

Performance study of Spindle, a web analytics
query engine implemented in Spark
CloudCom 2014

Brandon Amos* and David Tompkins, **Adobe Research**

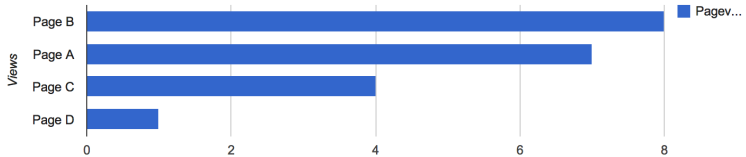
*Adobe summer intern, Ph.D. Student at Carnegie Mellon University.

December 19, 2014



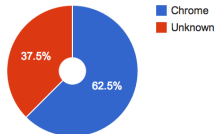
Motivation

- Adobe Marketing Cloud offers web analytics.



Page B

Total Pageviews: 8



Chrome	5
Unknown	3

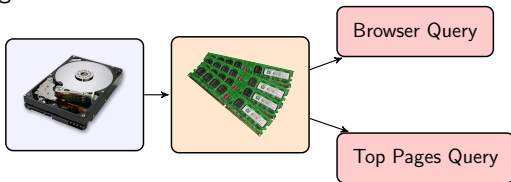
Motivation

- ▶ Adobe Marketing Cloud offers web analytics.
- ▶ Terabytes of data, thousands of servers.
- ▶ Trending general-purpose distributed data processing engines.
 - ▶ Apache Spark
 - ▶ Cloudera Impala
 - ▶ Google Dremel
- ▶ **Spindle** is an early investigation of the feasibility of Apache Spark for web analytics



Motivation

- ▶ Ideal Spark features
 - ▶ In-memory caching
 - ▶ Lineage



- ▶ **Problem:** Current performance studies do not show Spark's performance for interactive web analytics application.

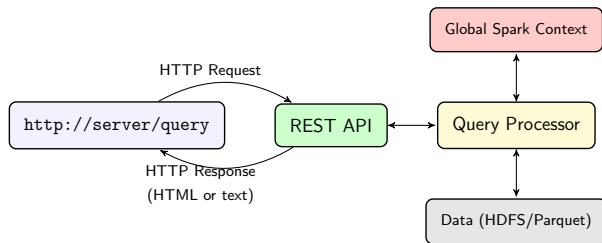


Spindle Architecture

Overview.

What is Spindle?

Parameters:
Date Range
Spark Tuning Options



Shorthand	Name
Q0	Pageviews
Q1	Revenue
Q2	RevenueFromTopReferringDomains
Q3	RevenueFromTopReferringDomainsFirstVisitGoogle
Q4	TopPages
Q5	TopPagesByBrowser
Q6	TopPagesByPreviousTopPages
Q7	TopReferringDomains

	Q0	Q1	Q2	Q3	Q4	Q5	Q6	Q7
post_pagename	x				x	x	x	
user_agent						x		
visit_referrer			x	x				
post_visid_high			x	x			x	x
post_visid_low			x	x			x	x
visit_num			x	x			x	x
visit_referrer								x
hit_time_gmt							x	
post_purchaseid		x	x	x				
post_product_list		x	x	x				
first_hit_referrer				x				



Spindle Architecture

Queries.

- Demo: <http://adobe-research.github.io/spindle/>



Spindle Architecture

Ad hoc queries.

- Spark SQL processes relational queries.

```
Press <tab> to see a list of available commands.
```

```
> sql select count(*) from all_data
```

```
[20]
```

```
> sql select post_pagename, hit_time_gmt from data_2014_08_16 order by hit_time_gmt
```

```
[Page D,1408187379]
```

```
[Page A,1408187380]
```

```
[Page B,1408187380]
```

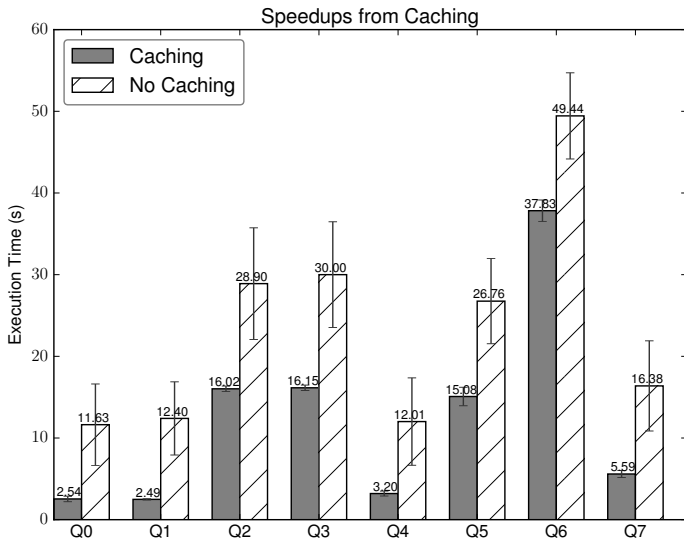


Empirical Results

Caching.

- ▶ Six cluster nodes, Spark and HDFS on each.
- ▶ 13.1GB of data.
- ▶ **Question:** How does caching in-memory improve performance?



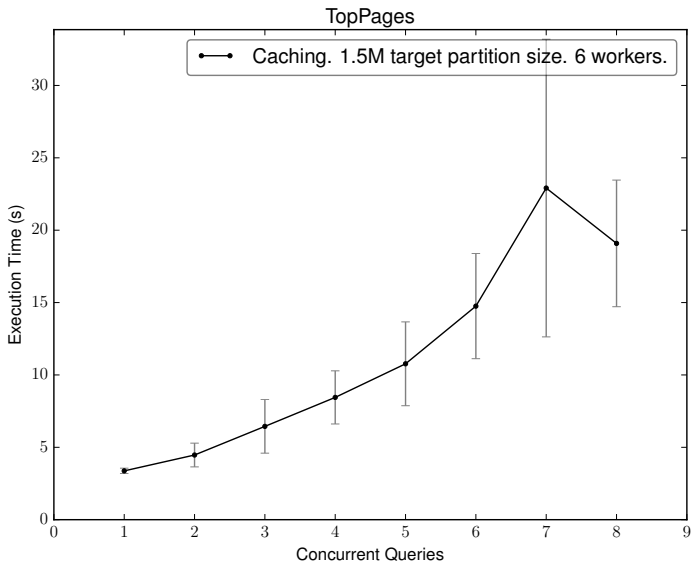


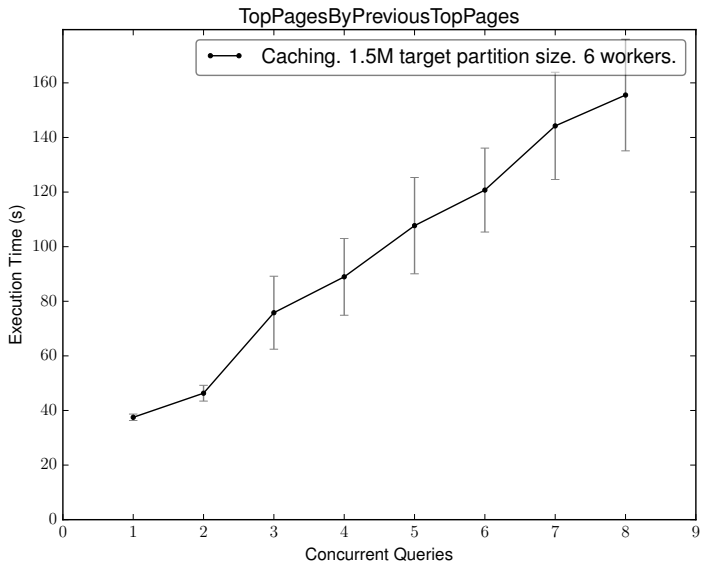
Empirical Results

Benchmarking concurrent queries.

- ▶ How much will Spindle's performance degrade if multiple users are utilizing it at the same time?
- ▶ Concurrently call a single query on the same data.
- ▶ Average the execution time.







Empirical Results

Remaining.

- ▶ Scaling Spark and HDFS workers.
- ▶ Intermediate data partitioning.



Conclusions

- ▶ Spark is a good candidate for real-time analytics processing.
- ▶ **Spindle** is an open-source prototype analytics processing engine.
 - ▶ Sample set of web analytics queries.
 - ▶ REST-based interface to tune parameters.
- ▶ Spindle's future work is on preprocessing.

Spindle Project		http://github.com/adobe-research/spindle
Brandon Amos		http://github.com/bamos
David Tompkins		http://github.com/DavidTompkins

