

Distributed Data Processing with Hadoop

Johan Oskarsson, johan@last.fm
Martin Dittus, martind@last.fm

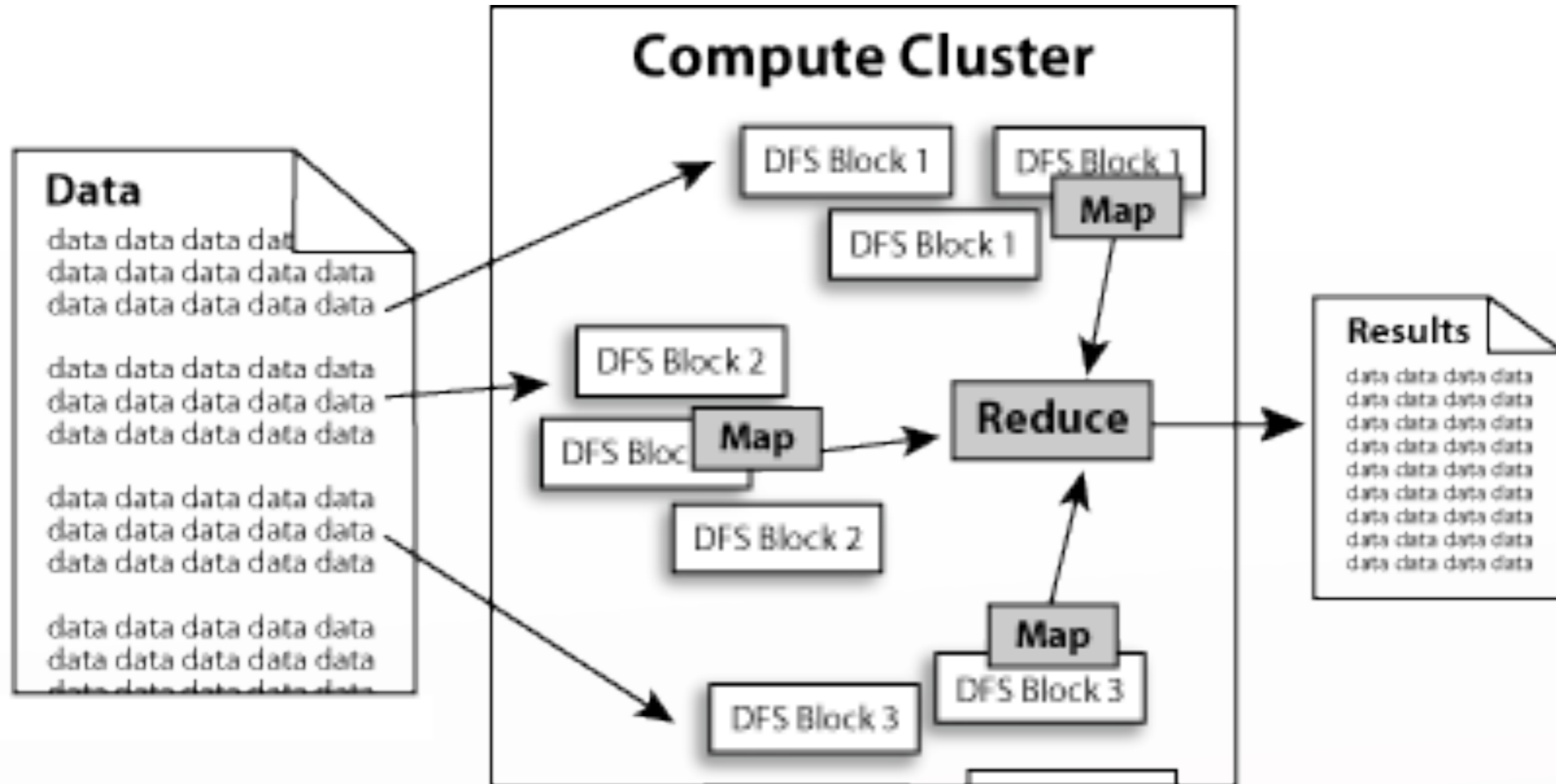


the social music revolution **last.fm**TM

Hadoop

- Apache Project
- Distributed Filesystem
- Distributed Job Execution
- Java Implementation of MapReduce

Hadoop: MapReduce



Hadoop Features

- Runs on Commodity Hardware
- Recovers from Hardware Failures
- Hadoop Programs in Java, Python, C++
- Hadoop Streaming
- S3/EC2
- ...

Hadoop Features

```
File Edit View Terminal Tabs Help
master1 Hadoop Map/Reduce Administration (p1 of 5)
master1 Hadoop Map/Reduce Administration

Started: Wed Nov 28 18:19:35 UTC 2007
Version: 0.14.3-dev, r587669
Compiled: Tue Nov 27 16:33:34 UTC 2007 by hadoop

-----

Cluster Summary

                Maps Reduces Tasks/Node Total Submissions Nodes
                104  174      5           66           35

-----

Running Jobs

Running Jobs
Jobid User Name Map % complete Map total Maps completed Reduce % complete Reduce total Reduces completed
job_200711281819_0062 hadoop distcp 42.66% 10000 4228 0.00% 0 0
job_200711281819_0063 hadoop elias-/user/hadoop/elias/abstats.tmp 0.00% 240 0 0.00% 73 0
job_200711281819_0064 hadoop elias-stationstats-trr2-/user/hadoop/elias/stationstats.tmp.2 0.00% 292 0 0.00% 73 0
job_200711281819_0065 hadoop metrics-/user/hadoop/processed/streaming_logjoin/flp/2007/11/25 0.00% 260 0 0.00% 73 0
job_200711281819_0066 hadoop elias-p1-/user/hadoop/elias-fnar-overall.tmp.1 0.00% 240 0 0.00% 73 0

-----

Completed Jobs

Completed Jobs
Jobid User Name Map % complete Map total Maps completed Reduce % complete Reduce total Reduces completed
job_200711281819_0001 hadoop sorter 100.00% 5600 5600 100.00% 175 175
job_200711281819_0002 hadoop elias-/user/hadoop/elias/abstats.tmp 100.00% 239 239 100.00% 73 73
job_200711281819_0003 hadoop elias-stationstats-mbpi-/user/hadoop/elias/stationstats.tmp 100.00% 240 240 100.00% 73

-- press space for next page --
Arrow keys: Up and Down to move. Right to follow a link; Left to go back.
H)elp O)ptions P)rint G)o M)ain screen Q)uit /=search [delete]=history list
0 -bash 1 -bash 2 -bash 3 -bash 4 -bash 5 -bash 6 -bash 7 -bash 8 -bash 9 -bash 10 -bash 11 -bash gimli 11/28
```

Hadoop Features

```
File Edit View Terminal Tabs Help
Hadoop job_200711281819_0062 on master1

REFRESH(10 sec): http://localhost:50030/jobdetails.jsp?jobid=job_200711281819_0062&refresh=10

Hadoop job_200711281819_0062 on master1

User: hadoop
Job Name: distcp
Job File: /home/hadoop/mapred/system/job_200711281819_0062/job.xml
Status: Running
Started at: Wed Nov 28 20:45:07 UTC 2007
Running for: 9mins, 34sec

-----

Kind % Complete Num Tasks Pending Running Complete Killed Failed/Killed
Task Attempts
map 48.57% 10000 5057 142 4801 0 39 / 0
reduce 0.00% 0 0 0 0 0 0 / 0

Job Counters
          Counter          Map  Reduce  Total
Launched map tasks          0      0  4,982
Data-local map tasks        0      0  4,982
Map-Reduce Framework Map input records 14,657  0 14,657
                          Map input bytes 704,973  0 704,973

-----

Change priority from NORMAL to: VERY_HIGH HIGH LOW VERY_LOW

-----

Go back to JobTracker
Hadoop, 2006.

Commands: Use arrow keys to move, '?' for help, 'q' to quit, '<' to go back.
Arrow keys: Up and Down to move. Right to follow a link; Left to go back.
H)elp O)ptions P)rint G)o M)ain screen Q)uit /=search [delete]=history list
0 -bash 1 -bash 2 -bash 3 -bash 4 -bash 5 -bash 6 -bash 7 -bash 8 -bash 9 -bash 10 -bash 11 -bash gimli 11/28
```

Application Examples

- Nutch/Lucene: Document Indices
- Last.fm: Music Charts
- Yahoo: Search, ...?
- Google, Facebook: ...?
- NY Times: Document Conversion
- "Forgot My Password Again"

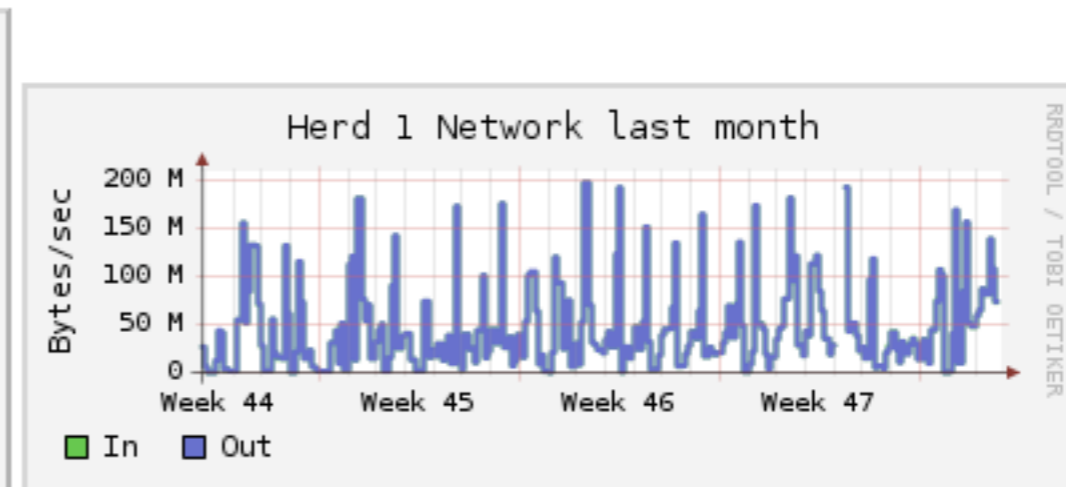
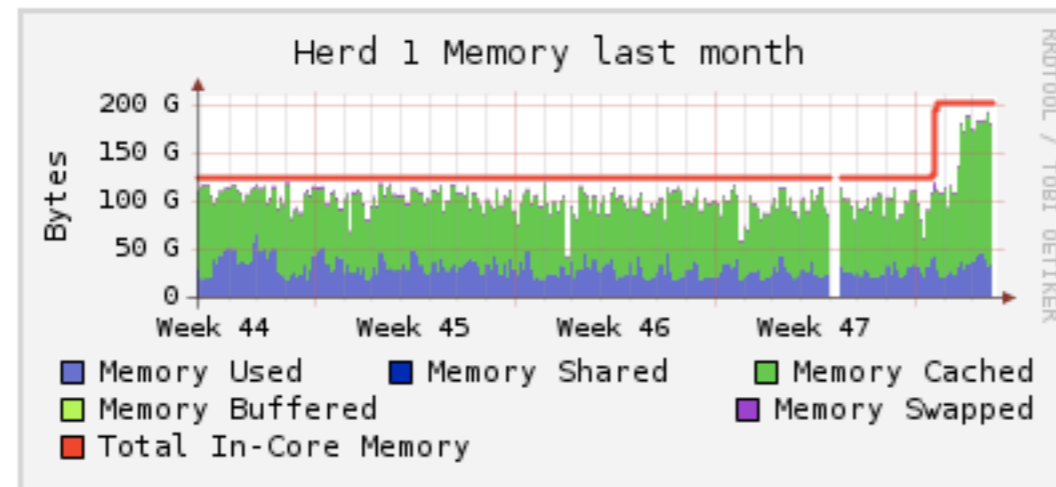
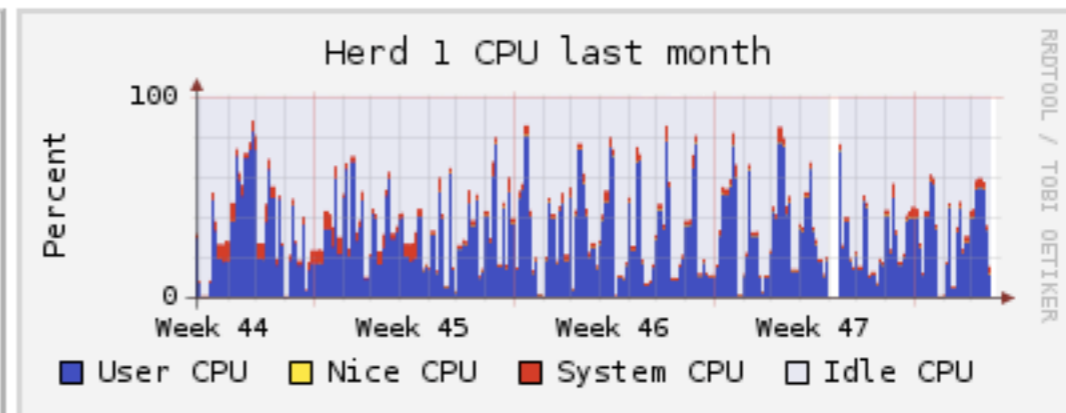
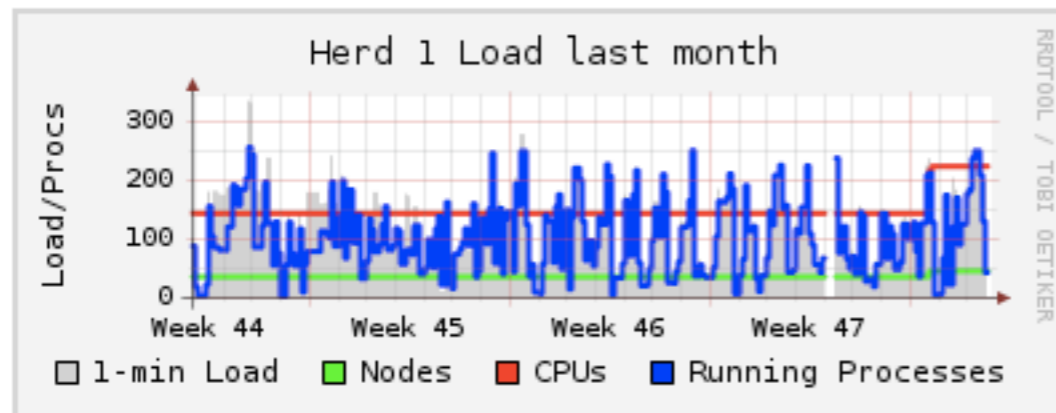
Related Projects

- HBase (BigTable Implementation)
- Hadoop On Demand
- Pig
- Hive
- MapReduce Tools for Eclipse
- ...

Hadoop at Last.fm

- DFS Stores 30 Terabytes of Data
- 2 Separate Clusters
- 5 Users
- ~200 Programs
- ~10k Jobs/month
 - Scheduled Jobs
 - Ad-Hoc Processing

Hadoop at Last.fm



Motivation for Last.fm

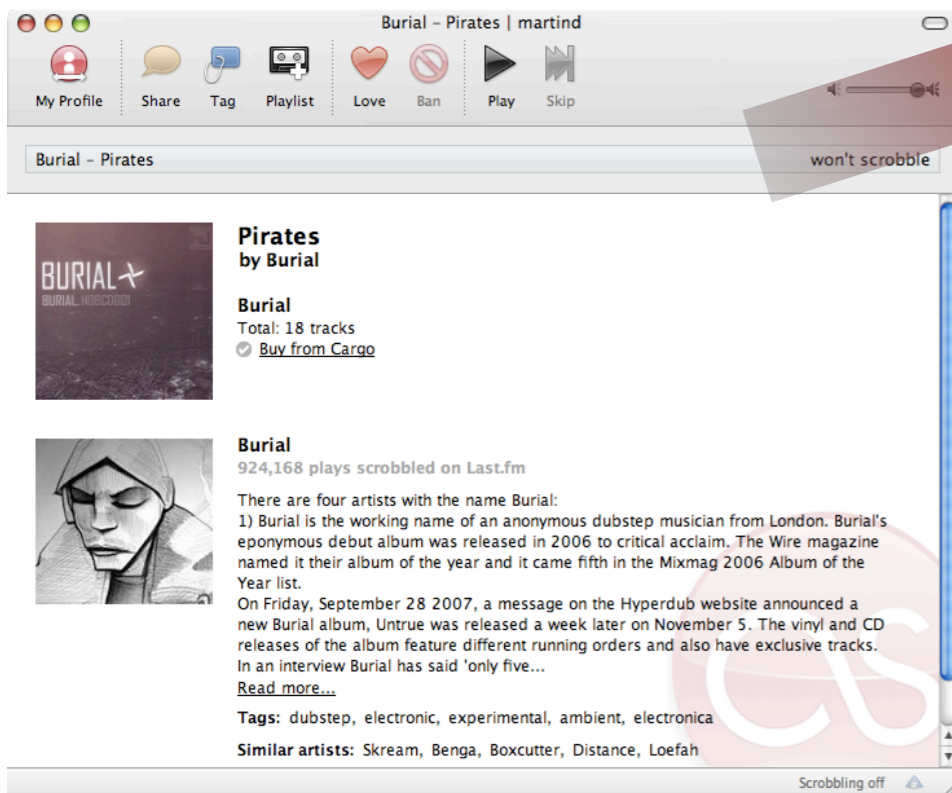
- Database was Bottleneck
 - Scaling DB Clusters: Expensive + Brittle
- vs. Hadoop
 - Throw More Cheap Machines at the Problem
 - Trade-off: Losing Real-time Updates

Uses: Charts

Show: for

Top Albums

	1	J Dilla – Donuts	66
	2	▶ Quasimoto – The Further Adventures of Lord Quas	48
	3	▶ Maurizio – M-Series	43
	4	▶ John Coltrane – Giant Steps	38
	5	▶ Verbal Kent – Move With the Walls	32
	6	Basic Channel – BCD	30
	7	▶ Burial – Burial	28
	7	▶ Oh No – Exodus Into Unheard Rhythms	28
	9	Braun and the Mob – As the Veneer of Dumbness	26
	10	▶ Freeform – Condensed	24
	10	▶ RZA – RZA as Bobby Digital in Stereo	24
	10	▶ Shank – Do	24
	13	▶ Wu-Tang Clan – Enter the Wu-Tang: 36 Chambers	23
	14	▶ ROOT 70 – Heaps Dub	20
	15	▶ Percee P – Perseverance	19
	16	▶ Various Artists – Difficult Easy Listening	18
	17	▶ Chris Clark – Clarence Park	17
	18	▶ Chris Clark – Empty the Bones of You	16
	19	▶ Burnt Friedman – Con Ritmo	15
	19	▶ Burnt Friedman & Jaki Liebezeit – Secret Rhythms 2	15
	19	▶ Various Artists – Replicant Rumba Rockers: A	15
	22	▶ John Coltrane – A Love Supreme Deluxe Edition	13



Burial - Pirates | martind

My Profile Share Tag Playlist Love Ban Play Skip

Burial - Pirates won't scrobble

Pirates
by Burial

Burial
Total: 18 tracks
Buy from Cargo

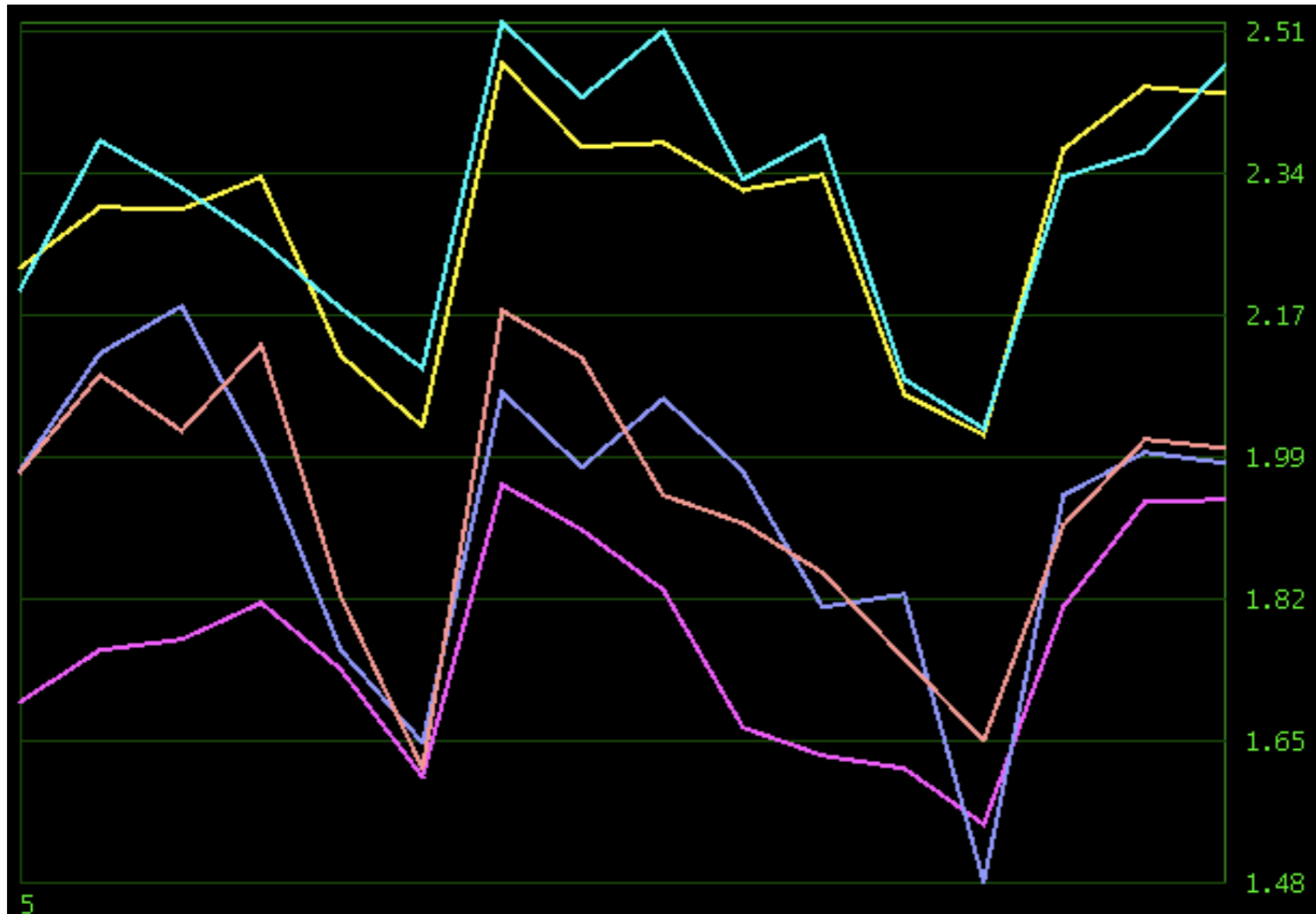
Burial
924,168 plays scrobbled on Last.fm

There are four artists with the name Burial:
1) Burial is the working name of an anonymous dubstep musician from London. Burial's eponymous debut album was released in 2006 to critical acclaim. The Wire magazine named it their album of the year and it came fifth in the Mixmag 2006 Album of the Year list.
On Friday, September 28 2007, a message on the Hyperdub website announced a new Burial album, Untrue was released a week later on November 5. The vinyl and CD releases of the album feature different running orders and also have exclusive tracks. In an interview Burial has said 'only five...'
[Read more...](#)

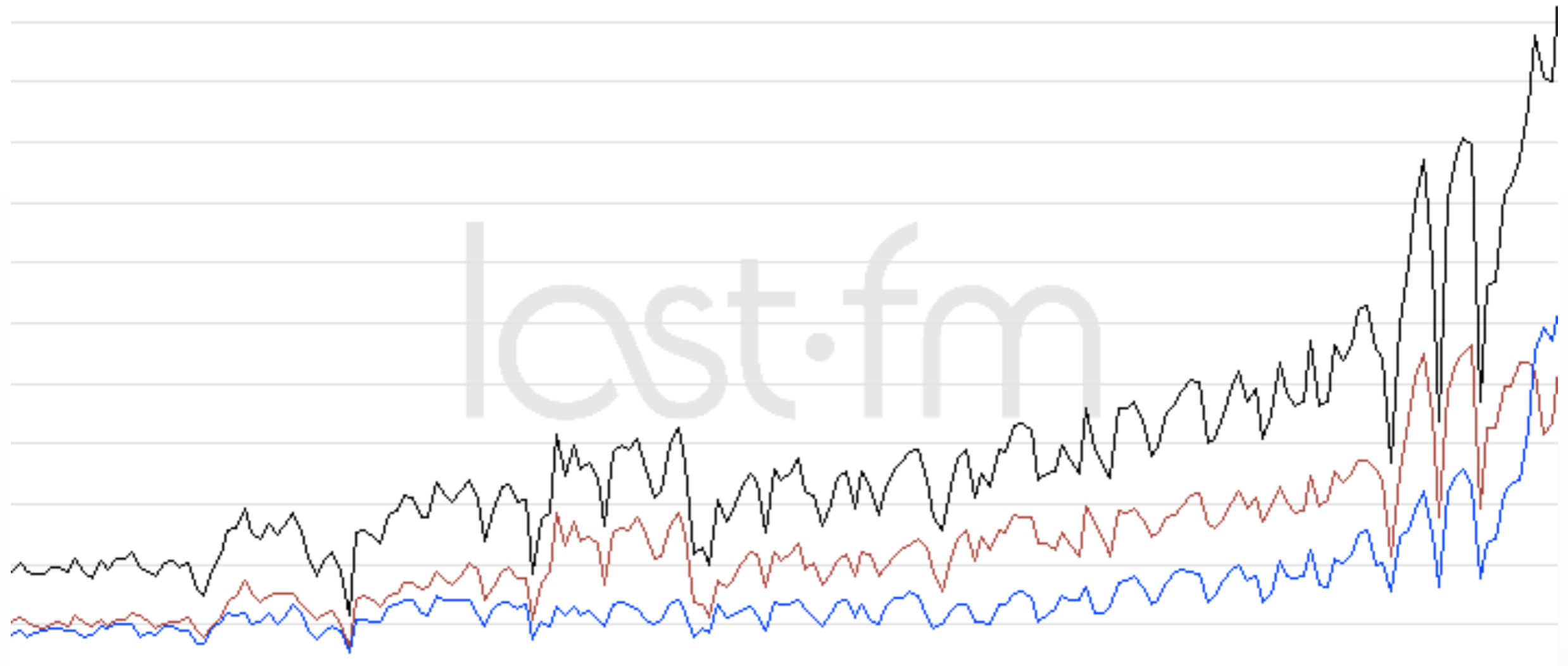
Tags: dubstep, electronic, experimental, ambient, electronica
Similar artists: Skream, Benga, Boxcutter, Distance, Loefah

Scrobbling off

Uses: QA, A/B Testing



Uses: Long-term Stats



Current Issues

- Scheduling Conflicts
- Debugging Distributed Applications
- Application Development Time

Questions?

- lucene.apache.org/hadoop/
- #hadoop on freenode.net

- We're Hiring!
 - last.fm/about/jobs
 - martind@last.fm