

Percolator

Smruti R. Sarangi

Department of Computer Science
Indian Institute of Technology
New Delhi, India

Outline

- 1 Motivation
 - Google's Search Algorithm
 - Requirements
- 2 Design
 - Structure
 - Algorithm
 - Details and Optimizations
- 3 Evaluation

Outline

- 1 Motivation
 - Google's Search Algorithm
 - Requirements
- 2 Design
 - Structure
 - Algorithm
 - Details and Optimizations
- 3 Evaluation

- Updating Google's web index continually is a major challenge.
 - Tens of petabytes of data
 - Billions of updates per day
 - Thousands of machines.
 - Cascading updates.

Google's Search Algorithm

- Every page has a “page rank”.
- The page rank of a popular page is supposed to be high.
- The page rank of a page is determined by the page rank of all the pages that link to it.
- For example:
 - If the New York Times website points to some link, then it has a high page rank. 😊
 - If Hauz Khas Times points to some website, it will have a very low page rank. 😊

Example of a Google Search Query

The screenshot shows a Google search interface. At the top, there's a navigation bar with links: +Smruti, Search, Images, Maps, Play, YouTube, News, Gmail, Drive, Calendar, and More. Below this is the Google logo and a search bar containing the text 'smruti sarangi'. To the right of the search bar is a blue search button with a magnifying glass icon. Below the search bar, there's a row of tabs: Web (selected), Images, Maps, More (dropdown), and Search tools. The search results are displayed below the tabs. At the top of the results, it says 'About 21,400 results (0.37 seconds)'. The first result is 'Homepage of Dr. Smruti Ranjan Sarangi' with a link to 'www.cse.iitd.ac.in/~srsarangi/'. Below this is a brief description: 'Smruti Ranjan Sarangi is an Assistant Professor in the computer science and engineering department at IIT Delhi since January 2011. He primarily works in ...'. There are four sub-links: 'Smruti R. Sarangi' (with a brief description), 'Should I do a Ph.D in the' (with a brief description), 'Teaching' (with a brief description), and 'Publications' (with a brief description). There are also two more sub-links: 'Students' (with a brief description) and 'Contact' (with a brief description). At the bottom of the results, there's a link to 'Smruti Sarangi - India | LinkedIn' with a brief description. At the very bottom, there's a link to 'Smruti Sarangi profiles | LinkedIn' with a brief description.

+Smruti Search Images Maps Play YouTube News Gmail Drive Calendar More ▾

Google smruti sarangi

Web Images Maps More ▾ Search tools

About 21,400 results (0.37 seconds)

[Homepage of Dr. Smruti Ranjan Sarangi](#)
www.cse.iitd.ac.in/~srsarangi/
Smruti Ranjan Sarangi is an Assistant Professor in the computer science and engineering department at IIT Delhi since January 2011. He primarily works in ...

[Smruti R. Sarangi](#)
Smruti R. Sarangi. Contact.
Information. J-203, Brigade ...

[Should I do a Ph.D in the](#)
Should I go to US or Europe for a MS/ Ph.D? This is a ... This ...

[Teaching](#)
teaching. 2012-2013. 2013 Spring Semester : CSL 860 ...

[Publications](#)
publications. (Please note that all of these papers are only ...

[Students](#)
research group members. Ph.D Students. Gayathri ...

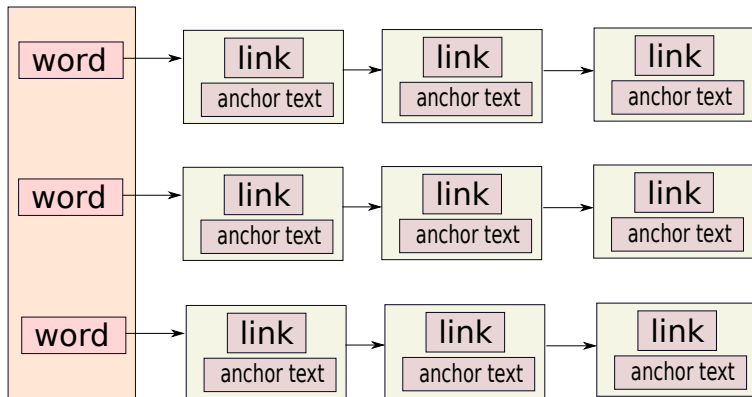
[Contact](#)
Please don't hesitate to contact me at anytime with queries of any ...

[More results from iitd.ac.in »](#)

[Smruti Sarangi - India | LinkedIn](#)
in.linkedin.com/pub/smruti-sarangi/6/28/112
New Delhi Area, India - Assistant Professor, Computer Science Dept, IIT Delhi
View **Smruti Sarangi's** (India) professional profile on LinkedIn. LinkedIn is the world's largest business network, helping professionals like **Smruti Sarangi** ...

[Smruti Sarangi profiles | LinkedIn](#)
www.linkedin.com/pub/dir/Smruti/Sarangi
View the profiles of professionals named **Smruti Sarangi** on LinkedIn. There are 5 professionals named **Smruti Sarangi**, who use LinkedIn to exchange ...

Structure of a Web Index



The Problem of Updates

- The links in the inverted list are arranged according to their page rank.
- If the page rank of a website changes then:
 - We need to update the inverted list to reflect the change.
 - The page rank of sites that it points to need to change.
 - This problem is known as **cascading update** .

Outline

- 1 Motivation
 - Google's Search Algorithm
 - Requirements
- 2 Design
 - Structure
 - Algorithm
 - Details and Optimizations
- 3 Evaluation

Requirements of a Solution

- Should provide ACID transaction semantics (do not want to corrupt database).
- Should have high throughput, and acceptable latency.
- Should be able to handle petabytes of data.
- Traditional DBMS systems are too slow → Need new technology
- Random access to data such that changes can **percolate**
- Consistency Model: Snapshot Isolation

Snapshot Isolation

- Assume two concurrent updates to a linked list.
 - If they do not access the same node or its parent, then they are disjoint.
 - Disjoint accesses can continue in parallel.
 - This is **different** from regular transaction semantics such as serializability.
- **Definition :**
 - When a transaction starts, it takes (appears to) a consistent **snapshot** of the entire database.
 - It then proceeds to update its private copy of the database.
 - The values are committed if they have not been changed by another transaction since the snapshot.

Outline

- 1 Motivation
 - Google's Search Algorithm
 - Requirements
- 2 Design
 - **Structure**
 - Algorithm
 - Details and Optimizations
- 3 Evaluation

Design of Percolator

- Built on top of Bigtable – Google's distributed storage engine
- Bigtable is a multidimensional database
 - Distributed key-value store
 - We save – row, column, timestamp
 - Atomic read-modify-write operations for each row
 - Meta data is stored in separate columns
- Observer framework
 - Any row has a set of observers.
 - They run specialized functions when data in the row changes.

Model of Transactions

- Provides support for ACID transactions
 - Hard to do in such a large database
 - Required: do not want to have Google's database in an inconsistent state
 - Uses timestamp for each data item
 - The set of timestamps at the beginning of a transaction is its snapshot.
- Transactions can include multiple rows across multiple BigTable tables
- Percolator implements its own lock service
- Percolator adds a special column to save locks.

Columns in BigTable

Column	Use
lock	contains a pointer to the lock
write	timestamp of committed data
data	data value
notify	list of observers
ack_ O	last timestamp at which observer O ran

Example

A transfers B 7\$

key	data	lock	write
A	6: 5:10\$	6: 5:	6:data@5 5:
B	6: 5:2\$	6: 5:	6:data@5 5:

key	data	lock	write
A	7: 3\$ 6: 5:10\$	7: primary 6: 5:	7: 6:data@5 5:
B	6: 5:2\$	6: 5:	6:data@5 5:

Example - II

key	data	lock	write
A	7:3\$	7: primary	7:
	6:	6:	6:data@5
	5:10\$	5:	5:
B	7: 9\$	7: primary@A	7:
	6:	6:	6:data@5
	5:2\$	5:	5:

key	data	lock	write
A	8:	8:	8: data @ 7
	7:3\$	7:	7:
	6:	6:	6:data@5
	5:10\$	5:	5:
B	7:9\$	7: primary@A	7:
	6:	6:	6:data@5
	5:2\$	5:	5:

Example - III

key	data	lock	write
A	8:	8:	8: data @ 7
	7:3\$	7:	7:
	6:	6:	6:data@5
	5:10\$	5:	5:
B	8:	8:	8: data @ 7
	7:9\$	7:	7:
	6:	6:	6:data@5
	5:2\$	5:	5:

Outline

- 1 Motivation
 - Google's Search Algorithm
 - Requirements
- 2 Design
 - Structure
 - **Algorithm**
 - Details and Optimizations
- 3 Evaluation

Algorithm: Begin Transaction

Algorithm 1: Begin Transaction

- 1 $startTs \leftarrow \text{oracle.getTimeStamp}()$
- 2 Set(W):
 $\text{writes.push}(W)$

Get Method

```
1 Get(row, column, value):
   while True do
2      $T \leftarrow \text{startTrans}(\text{row})$ 
     if  $T.\text{hasLock}(0, \text{startTs})$  then
3       backOffAndMaybeRemoveLock(row, col)
       continue
4     end
5      $\text{latestWrite} \leftarrow T.\text{read}(\text{row}, [0, \text{startTs}])$ 
     if ! $\text{latestWrite}$  then
6       return  $\phi$ 
7     end
8      $\text{dataTs} \leftarrow \text{latestWrite.timeStamp}$ 
     return ( $T.\text{read}(\text{row}, \text{"data"}, \text{dataTs})$ )
9 end
```

PreWrite

```
1 PreWrite(Write  $w$ , Write  $primary$ )  
   Column  $col \leftarrow w.col$   
    $T \leftarrow \text{startTransaction}(w.row)$   
  
2 if  $T.read(w.row, "write", [startTs, \infty])$  then  
3   |   return false  
4 end  
5 if  $T.read(w.row, "lock", [0, \infty])$  then  
6   |   return false  
7 end  
  
8  $T.write(w.row, "data", startTs, w.value)$   
    $T.write(w.row, "lock", startTs, \{primary.row, primary.col\})$   
   return  $T.commit()$ 
```

Commit - I

```
1 Commit()
   /* Prewrite all the entries */
2 (primary, secondaries)  $\leftarrow$  (writes[0], writes[1 ... n])
   if !PreWrite(primary, primary) then
3   |   return false
4 end
5 for Write w: secondaries do
6   |   if !PreWrite(w, primary) then
7   |   |   return false
8   |   end
9 end
10 commitTs  $\leftarrow$  oracle.getTimeStamp()
```

Commit - II

```
/* Commit the primary */
11  $T \leftarrow \text{startTransaction}(\text{primary.row})$ 

/* Test to see if aborted by somebody else */
12 if ! $T.\text{read}(\text{primary.row}, \text{"lock"}, \text{startTs})$  then
13   | return false
14 end

/* Write the primary and erase the lock */
15  $T.\text{write}(\text{primary.row}, \text{"write"}, \text{commitTs}, \text{"data@"} + \text{startTs})$ 
    $T.\text{erase}(\text{primary.row}, \text{"lock"}, \text{commitTs})$ 

/* Point of commit */
16 if ! $T.\text{Commit}()$  then
17   | return false
18 end
```


Commit - III

```
19 for Write w: secondaries do  
20   |   write(w.row, "write", commitTs, "data@"+startTs)  
    |   erase (w.row, "lock", commitTs)  
21 end  
22 return true
```

Outline

1

Motivation

- Google's Search Algorithm
- Requirements

2

Design

- Structure
- Algorithm
- Details and Optimizations

3

Evaluation

Timestamps

- The timestamp oracle needs to be able to sustain very high throughput.
- Possible to batch several RPC calls to the oracle to reduce network load.
- Needs to give out timestamps in increasing order.
- If it fails, then it needs to recover and issue timestamps that are greater than the ones it issued earlier.

Observers

- Each observer registers a set of columns, and a function.
- The function gets invoked, if any of the columns are updated.
- Possible to do **message collapsing**
- At most one observer's transaction will commit per column.
- Steps in running an observer
 - After an update to a column, Percolator sets the notify column.
 - A worker thread, ultimately picks up this information, and runs an observer.
 - If the latest timestamp of an observer run (**ack_O**) is less than the commit timestamp of the update, then run the observer.
 - Worker threads avoid **clumping** by scanning random parts of the database.

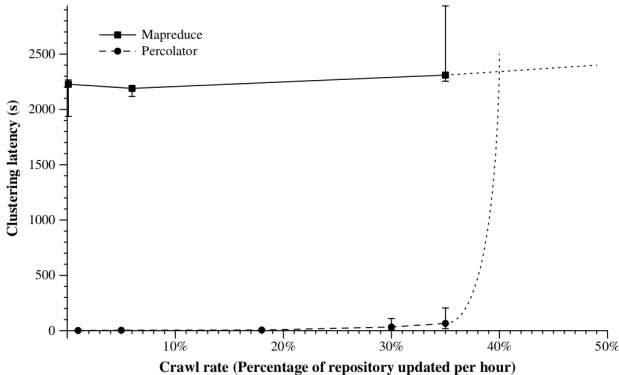
Performance Improvements

- Support for **read-modify-write** RPCs in BigTable.
- Create batches of RPC calls.
- Employ pre-fetching to reduce reads.
- Use blocking API calls, and a large number of threads to simplify the programming model.

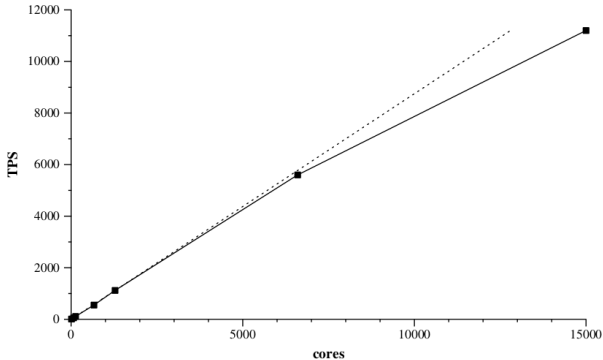
Setup

- Existing Setup:
 - Crawl billions of documents
 - Series of 100 map-reduces
 - document takes 2-3 days for getting indexed
- Percolator based indexing system – **Caffeine**
 - 100x faster
 - Average age of documents gets reduced to 50%

Performance vs Crawl Rate



Scalability for TPC/E benchmarks



Close to Linear Scaling



Large-scale Incremental Processing Using Distributed Transactions and Notifications by Daniel Peng and Frank Dabek, OSDI, 2010