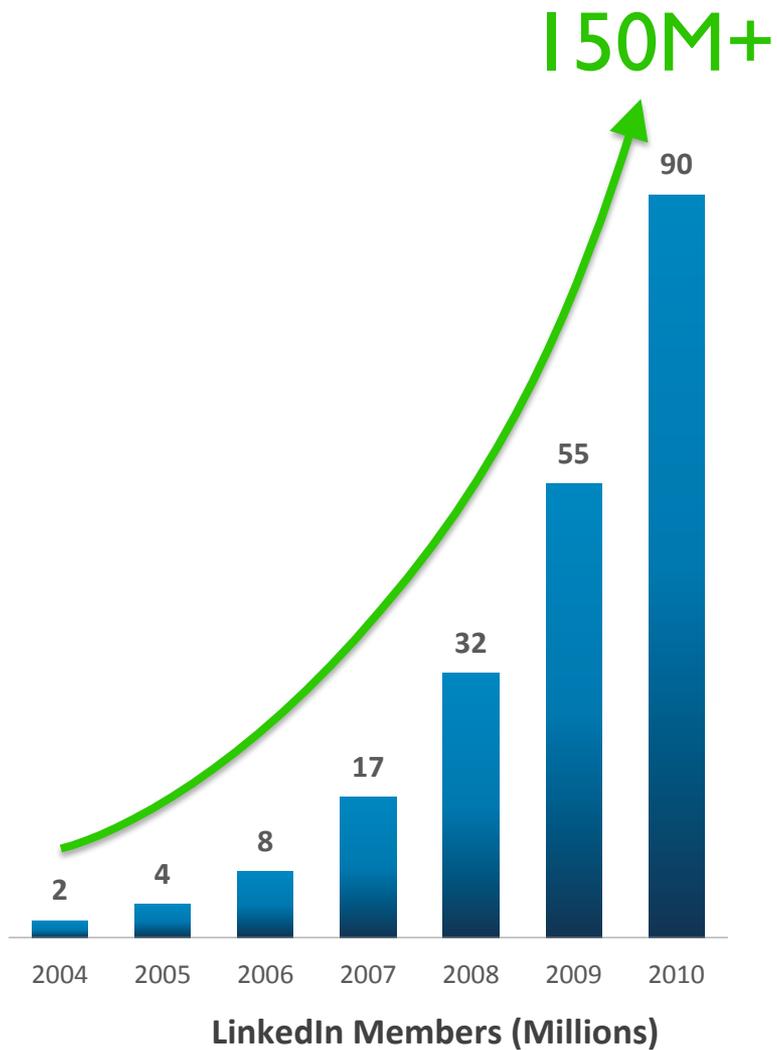
***Roshan Sumbaly***, Jay Kreps, Lei Gao, Alex Feinberg, Chinmay Soman, and Sam Shah

**LinkedIn**

# LinkedIn



**150M+**

LinkedIn Members (Millions)

| Year | Members |
|------|---------|
| 2004 | 2 |
| 2005 | 4 |
| 2006 | 8 |
| 2007 | 17 |
| 2008 | 32 |
| 2009 | 55 |
| 2010 | 90 |

**75%**
Fortune 100 Companies use LinkedIn to hire

**>2M**
Company Pages

**~4B**
Searches in 2011

# Data Features on LinkedIn

## People You May Know



## Viewers of this profile also viewed



## Related Searches



## Events you may be interested in



## LinkedIn Skills



## Jobs you may be interested in

## People You May Know



- Batch computed algorithms
  - MapReduce (Hadoop)
- Output
  - Large
  - Immutable
  - Key-value
  - Full refresh

How do we serve these massive outputs
to our 150 million members?

- Fast, available and elastic
- Bulk load massive data-sets
- Minimum time in error
- Easy to use
- Open-source

- **Fast, available and elastic**
- Bulk load massive data-sets
- Minimum time in error
- Easy to use
- Open-source

- **Distributed key-value system**

- Fast, available and elastic
- **Bulk load massive data-sets**
- Minimum time in error
- Easy to use
- Open-source

- Minimum performance impact during bulk loads
- Offload index construction to processing system

- Fast, available and elastic
- Bulk load massive data-sets
- **Minimum time in error**
- Easy to use
- Open-source

- Error in algorithm → Bulk load bad data → Bad state till next push
- Quick rollback capability

- Fast, available and elastic
- Bulk load massive data-sets
- Minimum time in error
- **Easy to use**
- Open-source

```
job.class=com.linkedin.jobs.BuildAndPushJob

build.input.path=/algorithm/output
push.store.name=people-you-may-know
push.cluster=tcp://testing-cluster-url:6666
build.replication.factor=1
```

- Fast, available and elastic
- Bulk load massive data-sets
- Minimum time in error
- **Easy to use**
- Open-source

```
job.class=com.linkedin.jobs.BuildAndPushJob

build.input.path=/algorithm/output
push.store.name=people-you-may-know
push.cluster=tcp://production-cluster-url:6666
build.replication.factor=2
```

- Fast, available and elastic
- Bulk load massive data-sets
- Minimum time in error
- Easy to use
- Open-source

- Apache License v2.0
- Project Voldemort – http://project-voldemort.com

- Background – Voldemort architecture
- Custom Voldemort storage engine
  - Minimal impact on live system
  - Fast rollback
  - Fast lookups
  - Easy rebalancing
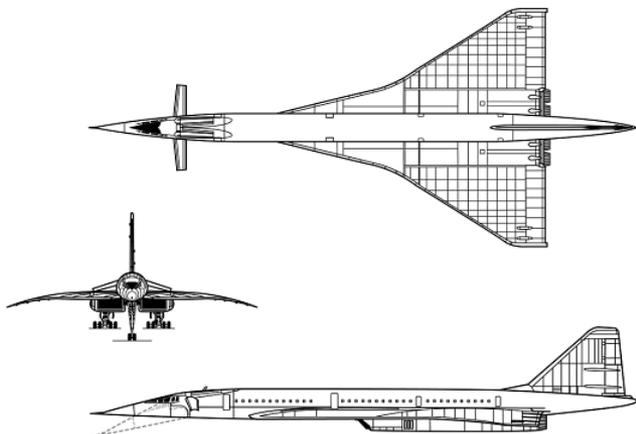- Performance

- **Background – Voldemort architecture**
- Custom Voldemort storage engine
    - Minimal impact on live system
    - Fast rollback
    - Fast lookups
    - Easy rebalancing
- Performance

- Amazon's Dynamo clone
  - Distributed key-value store
- Pluggable architecture
- Initially written for read-write storage engines
  - MySQL, BDB

Node 0

Node 1

Store 1

Store 1

Store 2

Store 2

- Multiple peers
- Multiple stores ( ~ tables )
- Different replication factor per store

# Project Voldemort

- Hash ring per store
- Ring split into partitions

**Node 0**

Store 1
Partition
0, 2, 4

Store 2
Partition
0, 2, 4

**Node 1**

Store 1
Partition
1, 3, 5

Store 2
Partition
1, 3, 5

$2^{32}-1$  0

Partition 5
(Node 1)

Partition 0
(Node 0)

Partition 4
(Node 0)

Partition 1
(Node 1)

Partition 3
(Node 1)

Partition 2
(Node 0)

- Latency degradation during load

Processing system (Hadoop)

Staging

LOAD DATA

Store 1

Staging

LOAD DATA

Store 1

T1    T2

View

T1    T2

View

Live requests

- Maintaining extra cluster

Processing system (Hadoop)

Staging

LOAD DATA

Store 1

Staging

LOAD DATA

Store 1

T1  T2

View

T1  T2

View

Live requests

- Maintaining extra cluster

# Existing approaches – Bulk load solution 3



Processing system (Hadoop)

M • • • M

R          R

HDFS

T1  T2      T1  T2

View        View

Live requests

- Multiple copy operations

- Background – Voldemort architecture
- Custom Voldemort storage engine
  - Minimal impact on live system
  - Fast rollback
  - Fast lookups
  - Easy rebalancing
- Performance

- Background – Voldemort architecture
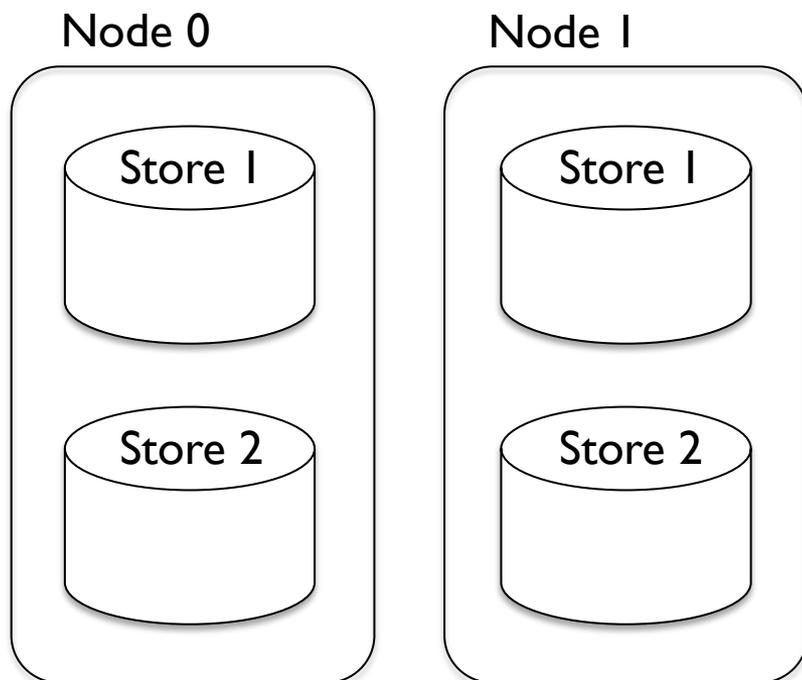- Custom Voldemort storage engine
  - Minimal impact on live system
  - Fast rollback
  - Fast lookups
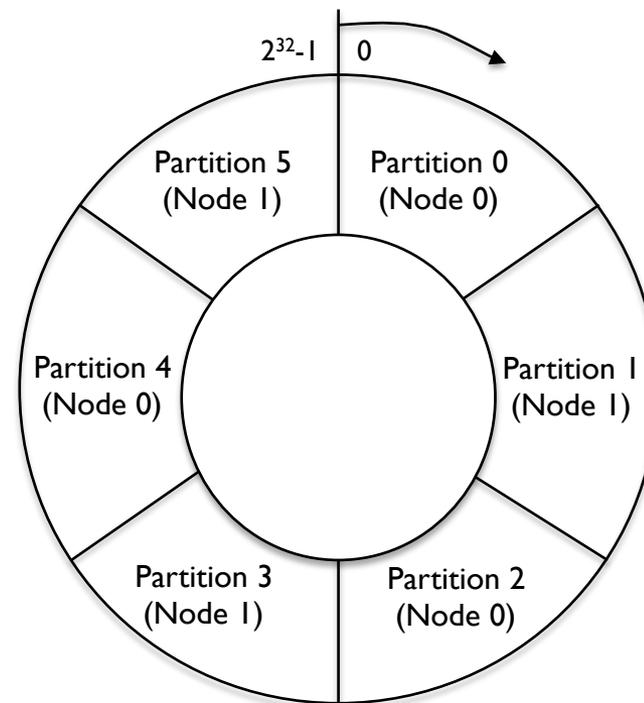  - Easy rebalancing
- Performance

HDFS

Hadoop flow

Construct store

- "Construct Store" step
  - Single MapReduce job
  - Map
    - Input - Output of algorithm
    - Output – Emit replication factor number of times
  - Partitioner
    - Redirect to appropriate reducer
  - Reducer
    - Output to Voldemort node based folders

# Bulk load extensions – Minimal impact on live system

- What is output in the reducer phase?
    - Store $\xrightarrow{\text{split into}}$ Partitions $\xrightarrow{\text{split into}}$ Chunk sets
    - One reducer = one chunk set
    - Chunk set = Index + data file

| Upper 8 bytes of MD5 of key | Offset into data file |
|---|---|
| | |

Sorted by top 8 bytes ↓

**Index file**

| Number of collided tuples | Key size | Value size | Key | Value | • • • |
|---|---|---|---|---|---|
| | | | | | • • • |

Tuple | Other collided tuples

**Data file**

- Background – Voldemort architecture
- Custom Voldemort storage engine
  - Minimal impact on live system
  - Fast rollback
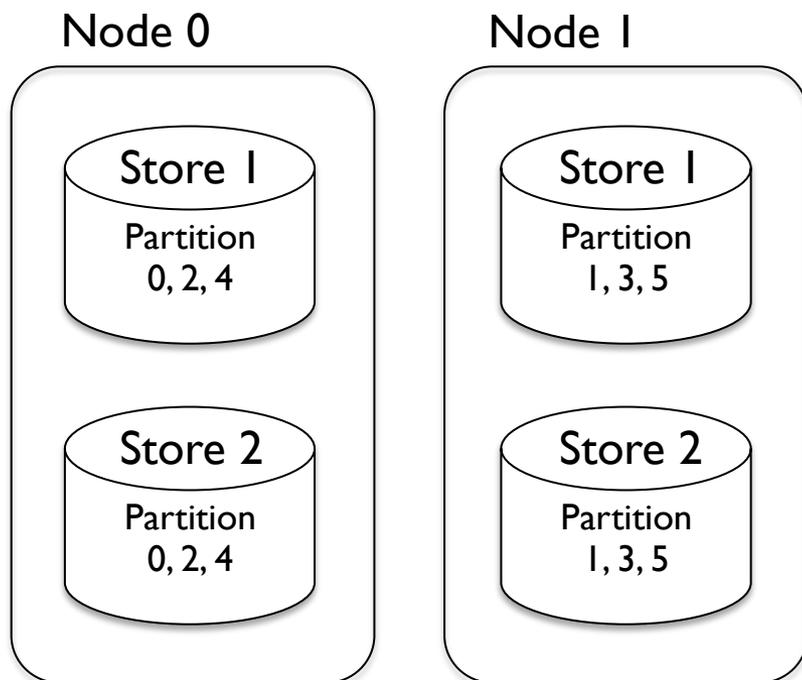  - Fast lookups
  - Easy rebalancing
- Performance

- **Construction**
  - **MapReduce job**
- Fetch
  - Pull chunk sets in parallel
  - Store into new version folder
- Swap
  - Close latest version's index files
  - Change latest version link
  - Memory map new version's index files
- Rollback
  - Close latest version's index file
  - Change latest version link
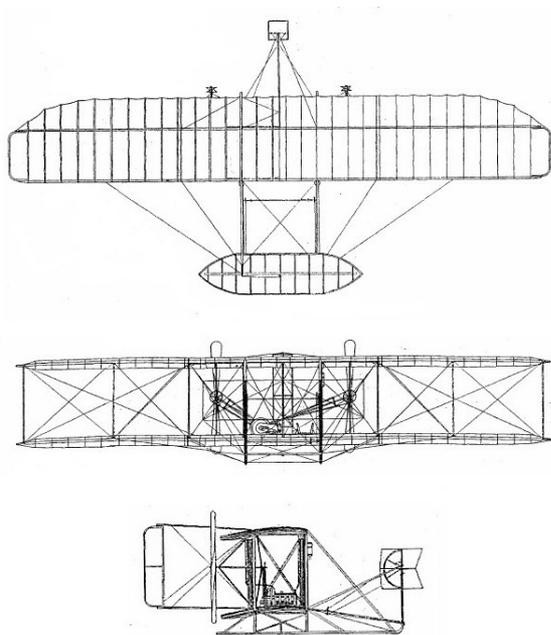  - Memory map old version's index file

Hadoop

HDFS

Trigger
Construction

Driver /
Scheduler

Voldemort
cluster

- Construction
  - MapReduce job
- **Fetch**
  - **Pull chunk sets in parallel**
  - **Store into new version folder**
- Swap
  - Close latest version's index files
  - Change latest version link
  - Memory map new version's index files
- Rollback
  - Close latest version's index file
  - Change latest version link
  - Memory map old version's index file

Hadoop

HDFS

Driver /
Scheduler

Parallel
pull

Trigger
Fetch

Voldemort
cluster

# Bulk load extensions - Rollback

Voldemort node

store-1

Chunk sets
/version-(i)

Chunk sets
/version-(i+1)

*latest* → version-(i+1)

/store-1

- Construction
  - MapReduce job
- **Fetch**
  - **Pull chunk sets in parallel**
  - **Store into new version folder**
- Swap
  - Close latest version's index files
  - Change latest version link
  - Memory map new version's index files
- Rollback
  - Close latest version's index file
  - Change latest version link
  - Memory map old version's index file

Hadoop

HDFS

Driver / Scheduler

Parallel pull

Trigger Fetch

Voldemort cluster

- Construction
  - MapReduce job
- Fetch
  - Pull chunk sets in parallel
  - Store into new version folder
- **Swap**
  - **Close latest version's index files**
  - **Change latest version link**
  - **Memory map new version's index files**
- Rollback
  - Close latest version's index file
  - Change latest version link
  - Memory map old version's index file

| Hadoop |
| --- |
| HDFS |

Driver / Scheduler

Trigger Swap

Voldemort cluster

Voldemort node

store-1

Chunk sets

/version-(i+1)

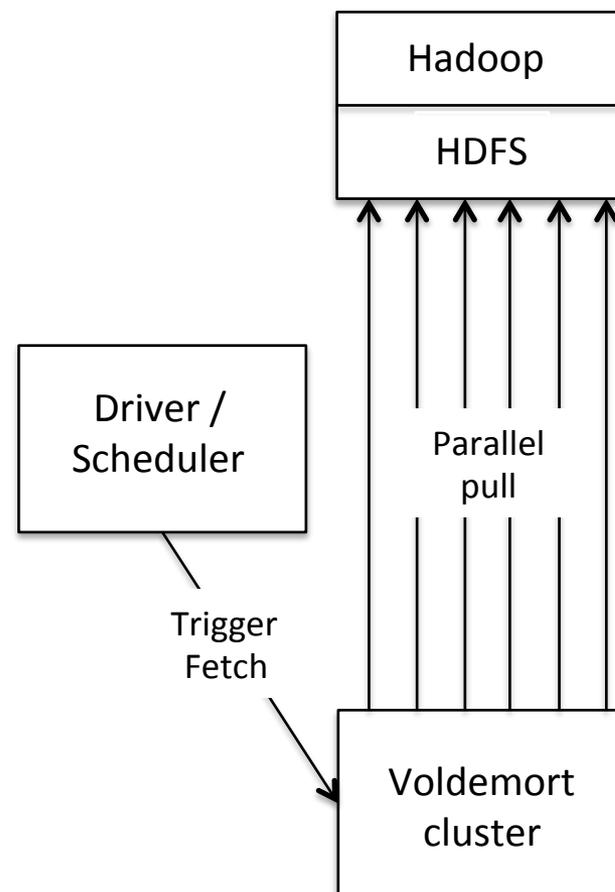Chunk sets

/version-(i+2)

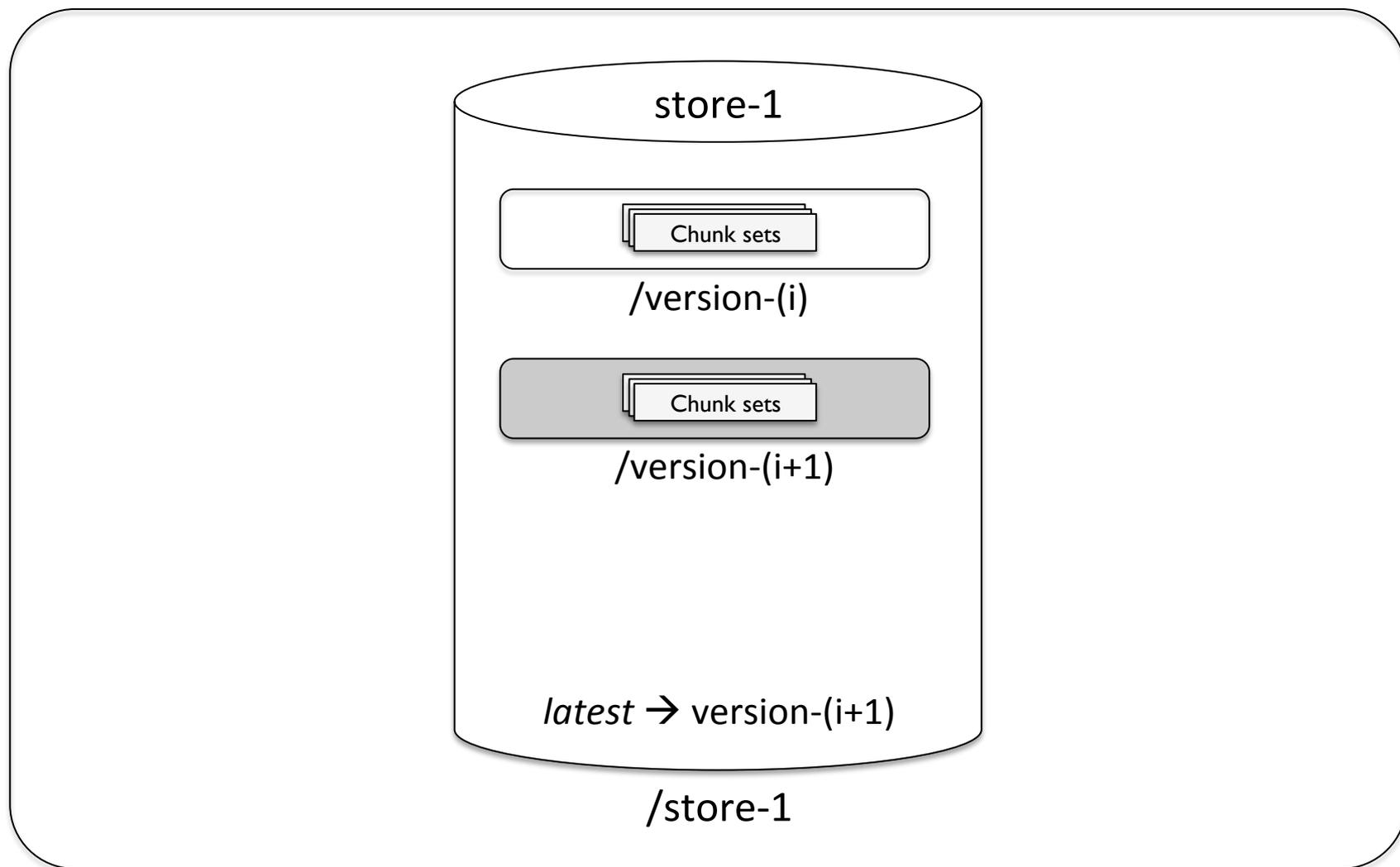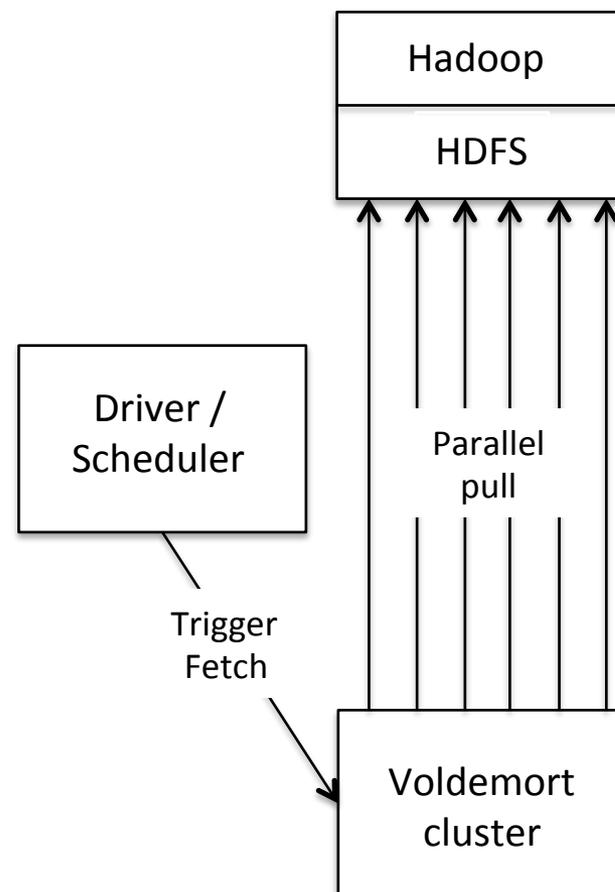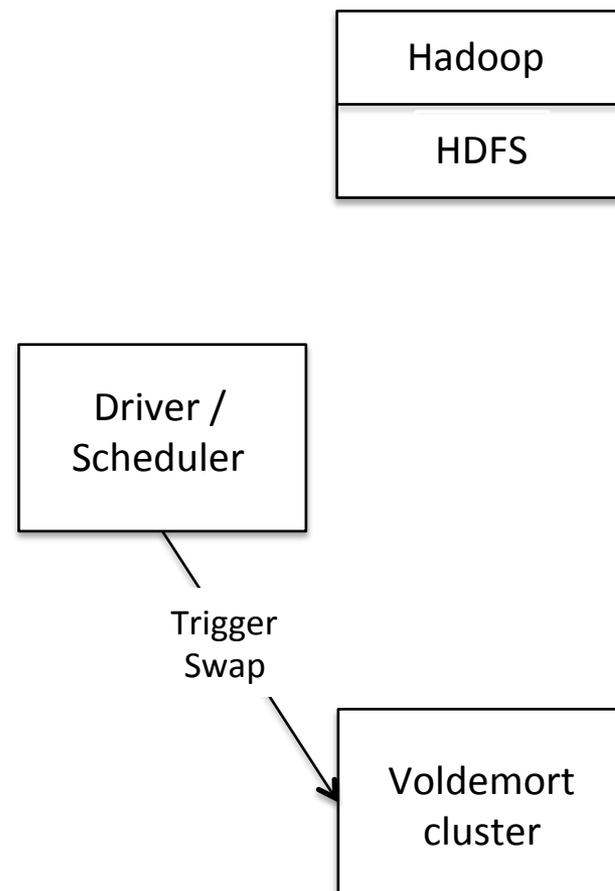*latest* → version-(i+2)

/store-1

# Bulk load extensions – Full pipeline

- Construction
  - MapReduce job
- Fetch
  - Pull chunk sets in parallel
  - Store into new version folder
- Swap
  - Close latest version's index files
  - Change latest version link
  - Memory map new version's index files
- **Rollback**
  - **Close latest version's index file**
  - **Change latest version link**
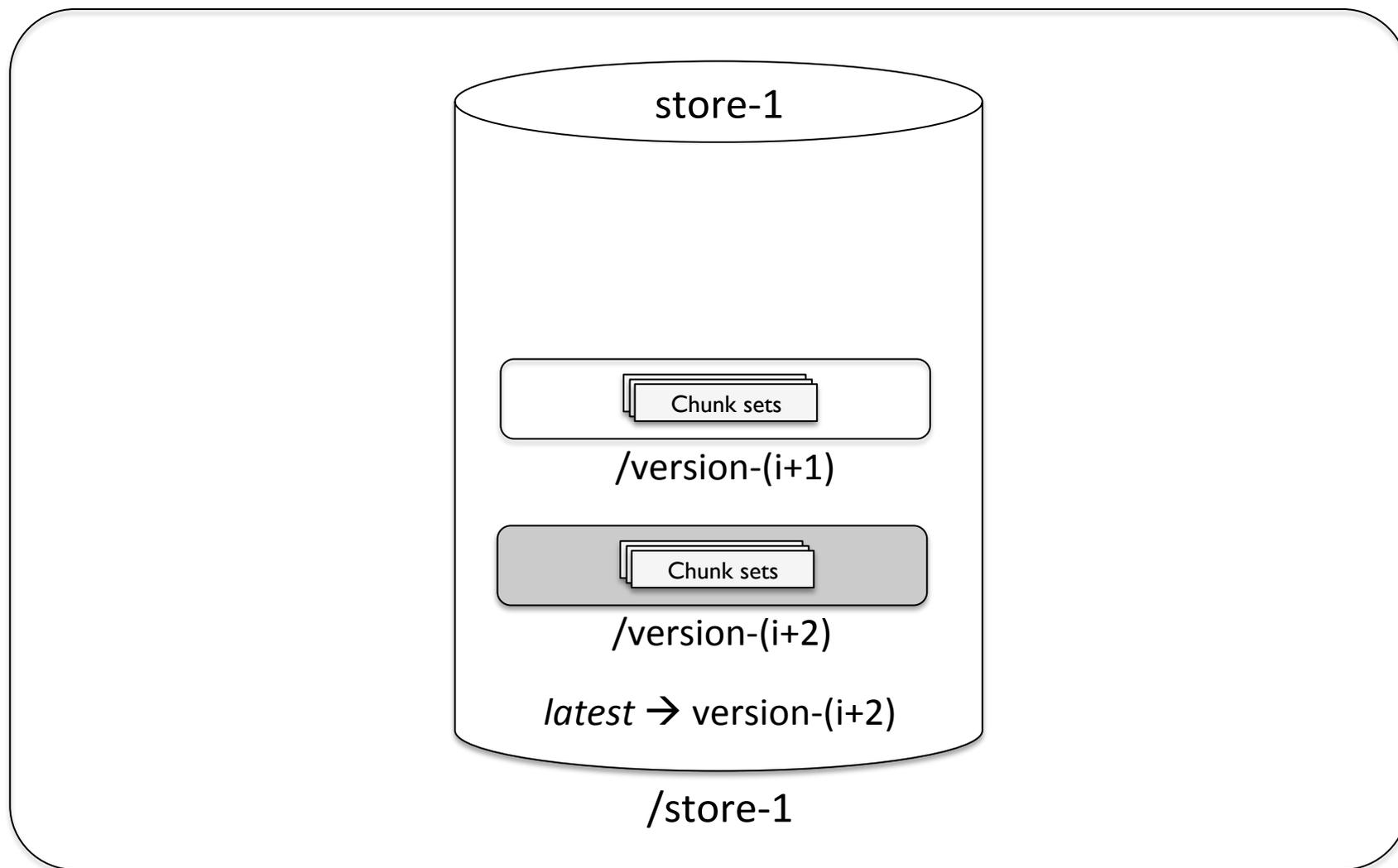  - **Memory map old version's index file**

Hadoop

HDFS

Driver / Scheduler

Trigger Rollback

Voldemort cluster

# Bulk load extensions - Rollback

Voldemort node

store-1

Chunk sets

/version-(i+1)

Chunk sets

/version-(i+2)

*latest* → version-(i+1)

/store-1

- Background – Voldemort architecture
- Custom Voldemort storage engine
  - Minimal impact on live system
  - Fast rollback
  - Fast lookups
  - Easy rebalancing
- Performance

# Bulk load extensions – Lookup

- Find partition and chunk set to read
- Binary search in index file of chunk set
- Jump to offset in data file
- Go through all collided tuples



| Upper 8 bytes of MD5 of key | Offset into data file |
| --- | --- |

Sorted by top 8 bytes

Index file

| Number of collided tuples | Key size | Value size | Key | Value | • • • |
| --- | --- | --- | --- | --- | --- |

Tuple

Other collided tuples

Data file

- Background – Voldemort architecture
- Custom Voldemort storage engine
  - Minimal impact on live system
  - Fast rollback
  - Fast lookups
  - Easy rebalancing
- Performance

- Adding new nodes with no downtime
- Change ownership of partitions to new nodes
  - Simple move of corresponding chunk sets + swap
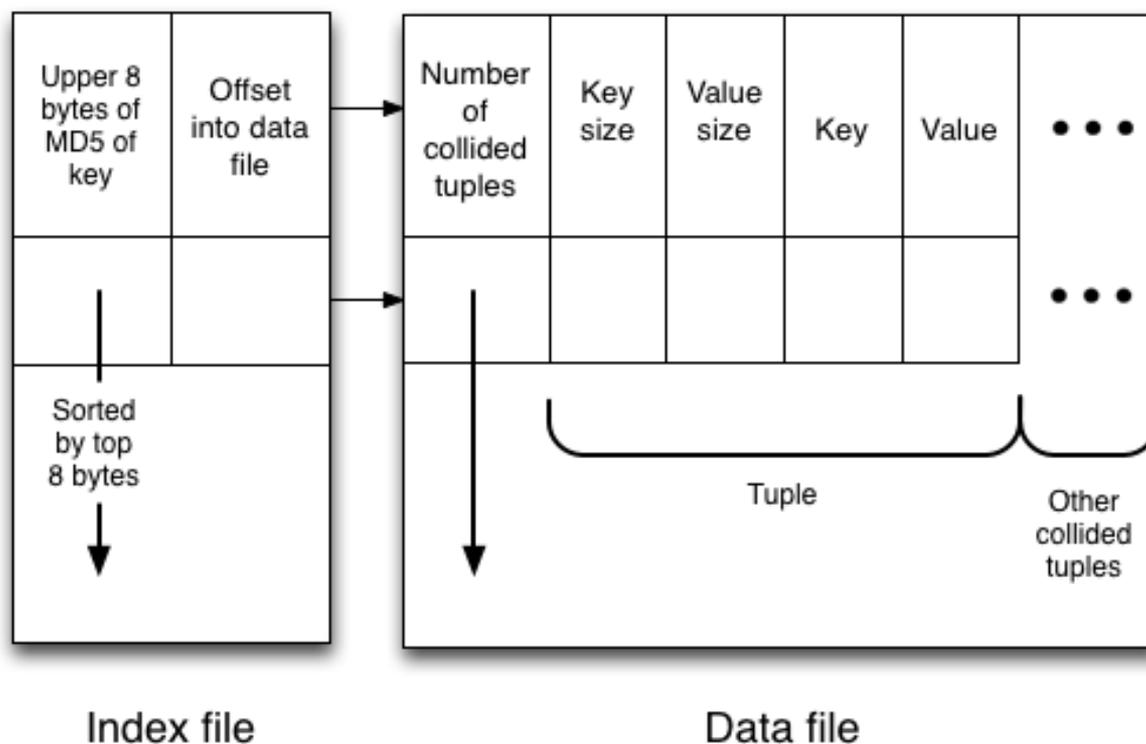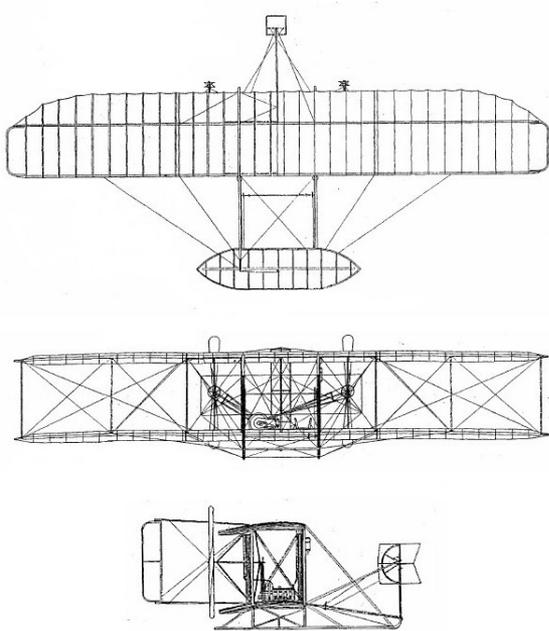
- Background – Voldemort architecture
- Custom Voldemort storage engine
  - Minimal impact on live system
  - Fast rollback
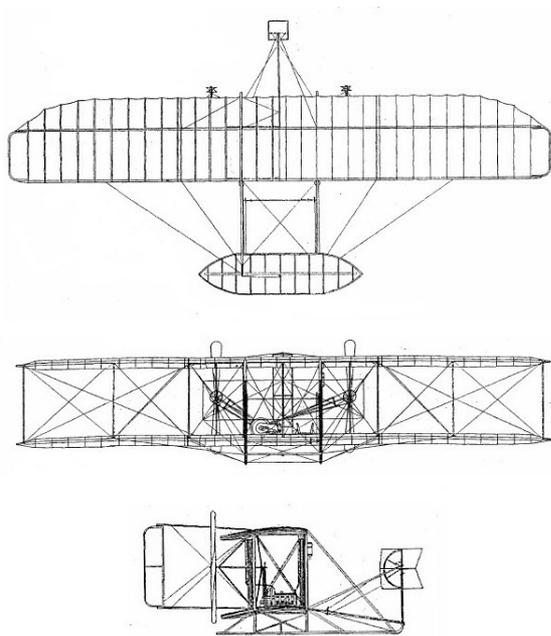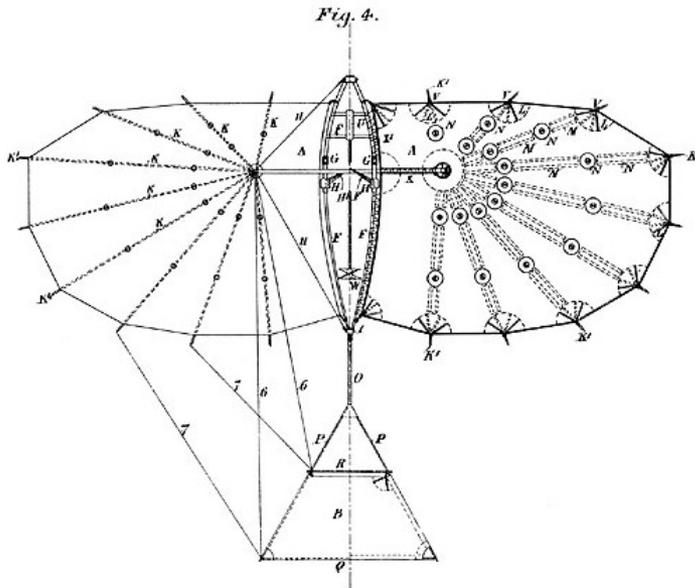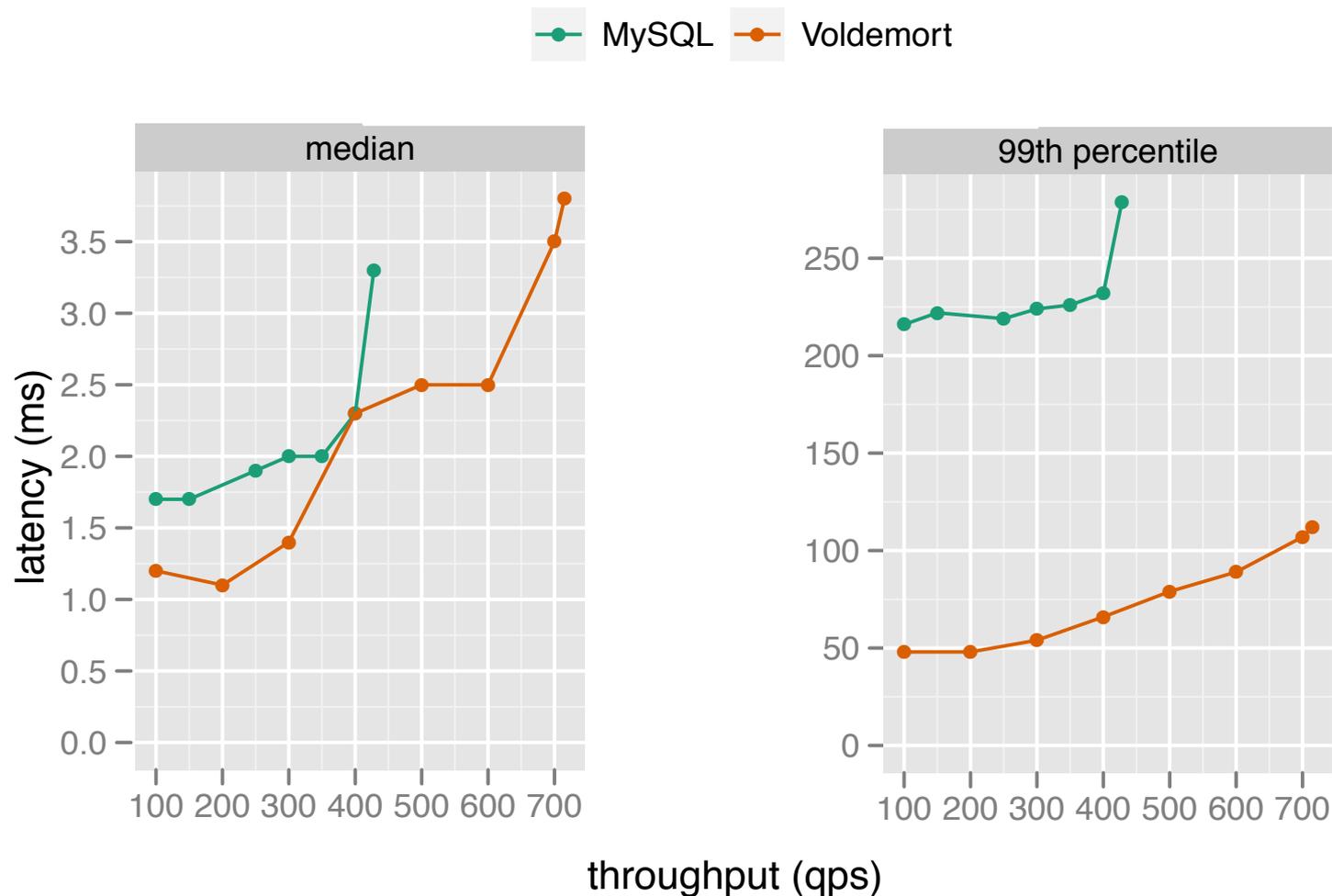  - Fast lookups
  - Easy rebalancing
- Performance

Fig. 4.

- Single node latency
- Multi-node latency
- Production

# Performance – Single node latency



100 GB data, 24 GB RAM

32 nodes client side latency

People You May Know – 25 nodes client side latency

- Shared nothing databases
    - Yahoo's PNUTS [Silberstein08]
        - Bulk load into range partitioned tables
        - Requires some compaction on serving system
    - HBase [Konstantinou11], Cassandra [Lebresne11]
        - Offline tablet construction
        - Expensive rollback
- Search systems
    - Build index offline and pull
        - MapReduce[Dean04] use-case

- LinkedIn
  - Serving ~120 stores in production for past 2 years
  - Fetching ~ 4 TB of data every day
  - 76 stores swapped every day
- Open-source
  - http://project-voldemort.com

# Bibliography

**Images**
- Slide 6 - www.aviastar.org/air/inter/inter_concorde.php
- Slide 15 - www.wright-brothers.org/Information_Desk/Just_the_Facts/Airplanes/Flyer_II.htm
- Slide 17 - www.youtube.com/watch?v=oz-7wJJ9HZ0
- Slide 48 - www.wright-brothers.org/History_Wing/History_of_the_Airplane/Century_Before/Road_to_Kitty_Hawk/Road_to_Kitty_Hawk.htm

**Papers**

[1]     Giuseppe DeCandia, Deniz Hastorun, Madan Jampani, Gunavardhan Kakulapati, Avinash Lakshman, Alex Pilchin, Swaminathan Sivasubramanian, Peter Vosshall, and Werner Vogels. 2007. Dynamo: Amazon's Highly Available Key-value Store. In Proceedings of 21st ACM SIGOPS symposium on Operating systems principles (SOSP '07)

[2]     Jeffrey Dean and Sanjay Ghemawat. MapReduce: Simplified Data Processing on Large Clusters. In Proceedings of the 6th Conference on Symposium on Operating Systems Design & Implementation (OSDI '04)

[3]     Adam Silberstein, Brian F. Cooper, Utkarsh Srivastava, Erik Vee, Ramana Yerneni, and Raghu Ramakrishnan. 2008. Efficient bulk insertion into a distributed ordered table. In Proceedings of the 2008 ACM SIGMOD international conference on Management of data (SIGMOD '08)

[4]     Ioannis Konstantinou, Evangelos Angelou, Dimitrios Tsoumakos, and Nectarios Koziris. Distributed Indexing of Web Scale Datasets for the Cloud. In Proceedings of the 2010 Workshop on Massive Data Analytics on the Cloud (MDAC '10)

[5]     Sasha Pachev. Understanding MySQL Internals. O'Reilly Media, 2007.

[6]     Tom White. Hadoop: The Definite Guide. O'Reilly Media, 2010

[7]     Sylvain Lebresne. Using the Cassandra Bulk Loader. http://www.datastax.com/dev/blog/bulk-loading