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

CloudCom 2014

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

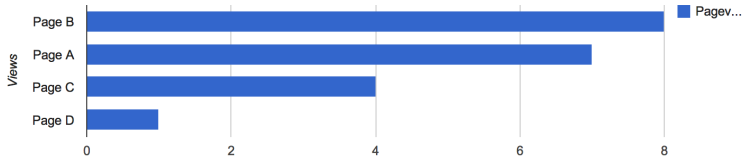[*]Adobe summer intern, Ph.D. Student at Carnegie Mellon University.

December 19, 2014

Motivation

- ▶ Adobe Marketing Cloud offers web analytics.
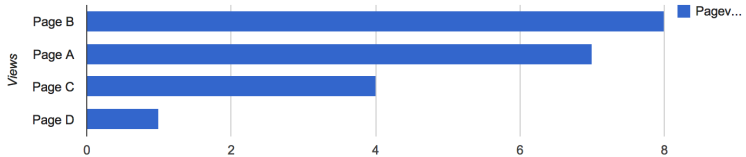
# Motivation

► Adobe Marketing Cloud offers web analytics.

# 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.

# Motivation

- Adobe Marketing Cloud offers web analytics.
- Terabytes of data, thousands of servers.
- Trending general-purpose distributed data processing engines.

# Motivation

- Adobe Marketing Cloud offers web analytics.
- Terabytes of data, thousands of servers.
- Trending general-purpose distributed data processing engines.
  - Apache Spark

# Motivation

- Adobe Marketing Cloud offers web analytics.
- Terabytes of data, thousands of servers.
- Trending general-purpose distributed data processing engines.
  - Apache Spark
  - Cloudera Impala

# 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

# 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

# Motivation

- Ideal Spark features
  - In-memory caching

# Motivation

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

# Motivation

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

# Motivation

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

# Motivation

- Ideal Spark features
  - In-memory caching
  - Lineage

# Motivation

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

## Motivation

- Ideal Spark features
  - In-memory caching
  - Lineage



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

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

**Overview.**
Queries.
Ad hoc queries.

# Spindle Architecture
Overview.

What is Spindle?

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

**Overview.**
Queries.
Ad hoc queries.

# Spindle Architecture

Overview.

What is Spindle?

```
http://server/query
```

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

**Overview.**
Queries.
Ad hoc queries.

# Spindle Architecture
Overview.

### What is Spindle?

**Parameters:**
Date Range
Spark Tuning Options

`http://server/query`

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

**Overview.**
Queries.
Ad hoc queries.

# Spindle Architecture
Overview.

### What is Spindle?

**Parameters:**
Date Range
Spark Tuning Options

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

**Overview.**
Queries.
Ad hoc queries.

# Spindle Architecture
Overview.

### What is Spindle?

**Parameters:**
Date Range
Spark Tuning Options

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

**Overview.**
Queries.
Ad hoc queries.

# Spindle Architecture
Overview.

### What is Spindle?

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

**Overview.**
Queries.
Ad hoc queries.

# Spindle Architecture
Overview.

### What is Spindle?



**Parameters:**
Date Range
Spark Tuning Options

HTTP Request

`http://server/query`

REST API

HTTP Response
(HTML or text)

Global Spark Context

Query Processor

Data (HDFS/Parquet)

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
**Queries.**
Ad hoc queries.

| Shorthand | Name |
|---|---|
| Q0 | Pageviews |

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
**Queries.**
Ad hoc queries.

| Shorthand | Name |
| --- | --- |
| Q0 | Pageviews |
| Q1 | Revenue |

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
**Queries.**
Ad hoc queries.

| Shorthand | Name |
|-----------|------|
| Q0 | Pageviews |
| Q1 | Revenue |
| Q2 | RevenueFromTopReferringDomains |

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
**Queries.**
Ad hoc queries.

| Shorthand | Name |
|-----------|------|
| Q0 | Pageviews |
| Q1 | Revenue |
| Q2 | RevenueFromTopReferringDomains |
| Q3 | RevenueFromTopReferringDomainsFirstVisitGoogle |

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
**Queries.**
Ad hoc queries.

| Shorthand | Name |
| --- | --- |
| Q0 | Pageviews |
| Q1 | Revenue |
| Q2 | RevenueFromTopReferringDomains |
| Q3 | RevenueFromTopReferringDomainsFirstVisitGoogle |
| Q4 | TopPages |

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
**Queries.**
Ad hoc queries.

| Shorthand | Name |
|-----------|------|
| Q0 | Pageviews |
| Q1 | Revenue |
| Q2 | RevenueFromTopReferringDomains |
| Q3 | RevenueFromTopReferringDomainsFirstVisitGoogle |
| Q4 | TopPages |
| Q5 | TopPagesByBrowser |

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
**Queries.**
Ad hoc queries.

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

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
**Queries.**
Ad hoc queries.

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

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
**Queries.**
Ad hoc queries.

| 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

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
**Queries.**
Ad hoc queries.

| 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 | × | | | | × | × | × | |

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
**Queries.**
Ad hoc queries.

| 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 | × |  |  |  | × | × | × |  |
| user_agent |  |  |  |  |  | × |  |  |

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
**Queries.**
Ad hoc queries.

| 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 | × | | | | × | × | × | |
| user_agent | | | | | | × | | |
| visit_referrer | | | × | × | | | | |

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
**Queries.**
Ad hoc queries.

| 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 | × |  |  |  | × | × | × |  |
| user_agent |  |  |  |  |  | × |  |  |
| visit_referrer |  |  | × | × |  |  |  |  |
| post_visid_high |  |  | × | × |  |  | × | × |

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
**Queries.**
Ad hoc queries.

| 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 | × | | | | × | × | × | |
| user_agent | | | | | | × | | |
| visit_referrer | | | × | × | | | | |
| post_visid_high | | | × | × | | | × | × |
| post_visid_low | | | × | × | | | × | × |

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
**Queries.**
Ad hoc queries.

| 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 | × |  |  |  | × | × | × |  |
| user_agent |  |  |  |  |  | × |  |  |
| visit_referrer |  |  | × | × |  |  |  |  |
| post_visid_high |  |  | × | × |  |  | × | × |
| post_visid_low |  |  | × | × |  |  | × | × |
| visit_num |  |  | × | × |  |  | × | × |

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
**Queries.**
Ad hoc queries.

| 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 | × | | | | × | × | × | |
| user_agent | | | | | | × | | |
| visit_referrer | | | × | × | | | | |
| post_visid_high | | | × | × | | | × | × |
| post_visid_low | | | × | × | | | × | × |
| visit_num | | | × | × | | | × | × |
| visit_referrer | | | | | | | | × |

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
**Queries.**
Ad hoc queries.

| 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 | × | | | | × | × | × | |
| user_agent | | | | | | × | | |
| visit_referrer | | | × | × | | | | |
| post_visid_high | | | × | × | | | × | × |
| post_visid_low | | | × | × | | | × | × |
| visit_num | | | × | × | | | × | × |
| visit_referrer | | | | | | | | × |
| hit_time_gmt | | | | | | | × | |

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
**Queries.**
Ad hoc queries.

| 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 | × | | | | × | × | × | |
| user_agent | | | | | | × | | |
| visit_referrer | | | × | × | | | | |
| post_visid_high | | | × | × | | | × | × |
| post_visid_low | | | × | × | | | × | × |
| visit_num | | | × | × | | | × | × |
| visit_referrer | | | | | | | | × |
| hit_time_gmt | | | | | | | × | |
| post_purchaseid | | × | × | × | | | | |

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
**Queries.**
Ad hoc queries.

| 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 | × |  |  |  | × | × | × |  |
| user_agent |  |  |  |  |  | × |  |  |
| visit_referrer |  |  | × | × |  |  |  |  |
| post_visid_high |  |  | × | × |  |  | × | × |
| post_visid_low |  |  | × | × |  |  | × | × |
| visit_num |  |  | × | × |  |  | × | × |
| visit_referrer |  |  |  |  |  |  |  | × |
| hit_time_gmt |  |  |  |  |  |  | × |  |
| post_purchaseid |  | × | × | × |  |  |  |  |
| post_product_list |  | × | × | × |  |  |  |  |

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
**Queries.**
Ad hoc queries.

| 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 | × | | | | × | × | × | |
| user_agent | | | | | | × | | |
| visit_referrer | | | × | × | | | | |
| post_visid_high | | | × | × | | | × | × |
| post_visid_low | | | × | × | | | × | × |
| visit_num | | | × | × | | | × | × |
| visit_referrer | | | | | | | | × |
| hit_time_gmt | | | | | | | × | |
| post_purchaseid | | × | × | × | | | | |
| post_product_list | | × | × | × | | | | |
| first_hit_referrer | | | | × | | | | |

Motivation
Spindle Architecture
Empirical Results
Conclusions

Overview.
Queries.
Ad hoc queries.

# Spindle Architecture
Queries.

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

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
Queries.
**Ad hoc queries.**

# Spindle Architecture
Ad hoc queries.

- ▶ Spark SQL processes relational queries.

Motivation
Spindle Architecture
Empirical Results
Conclusions

Overview.
Queries.
Ad hoc queries.

# 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
```

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
Queries.
**Ad hoc queries.**

# 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]
```

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
Queries.
**Ad hoc queries.**

# 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
```

Motivation
**Spindle Architecture**
Empirical Results
Conclusions

Overview.
Queries.
**Ad hoc queries.**

# 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]
```

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

**Caching.**
Benchmarking concurrent queries.
Remaining.

# Empirical Results
Caching.

- ▶ Six cluster nodes, Spark and HDFS on each.

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

**Caching.**
Benchmarking concurrent queries.
Remaining.

# Empirical Results
Caching.

- Six cluster nodes, Spark and HDFS on each.
- 13.1GB of data.

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

**Caching.**
Benchmarking concurrent queries.
Remaining.

# Empirical Results
Caching.

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

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

**Caching**.
Benchmarking concurrent queries.
Remaining.

Speedups from Caching

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

**Caching**.
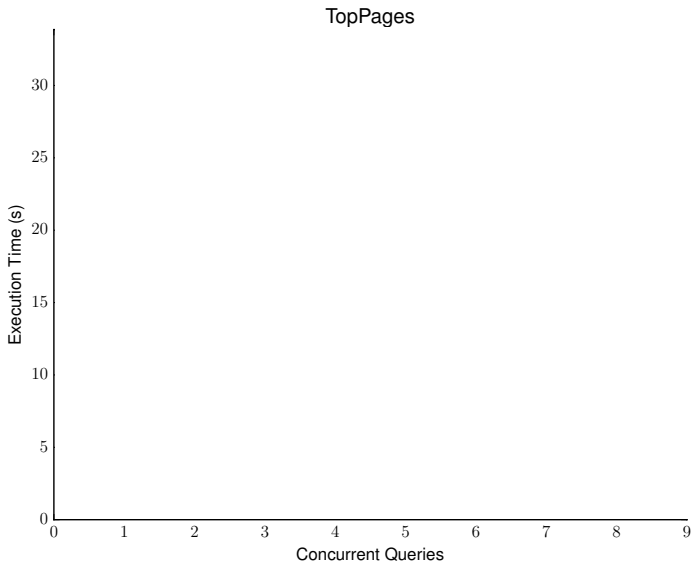Benchmarking concurrent queries.
Remaining.

## Speedups from Caching

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

**Caching.**
Benchmarking concurrent queries.
Remaining.



Speedups from Caching

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

**Caching.**
Benchmarking concurrent queries.
Remaining.

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

**Caching.**
Benchmarking concurrent queries.
Remaining.

Speedups from Caching

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

**Caching.**
Benchmarking concurrent queries.
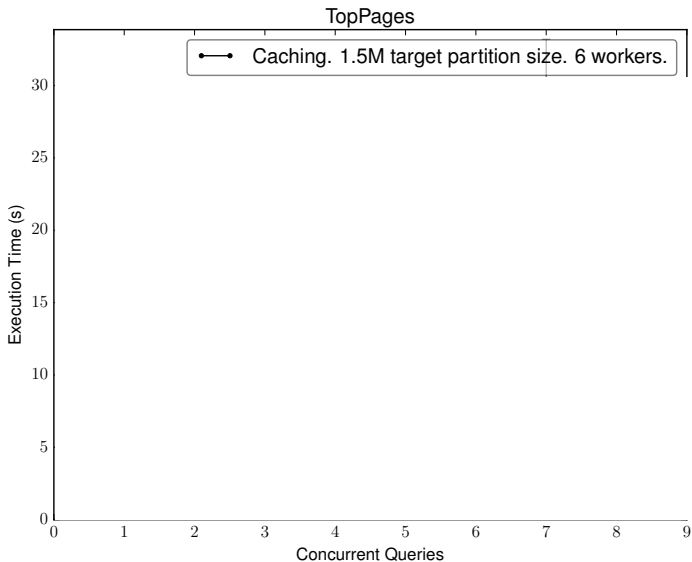Remaining.

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

Caching.
**Benchmarking concurrent queries.**
Remaining.

# Empirical Results
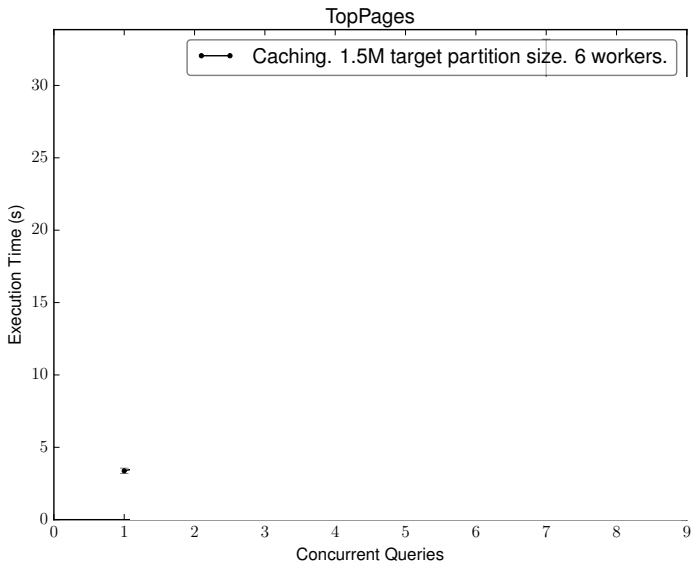Benchmarking concurrent queries.

- ▶ How much will Spindle's performance degrade if multiple users are utilizing it at the same time?

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

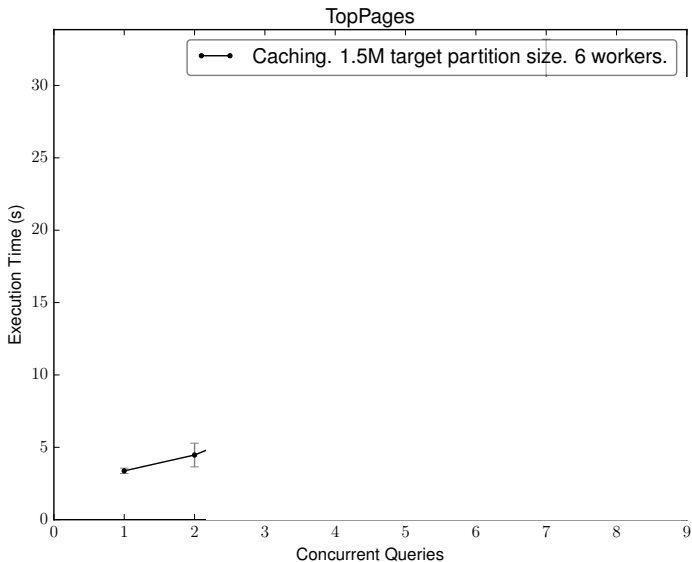Caching.
**Benchmarking concurrent queries.**
Remaining.

# Empirical Results
Benchmarking concurrent queries.

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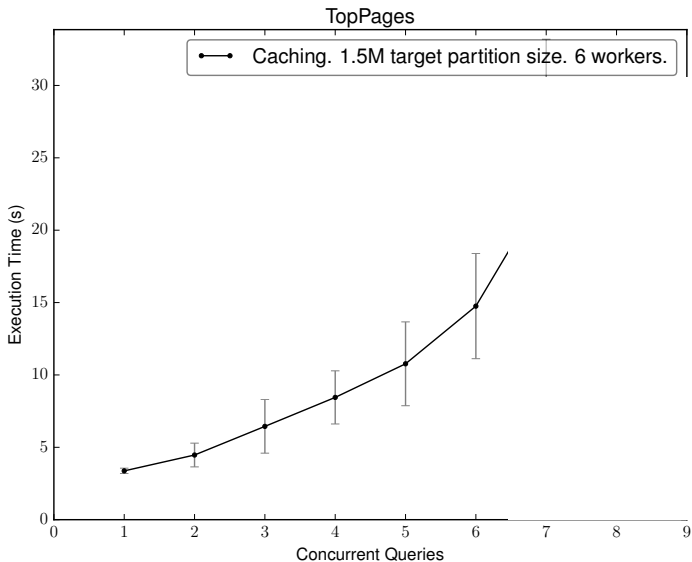- ▶ Concurrently call a single query on the same data.

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

Caching.
**Benchmarking concurrent queries.**
Remaining.

# 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.
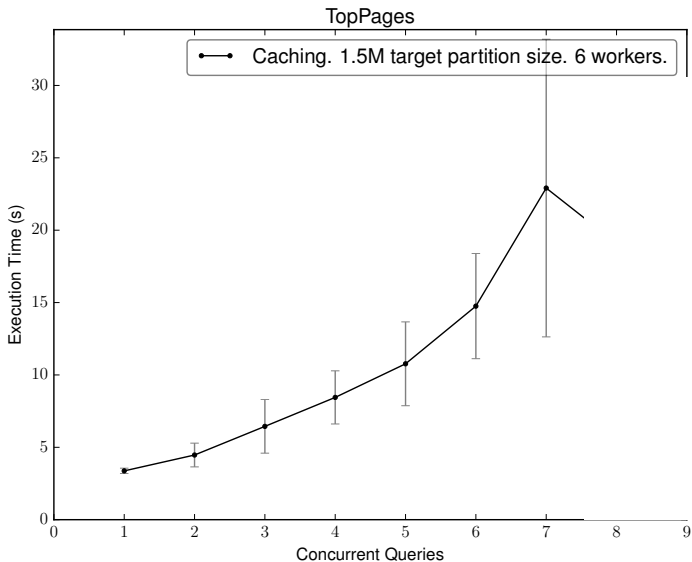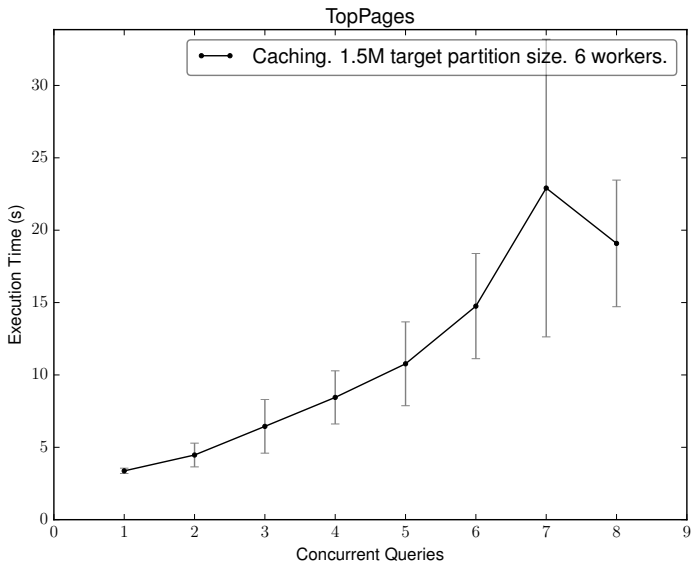- Average the execution time.

TopPages

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

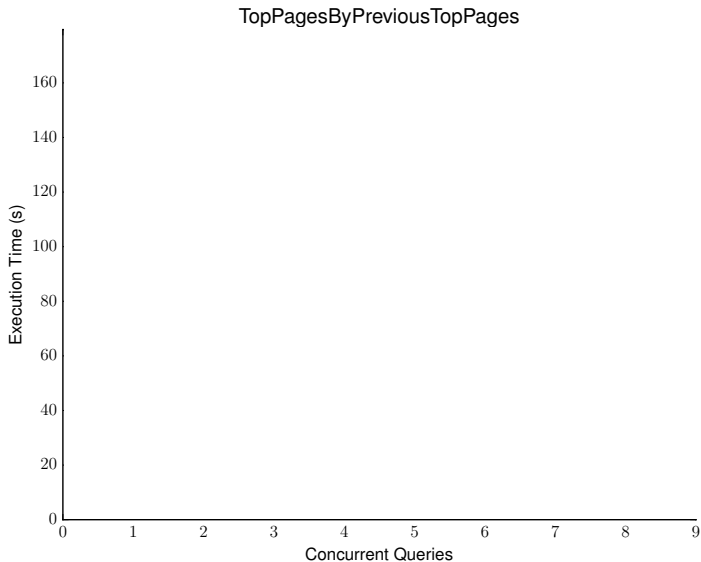Caching.
**Benchmarking concurrent queries.**
Remaining.

TopPages

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

Caching.
**Benchmarking concurrent queries.**
Remaining.



TopPages

Execution Time (s)

Concurrent Queries

Caching. 1.5M target partition size. 6 workers.

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

Caching.
**Benchmarking concurrent queries.**
Remaining.



TopPages

Caching. 1.5M target partition size. 6 workers.

Motivation
Spindle Architecture
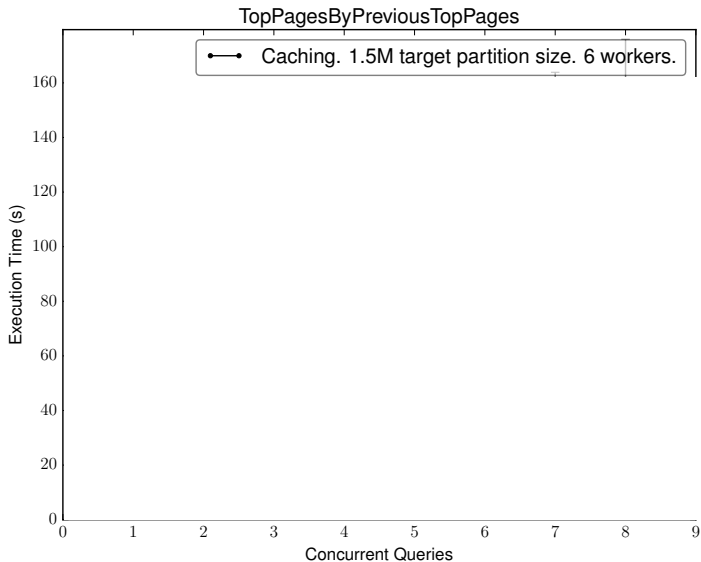**Empirical Results**
Conclusions

Caching.
**Benchmarking concurrent queries.**
Remaining.

TopPages

Caching. 1.5M target partition size. 6 workers.

Execution Time (s)

Concurrent Queries

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

Caching.
**Benchmarking concurrent queries.**
Remaining.

TopPages

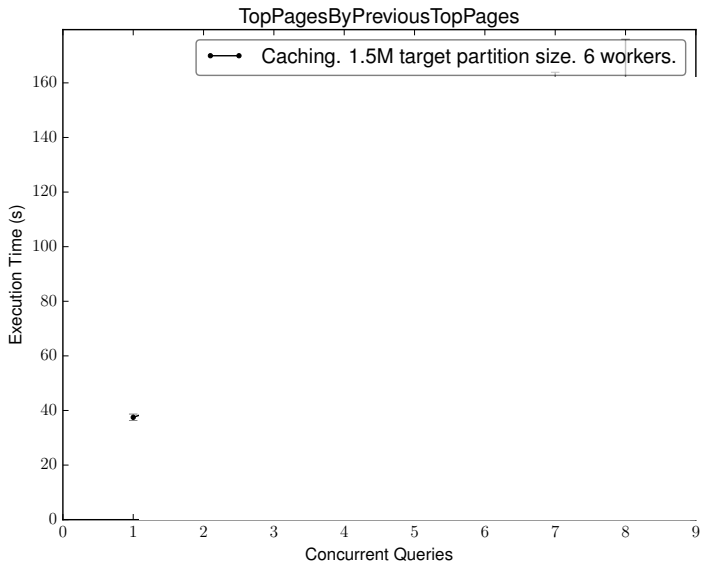Caching. 1.5M target partition size. 6 workers.

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

Caching.
**Benchmarking concurrent queries.**
Remaining.



TopPages

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

Caching.
**Benchmarking concurrent queries.**
Remaining.

TopPages

Caching. 1.5M target partition size. 6 workers.

Execution Time (s)

Concurrent Queries

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

Caching.
**Benchmarking concurrent queries.**
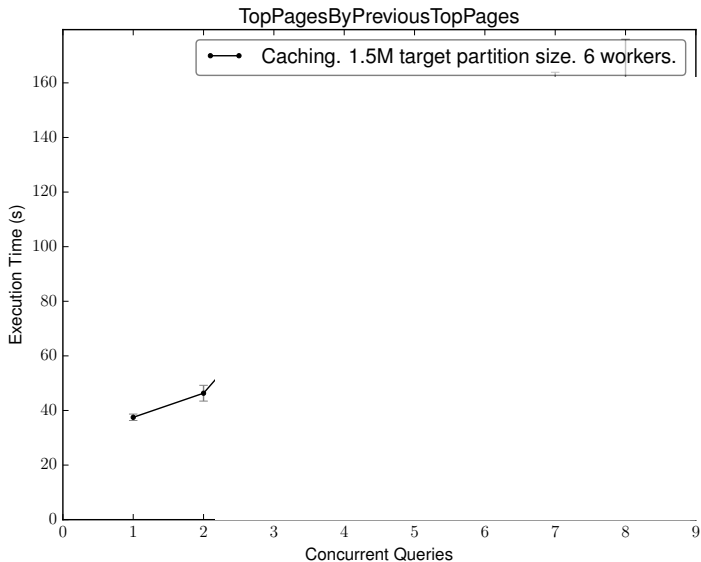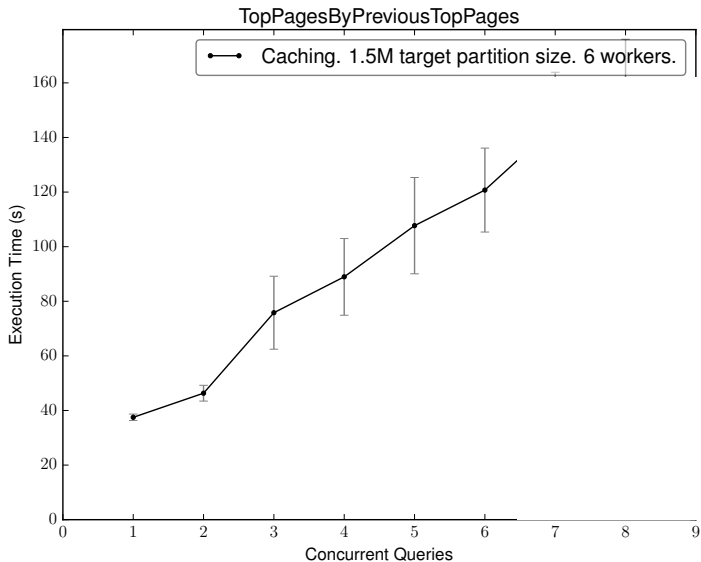Remaining.

TopPagesByPreviousTopPages

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

Caching.
**Benchmarking concurrent queries.**
Remaining.

TopPagesByPreviousTopPages

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

Caching.
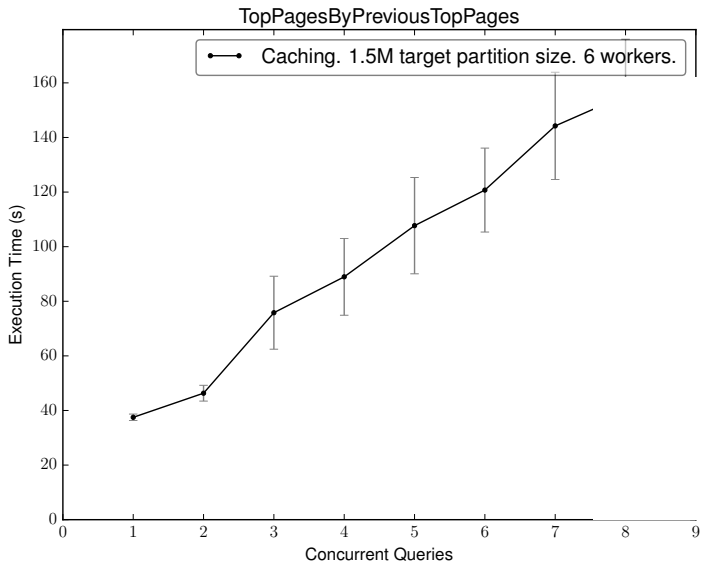**Benchmarking concurrent queries.**
Remaining.

TopPagesByPreviousTopPages

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

Caching.
**Benchmarking concurrent queries.**
Remaining.



TopPagesByPreviousTopPages

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

Caching.
**Benchmarking concurrent queries.**
Remaining.



TopPagesByPreviousTopPages

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

Caching.
**Benchmarking concurrent queries.**
Remaining.



TopPagesByPreviousTopPages

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

Caching.
**Benchmarking concurrent queries.**
Remaining.



TopPagesByPreviousTopPages

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

Caching.
**Benchmarking concurrent queries.**
Remaining.



TopPagesByPreviousTopPages

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

Caching.
Benchmarking concurrent queries.
**Remaining**.

# Empirical Results
Remaining.

- Scaling Spark and HDFS workers.

Motivation
Spindle Architecture
**Empirical Results**
Conclusions

Caching.
Benchmarking concurrent queries.
**Remaining**.

# Empirical Results
Remaining.

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

# Conclusions

- Spark is a good candidate for real-time analytics processing.

# Conclusions

- Spark is a good candidate for real-time analytics processing.
- **Spindle** is an open-source prototype analytics processing engine.

# 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.

# 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.

# 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.

# 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.

# 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` |