

# Pinot

Real-time Analytics at LinkedIn

*Greg Brandt (2015-09-11)*

## Agenda

- Pinot at LinkedIn
- Design and Architecture
- Segment Layout
- StarTree for Fast Aggregation
- ThirdEye for Business Monitoring

- Who's Viewed My Profile  
– (member)



- XLNT
  - A/B testing
  - (internal)

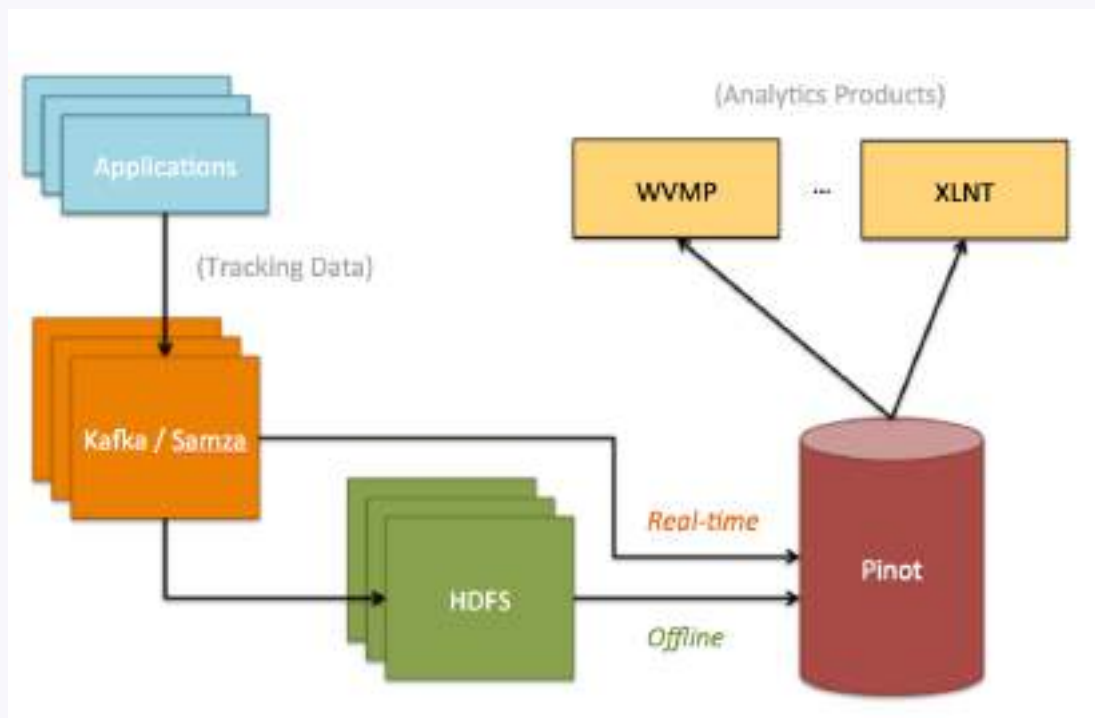


## Use Cases at LinkedIn

- Ad campaigns – (customers)



# Data Analytics Ecosystem at LinkedIn



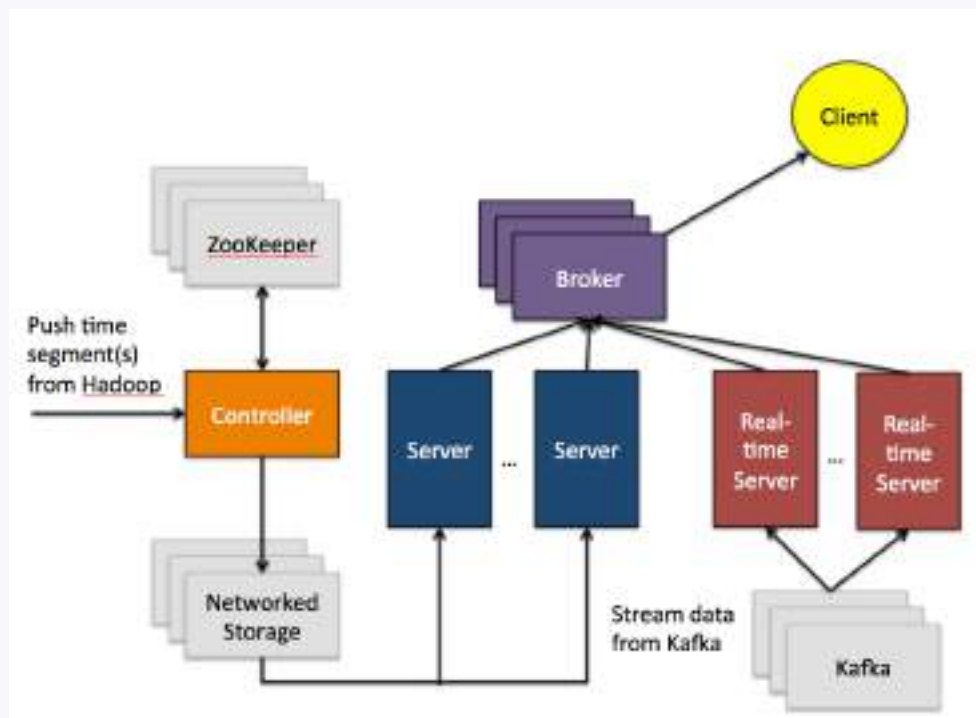
### Goals

- SQL-like interface
- Columnar storage and indexing
- Real-time data load

### Non-goals

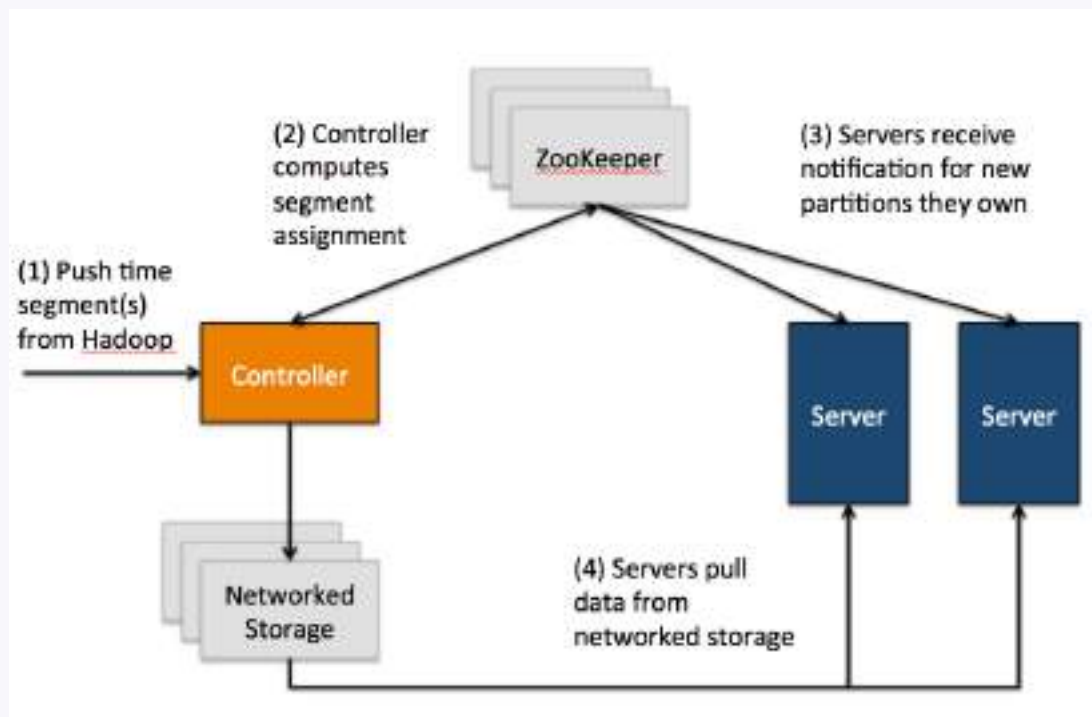
- JOIN, UDF (unpredictable performance)
- Not source of truth
- Mutability
  - Immutability allows trivial data distribution

# Pinot Architecture



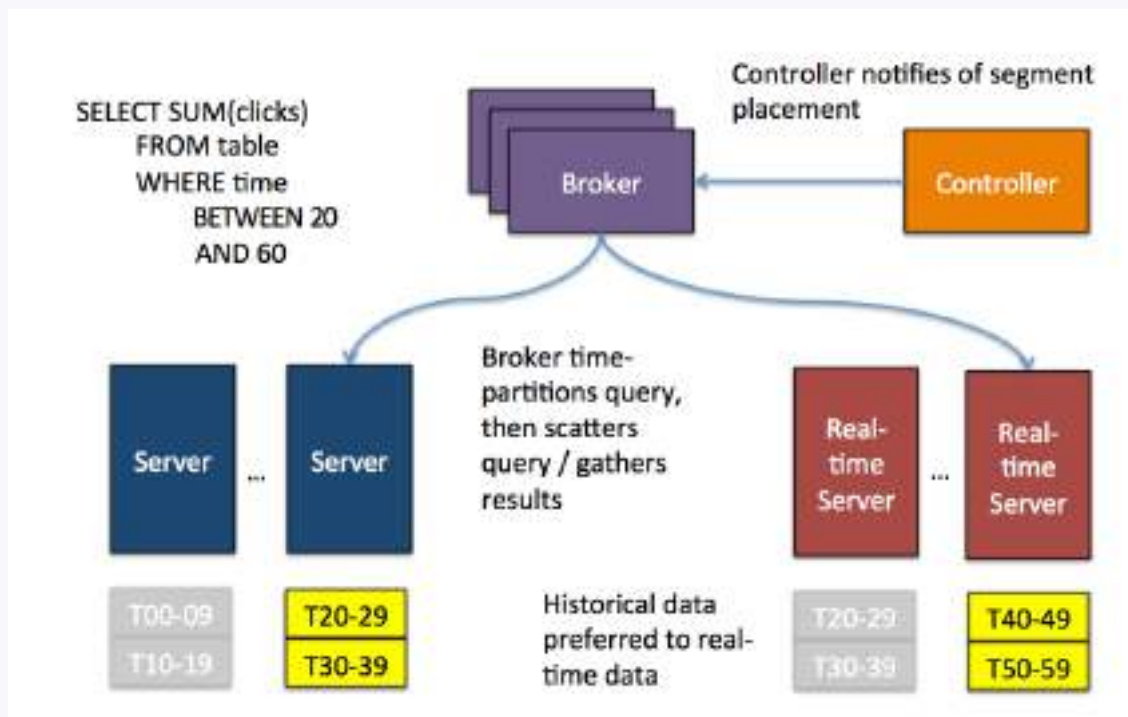


## Segment Assignment



- A framework for building distributed, stateful applications
- Developed at LinkedIn for Espresso, (distributed OLTP database)
- Leveraged in Pinot for cluster management (Controller uses Helix libraries)
- <http://helix.apache.org>

## Query Execution: Distributed



# 1. Parser

- SELECT COUNT(\*) ... => Abstract Syntax Tree

# 2. Logical Planner

- Operators like COUNT, WHERE, etc.

# 3. Physical Planner

- Consider data layout (sort key, dictionary encoding, compression, etc.)

## Layout: Columnar Storage

Doc Id	Model	Year	Price	tags
1	Toyota	1997	7500	moon-roof, reliable
2	Toyota	2001	7500	mp3, compact
3	Mazda	2005	1500	Manual, compact
4	Honda	2010	8000	leather, cool
5	Honda	2011	9000	automatic
-	-	-	-	-
-	-	-	-	-
-	-	-	-	-
-	-	-	-	-
10000	bmw	-	25000	moon-roof, reliable

RAW DATA



INDEX  
GENERATION

Doc Id	Model	Year	Price	tags
1	Toyota	1997	7500	moon-roof, reliable
2	Toyota	2001	7500	mp3, compact
3	Mazda	2005	1500	Manual, compact
4	Honda	2010	8000	leather, cool
5	Honda	2011	9000	automatic
-	-	-	-	-
-	-	-	-	-
-	-	-	-	-
-	-	-	-	-
10000	bmw	-	25000	moon-roof, reliable

COLUMNAR DATA

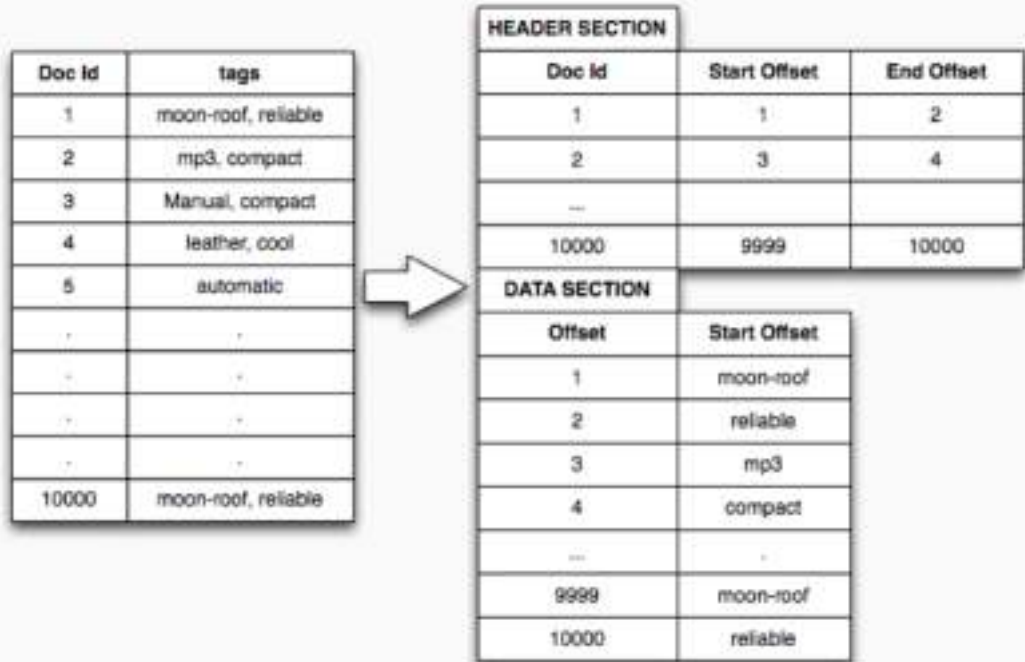
## Layout: Sorted Forward Index

Doc Id	Model
1	Acura
...	Acura
100	Acura
101	Bmw
..	Bmw
500	Bmw
.	.
.	.
99,999,900	Zagato
100,000,000	Zagato



Model	Start DocId	End DocId
Acura	1	100
Bmw	101	500
...		
..	.	.
...	.	.
...		
Zagato	99,999,900	100,000,000

- Indexes
  - Bitmap,
  - Roaring Bitmap,
  - Inverted
- Compression
  - Dictionary encoding,
  - P4Delta
- Multi-value columns



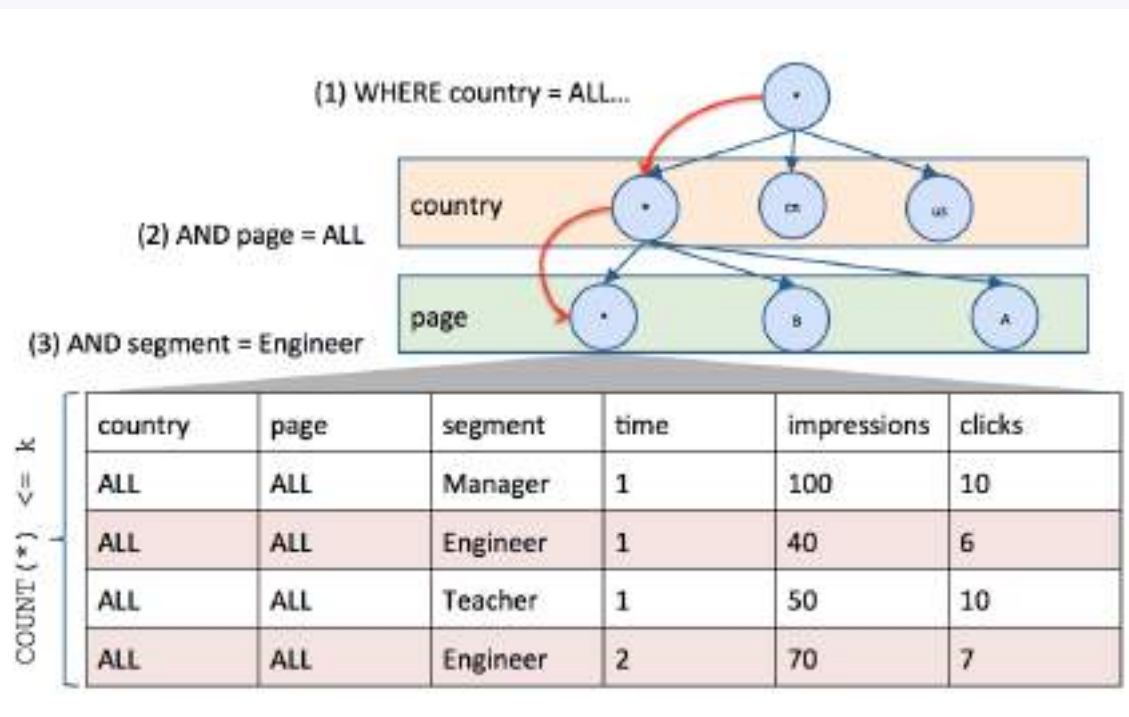
Layout: StarTree Index

Based on “Iceberg-Cubes” (Bottom-Up)

Invariant: scan  $\leq k$  records to answer any aggregation query



## StarTree Index Structure



- Consider a data set with two metrics (**clicks**, **impressions**), and four dimensions (**country**, **page**, **user segment**, **time**)
- ~3B raw events per day
- 196 countries, 100 pages, 25 user segments
- Store documents at 1 minute granularity

Let' s calculate aggregate time series for one week...

## StarTree Case Study

```
SELECT SUM(clicks),SUM(impressions)
FROM table
WHERE time BETWEEN
      2015-08-01 AND
      2015-08-08
GROUP BY (country, time)
```

*StarTree space complexity largely impacted by data skew; **1.5 to 2 times raw data size** in practice*

*Assuming scan speed of ~40M documents per second*

Method	# Docs Scanned	Time Complexity	Space Complexity
Ad-hoc Pig / MR	21B	<b>8.75 minutes</b>	Constant
Time roll-up	705M	<b>2.04 minutes</b>	Constant
Full cubing	Constant	Constant	<b>128.66 GB / week</b>
StarTree (k = 100)	28M	<b>0.7 seconds</b>	O(leaves * k)

1. Take highest cardinality dimension,  $D_h$
2. Compute aggregates after removing  $D_h$
3. Append aggregates to segment, sort on  $D_h$
4. For each value of  $D_h$ , *if there are greater than  $k$  documents in sub-table*, repeat on sub-table
  - Including the new aggregate value, “\*”

## Building StarTree in Pinot

```
>> tree testTable_0_127_1/
testTable_0_127_1/
├── D0.dict
├── D0.sv.unsorted.fwd
├── M0.dict
├── M0.sv.unsorted.fwd
├── ...
├── creation.meta
├── daysSinceEpoch.dict
├── daysSinceEpoch.sv.unsorted.fwd
├── metadata.properties
├── startree
│   ├── D0.dict
│   ├── D0.sv.unsorted.fwd
│   ├── M0.dict
│   ├── M0.sv.unsorted.fwd
│   ├── ...
│   ├── daysSinceEpoch.dict
│   ├── daysSinceEpoch.sv.unsorted.fwd
│   ├── metadata.properties
│   └── startree.ser
```

The original Pinot segment, now partially sorted on dimension combination

The aggregate Pinot segment, containing partial aggregates (i.e. any StarTree prefix involving a “\*” node)

## StarTree speeds up Pinot queries two ways...

1. Aggregation queries can use partial aggregates, bound scan time by  $k$
2. Raw segment is partially sorted by StarTree path dimension values, improves filter queries

- Leverage StarTree to to monitor highly-dimensional, aggregate time-series business metrics in real-time
- Interactive visualization tools
  - Drill-down, guided root-cause analysis
- Automated anomaly detection algorithms
  - Data-driven OLAP cube exploration
  - Detect spikes, step-level shifts

## Visualization: Heat Map

- Used in financial sector
- E.g. Showing y/3y on 2008-08-08, for S&P 500
- Click on cell to fix `ticker = 'aig'` in `WHERE` clause, drill down

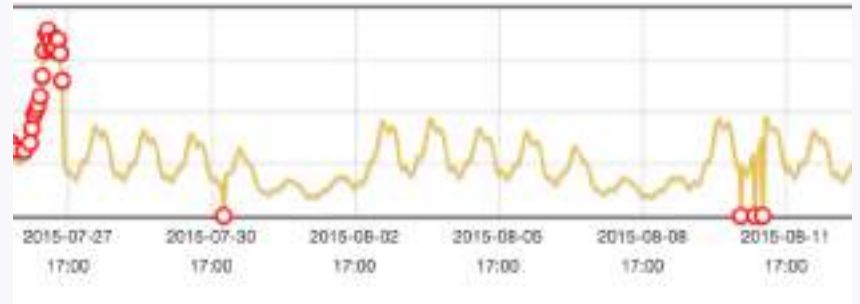
Financial services  
like AIG and  
CitiGroup





- Explore OLAP cube top-down (start at overall aggregate metric)
- Use heuristic (like total volume) to prune search, determine dimension exploration order
- Compute GROUP BY (next\_dimension, time) at each level, and run anomaly detection algorithm on each time series

- Model time series structural components: seasonality + trend + noise
- Best for detecting spikes/dips
- Applied in Ads metric monitoring



- Hypothesis testing to determine whether the pattern within time series interval differs from the rest
- Best for detecting step-level changes

Anomalies in Home Page download time, in India

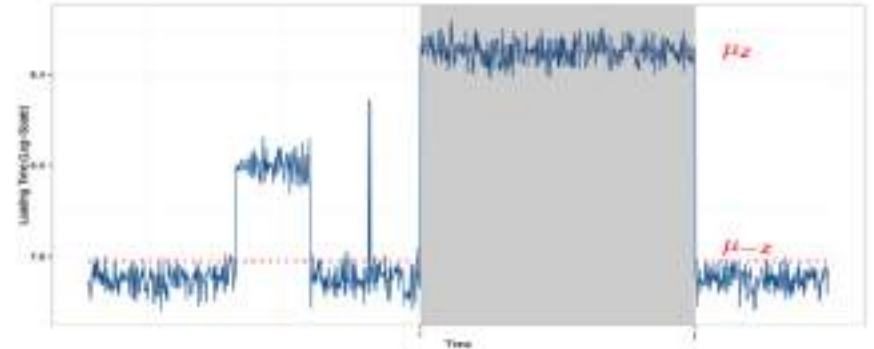


Figure 3 Scan windows and hypothesis testing.

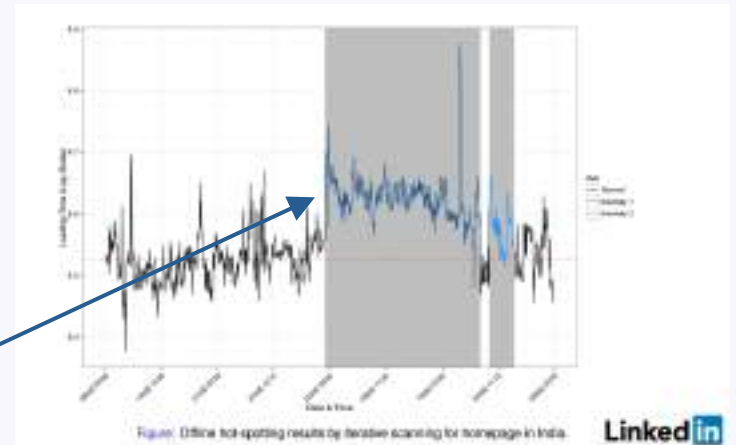


Figure: Offline for-spotting results by iterative scanning for homepage in India

Questions?

[gbrandt@linkedin.com](mailto:gbrandt@linkedin.com)



TALKINGDATA GLOBAL MOBILE  
**DATA SUMMIT**