



Big Data at Netflix: Faster and Easier

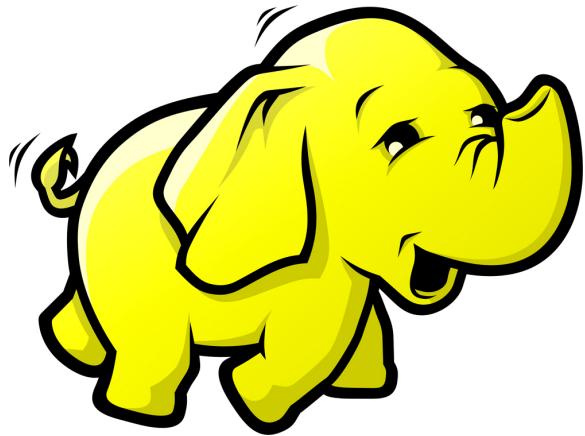
Kurt Brown

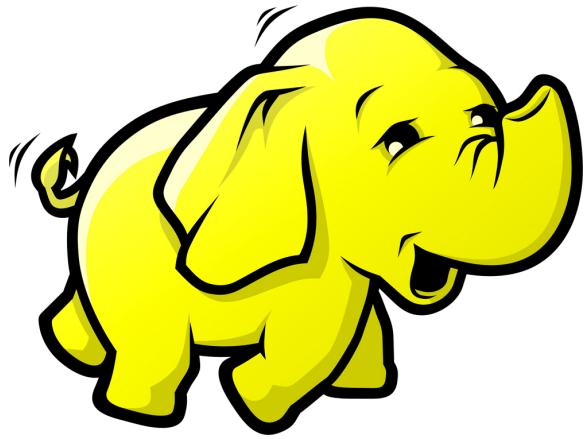
NETFLIX

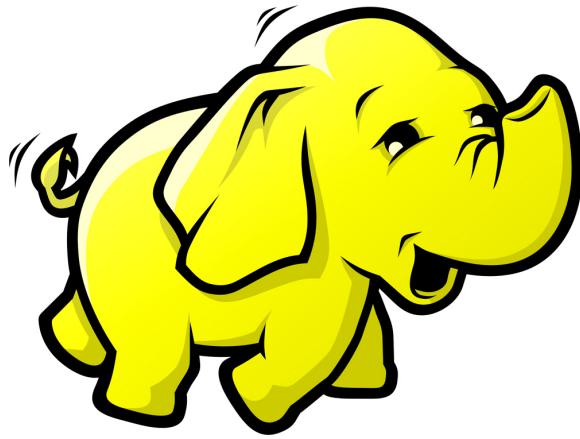
- 1/3+ Internet traffic (NA / peak)
- 100+ Million hours per day
- 65+ Million members / 50+ countries
- Outside the US – streaming only









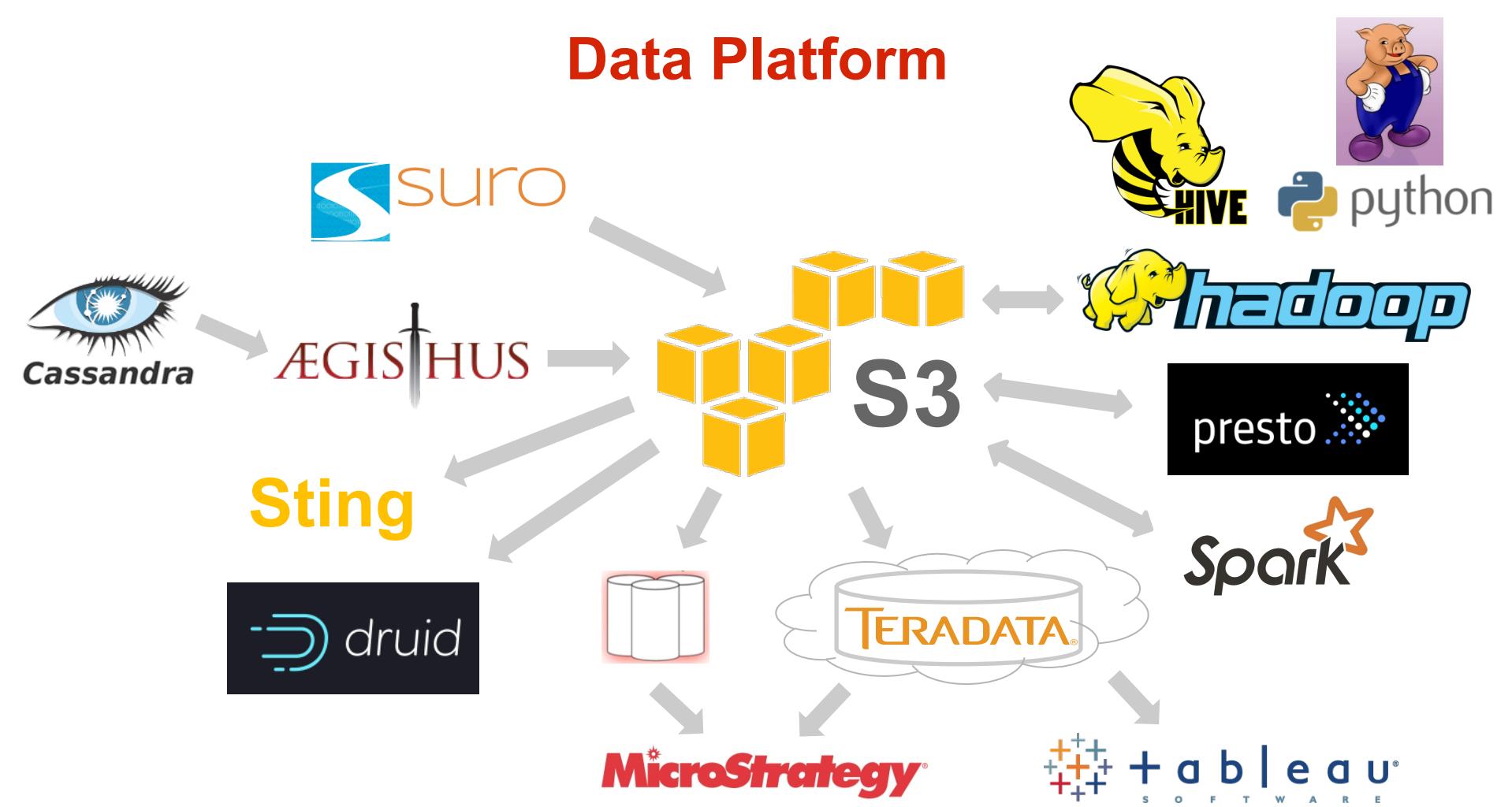




500 Billion Events / Day



Data Platform



Faster

presto 

Why Presto (vs. Alternatives)?

- AWS Interplay
- Open Source
- Java
- Real-world Big Data
- Future

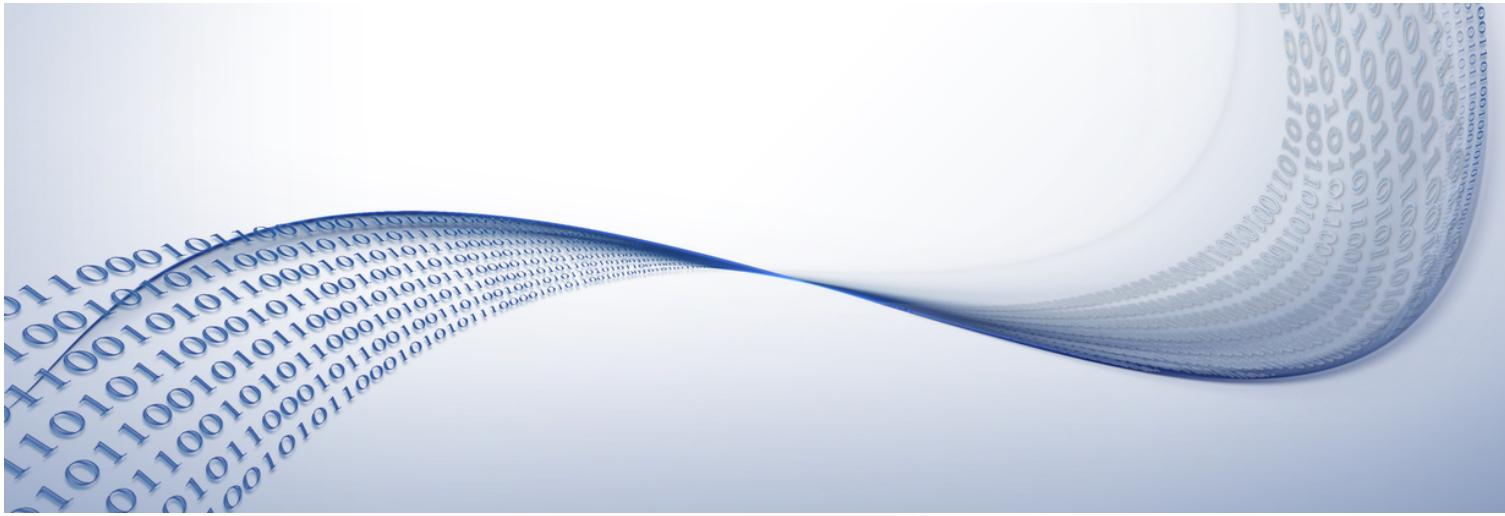


Takeaways – “Good”

- Cohesive Environment
- Multiple Language Support
- Performance (100X???)
- Community
- Future

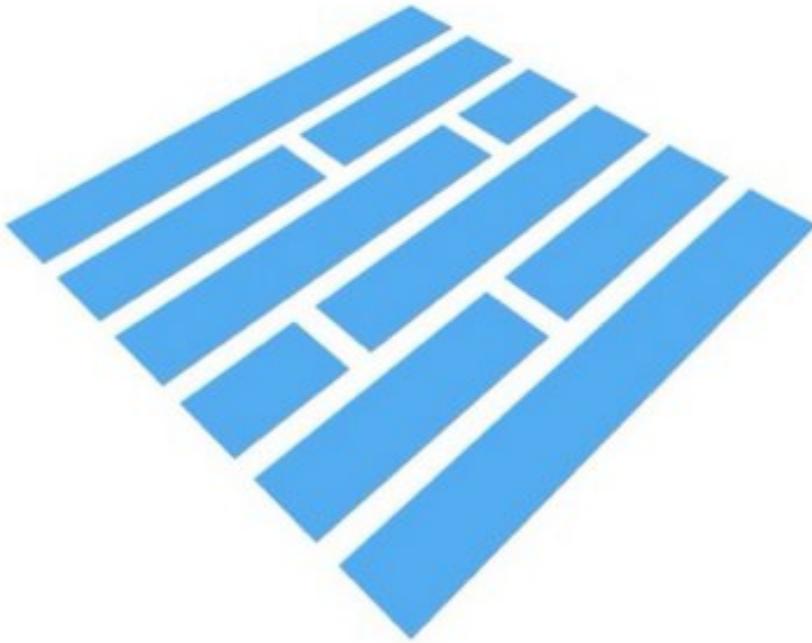
Takeaways – “Bad”

- Maturity
- Running at Scale
- Multi-tenancy / Concurrency
- Tuning?
- Scala?
- Current Investments?



Stream Processing

- Consistency / Ready for use
- Latency
- Scaling
- Storm? Samza? Spark?



Parquet

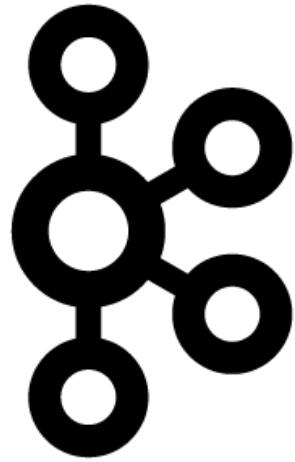
Parquet

- Columnar
- Cross-project
- Community
- Contributions



TEZA

Easier

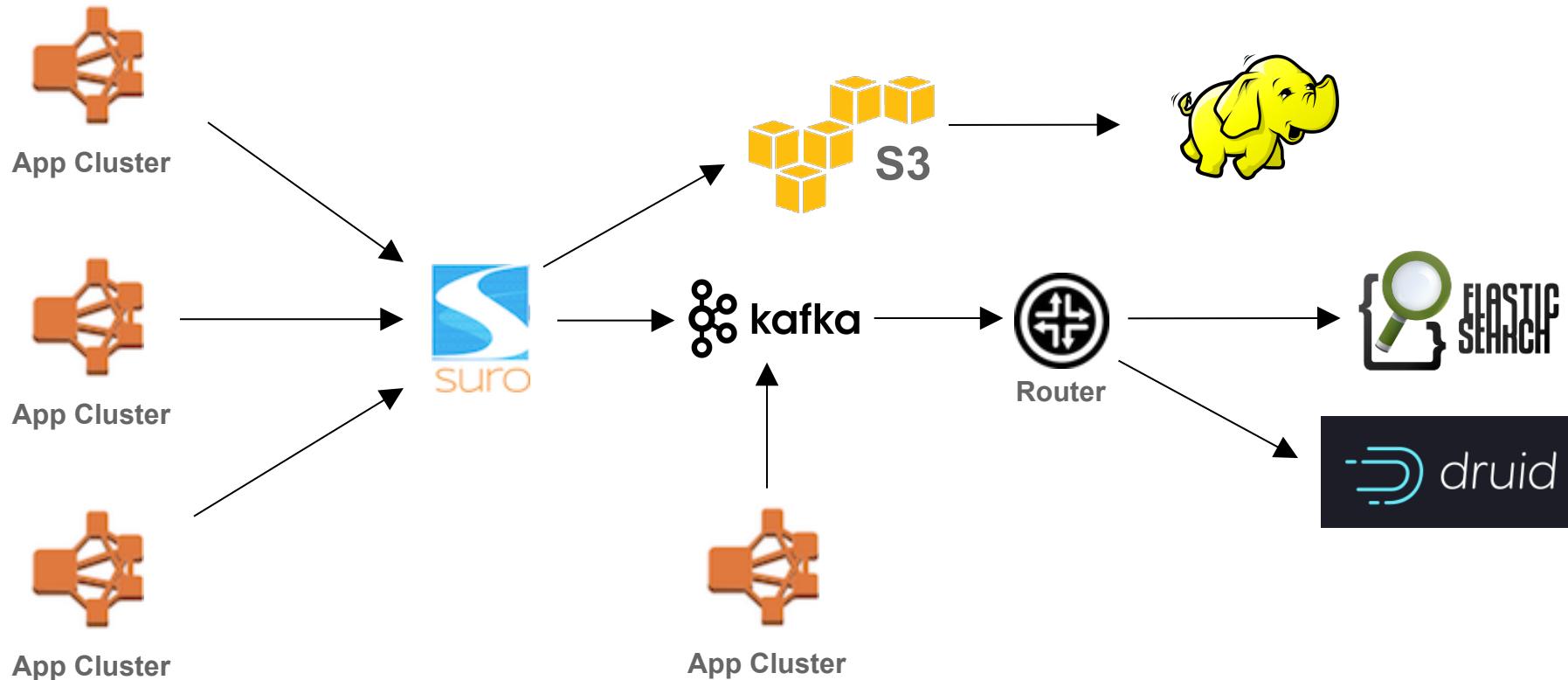


kafka

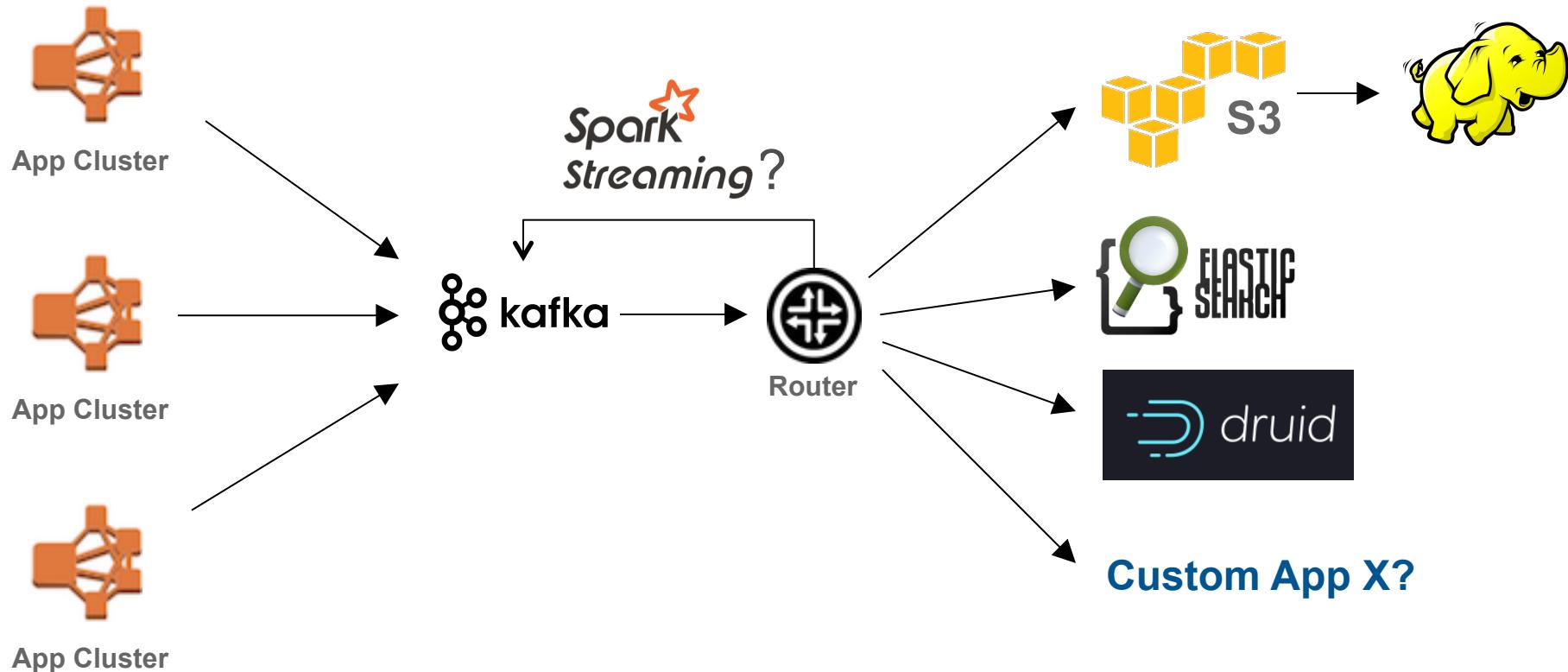
Kafka

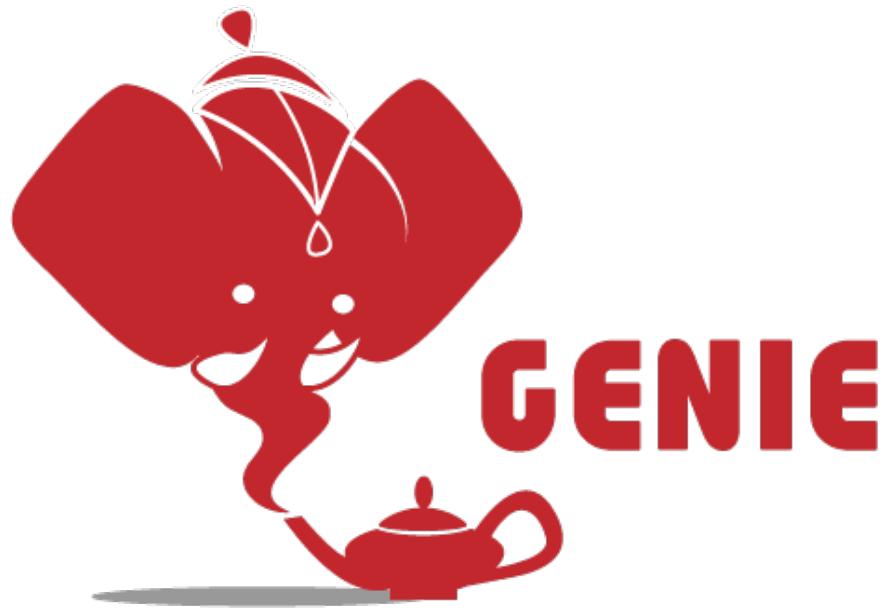
- Simplify our Architecture
- Durability
- Community

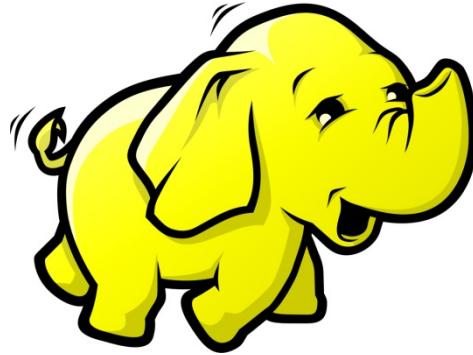
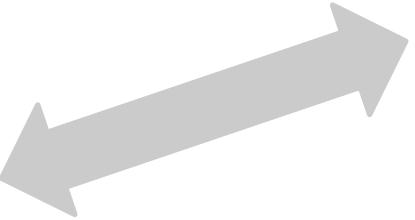
Current Data Pipeline

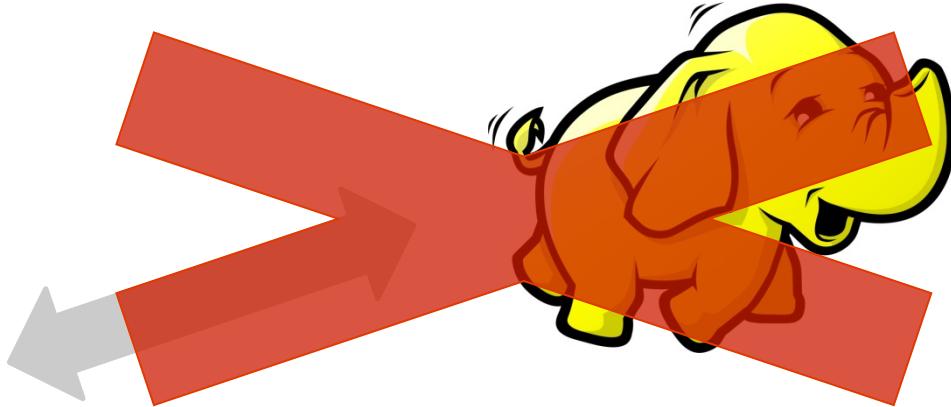


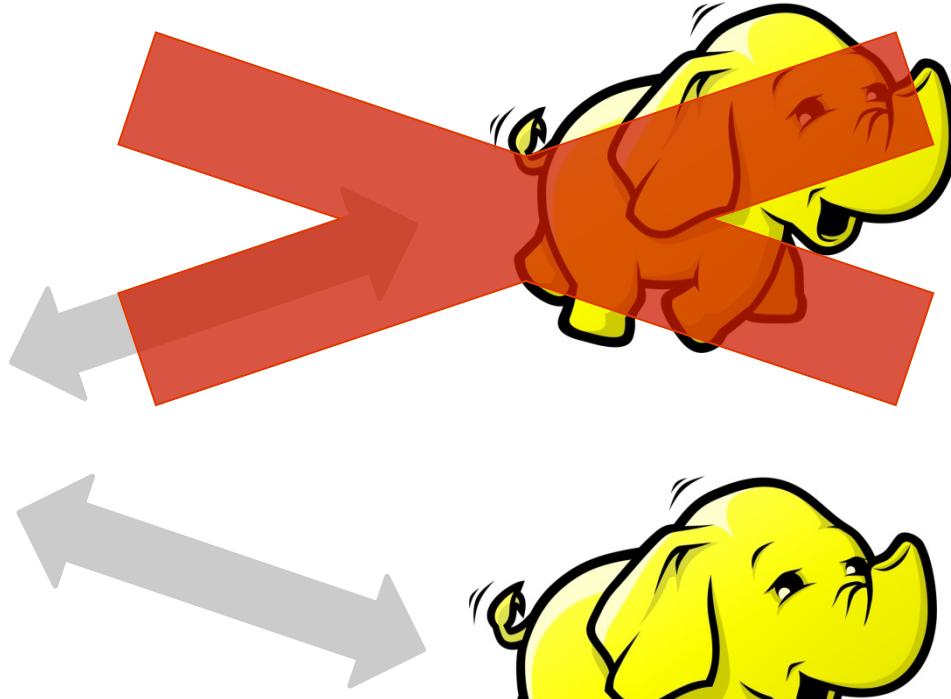
Future Data Pipeline

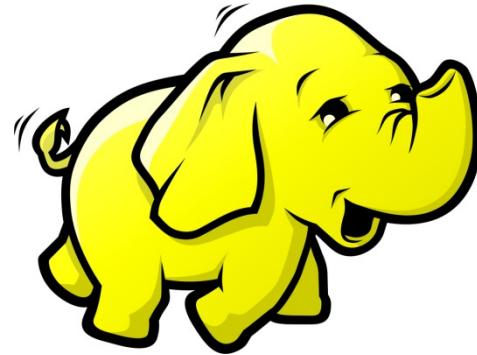
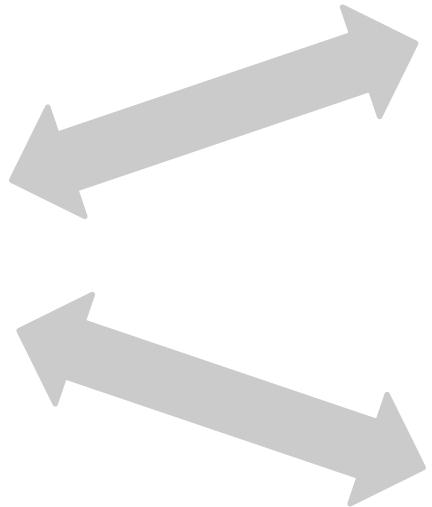




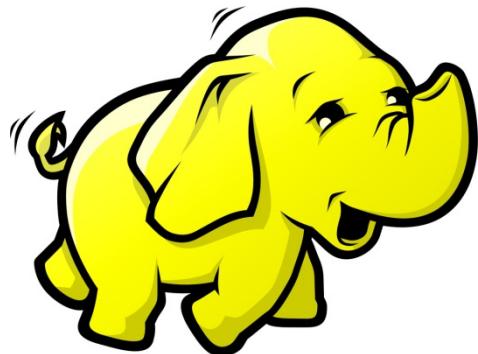








High SLA



Query

inviso



2015-02-09 00:12:04 (UTC -0800)

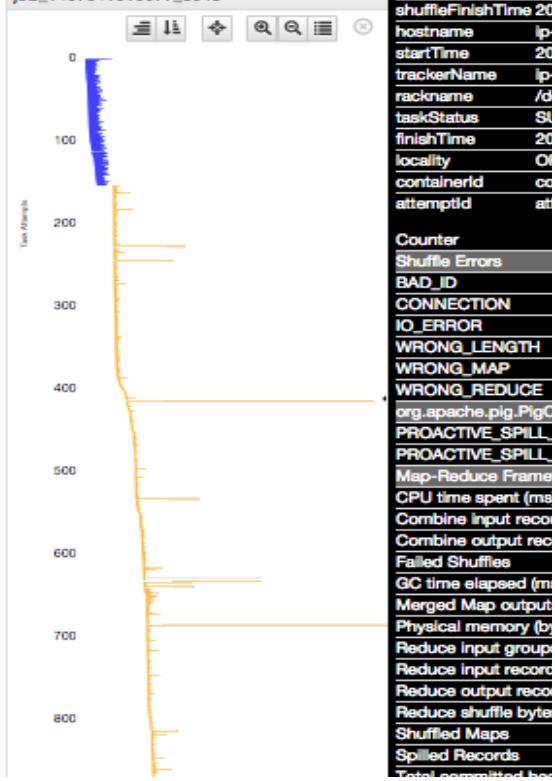
2015-02-09 02:06:24 (UTC -0800)

2015-02-09 04:00:44 (UTC -0800)

2015-02-09 05:55:04 (UTC -0800)

2015-02-09 07:49:24 (UTC -0800)

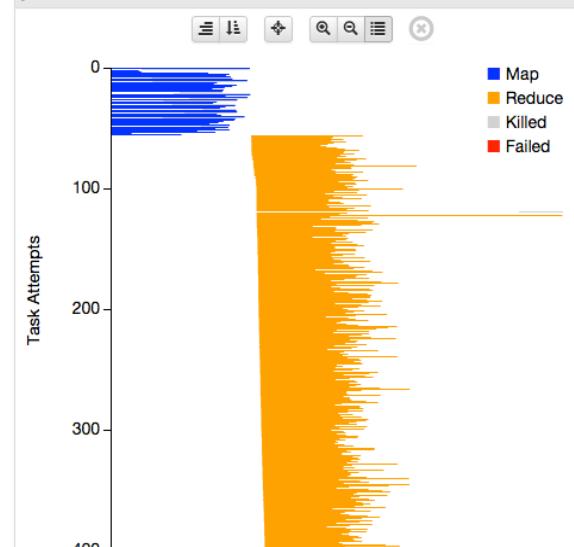
job_1407341916077_3948

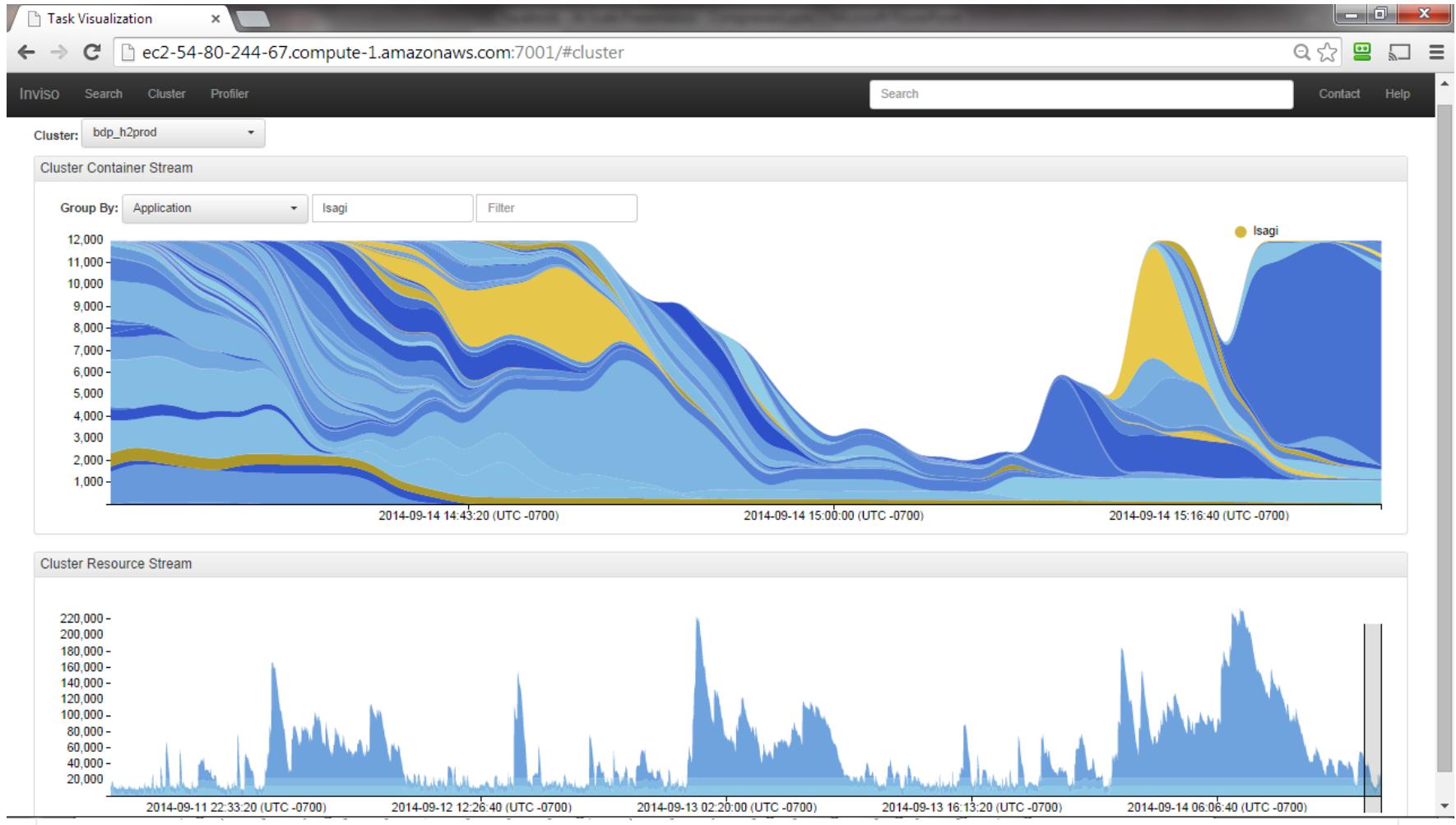


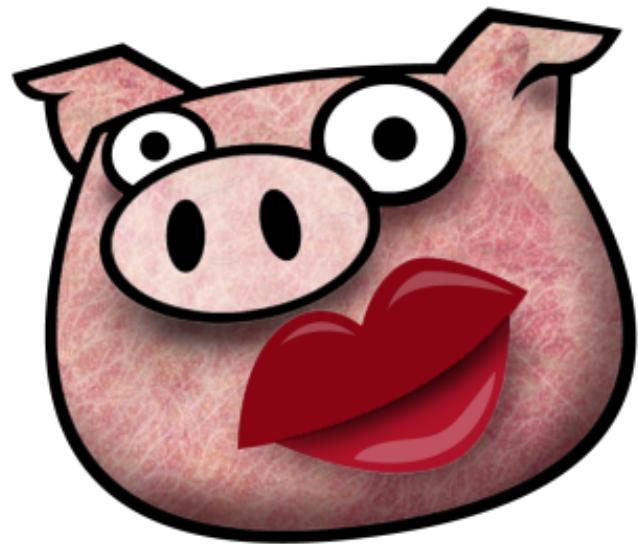
job_1407268826556_5014

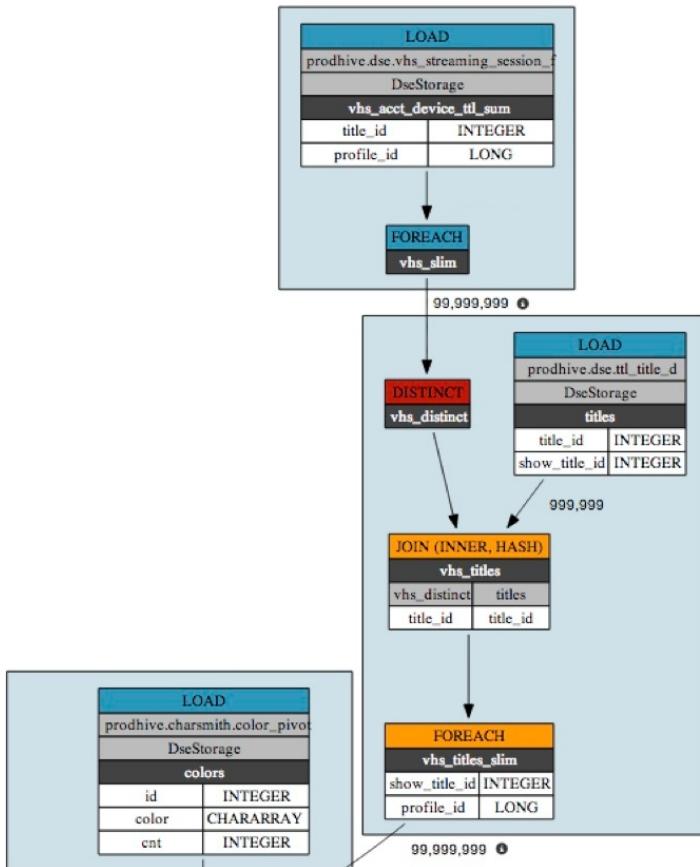


job_1406918339206_88197







**Script Status:**

Status:	finished
Start:	6/9/2013 10:04:39 PM
End:	6/9/2013 10:35:10 PM
Heartbeat:	6/9/2013 10:35:10 PM

Jobs	Map	Reduce
job_201306032155_18654	100%	100%
job_201306032155_18666	100%	100%
job_201306032155_18670	100%	100%
job_201306032155_18674	100%	100%
job_201306032155_18677	100%	100%

FOREACH	
session_p	
account_id	LONG
profile_id	LONG
country_iso_code	CHARARRAY
title_id	INTEGER
package_id	LONG
user_session_id	CHARARRAY
device_type_id	LONG
play_location_id	INTEGER
playtime_source	CHARARRAY
profile_language_locale_cd	CHARARRAY
audio_language_locale_cd	CHARARRAY
text_language_locale_cd	CHARARRAY
profile_language_iso_code	CHARARRAY
audio_language_iso_code	CHARARRAY
text_language_iso_code	CHARARRAY
standard_end_position_sec	INTEGER
standard_sanitized_duration_sec	LONG
browse_sanitized_duration_sec	LONG
total_sanitized_duration_sec	LONG
session_cnt	LONG
view_dateint	LONG
view_hour	INTEGER

↓

JOIN (INNER, REPLICATED)	
session_dt_p	
reg_dt_p	session_p
utc_date	view_dateint
utc_hour	view_hour
country_iso_code	country_iso_code

↓

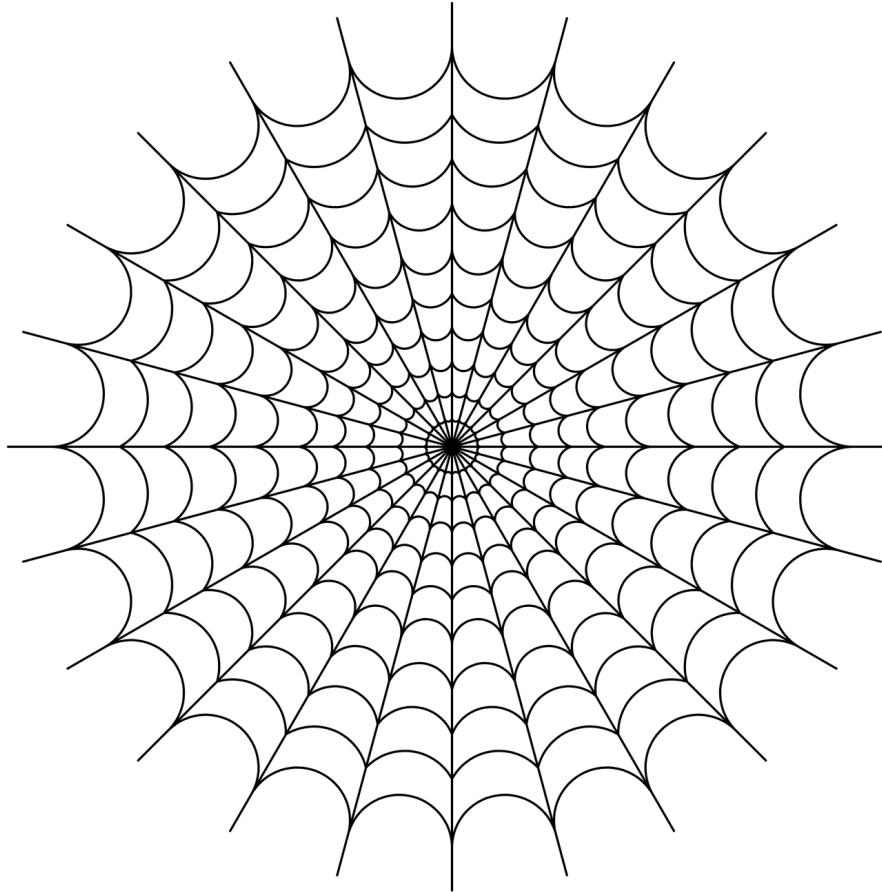
FOREACH	
reg_dt_p	
region_dateint	LONG
utc_date	LONG
utc_hour	INTEGER
country_iso_code	CHARARRAY

28,494,000 ⓘ

dt::calendar_date >= window_ttl_f::window_start

Metacat

Charlotte





Dependencies

Column: --ALL-- Type: --ALL-- Username: --ALL--

Upstream (1)

Name	Username	Type	Frequency
PigLatin:vhs_acct_device_ttl_sum.pig		etl	53

Downstream (229)

Name	Username	Type	Frequency
INSERT OVERWRITE TABLE vhs_acct_device_ttl_sum AS SELECT * FROM vhs_acct_device_ttl_sum WHERE date = '2015-02-24' AND c_type = 'TV'		adhoc	3228
INSERT OVERWRITE TABLE vhs_acct_device_ttl_sum AS SELECT * FROM vhs_acct_device_ttl_sum WHERE date = '2015-02-24' AND c_type = 'PC'		adhoc	1522
INSERT OVERWRITE TABLE vhs_acct_device_ttl_sum AS SELECT * FROM vhs_acct_device_ttl_sum WHERE date = '2015-02-24' AND c_type = 'Mobile'		adhoc	1414
PigLatin:PVR Show Pop - 20150204		etl	582
STING.Kids Originals Title SOP SOD		sting	575
CREATE TABLE mjarley.title_streaming_...null(Stage-30)		adhoc	574
CREATE TABLE iyohai.title_streaming_e...null(Stage-15)		adhoc	437

Big Data Portal

Netflix Big Data Portal [Close]

bigdataportal.dynprod.netflix.net:7001/#tables/prodhive/dse/vhs_acct_device_ttl_sum

NETFLIX Big Data Portal

kubrown

Home - Query

Dashboard

Schema Browser

S3 Browser

Quick Schema Browser ?

No tags applied Tags

prodhive/dse/vhs_acct_device_ttl_sum Star

[/prodhive/dse/vhs_acct_device_ttl_sum](#) Star

Location: prodhive/dse/vhs_acct_device_ttl_sum Tags: No tags applied Tags

Schema Details Dependencies Notes

Table Info

Tags core Tags

Location s3n://netflix-dataoven-prod-users/hive/warehouse/dse.db/vhs_acct_device_ttl_sum.gz

Partition Keys region_dateint

Created By [Redacted]

Created Date 5/23/2012 1:55:40 PM

Last Modified By [Redacted]

Modified Date 3/4/2013 3:18:40 PM

Lifetime Not set Set Lifetime

View in Franklin [Franklin icon]

Columns:

country_iso_code	Flame icon
account_id	Flame icon

NETFLIX Big Data Portal

bigdataportal.dynprod.netflix.net:7001

Contact Us What's New?

PRESTO Untitled Query

```
1 select * from dse.geo_country_d
```

Presto prohive RUN (ctrl-enter or F5) Save Clear Show Options

Query History Data Viewer

Recent Queries

location

filter

prodhive

Recent Queries ①

Status Filter: All ▾

Keyword Filter:

⌚ select * from dse.location_d 10:29

Query Text:
select * from dse.location_d

Parameters:

Status:	RUNNING	KILL JOB
Job Type:	PrestoJob	
Genie Job Id:	48991974-b865-11e4-96e1-0ab5e51a10cb ⓘ	
Start Time:	2/19/2015 10:29:47 AM	
Update Time:	2/19/2015 10:29:48 AM	

Links:

Saved Queries ①

Owned By You

	Chukwa Synthetic Events Compare 1	
	kurttest	

Shared With You

	presto demo	
	geo_country_d	

Public

	Query Joining map command audit and request audit tables	
--	--	--

NETFLIX Big Data Portal

bigdataportal.dynprod.netflix.net:7001

Recent Queries

Documentation

API

To reattach to this job using [Kragle](#):

```
import kragle as kg

job = kg.genie.reattach_job('9a3b8aacfea9-11e4-8201-12f997ac03a5')
```

To fetch the first 100 rows of data from genie stdout:

```
kg.transport.Transporter().source(kg.genie.reattach_job('9a3b8aacfea9-11e4-8201-12f997ac03a5')).execute().r
ow_map(limit=100)
```

To resubmit this query as a new job:

```
job = kg.genie.HiveJob().script('''select * from dse.geo_country_d''')

running_job = job.execute()
```

Big Data API

(aka Kragle)

IP[y]: Notebook

Untitled2 (unsaved changes)

File Edit View Insert Cell Kernel Help



Cell Toolbar: None

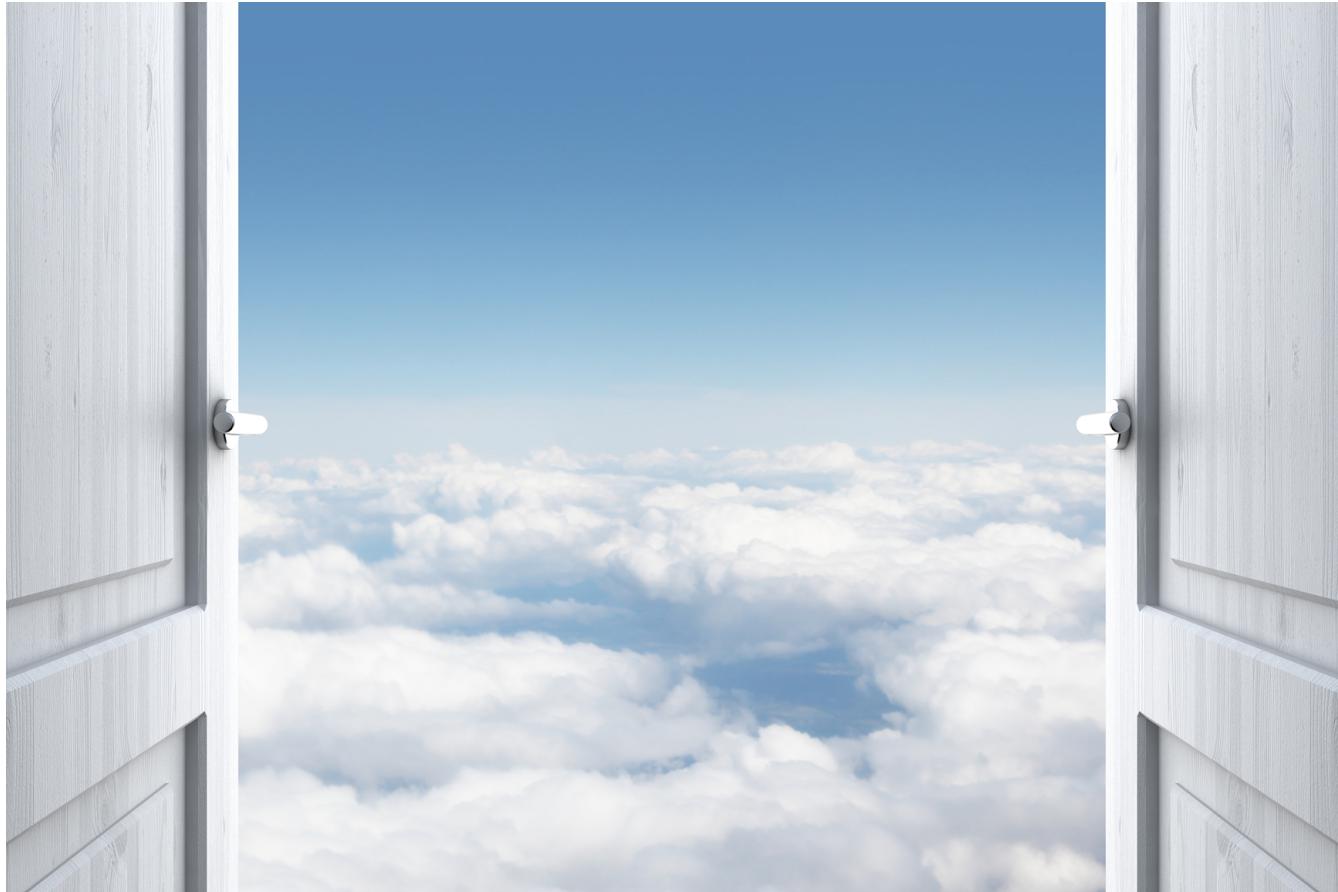
In [1]: `import krangle as kg`In [2]: `job = kg.genie.PrestoJob() \
 .script('select * from dse.geo_country_d') \
 .headers() \
 .execute()`In [3]: `job.pandas()`

Out[3]:

	country_iso_code	country_sk	country_desc	country_full_desc	company_code	subregion_sk	subregion_pd_sk	content_subregion_sk	content_pd_sk
0	FR	8728	France	France	65	14	14	15	France
1	CM	8700	Cameroon	Cameroon	-1	9	3	9	Other
2	AR	6582	Argentina	Argentina	33	8	4	8	Latin A
3	MD	8788	Moldova	Moldova, Republic of	-1	9	3	9	Other

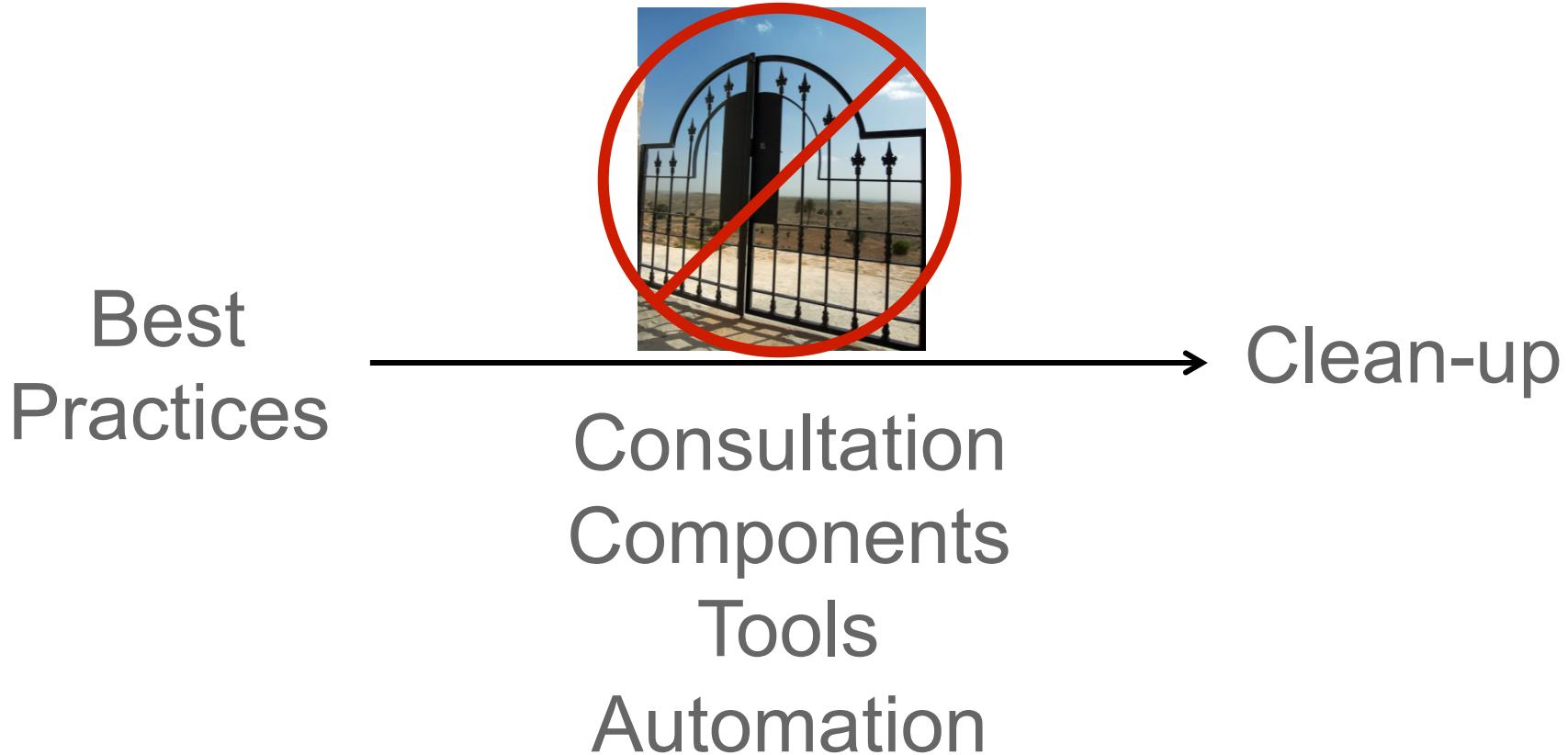
All about the tech?

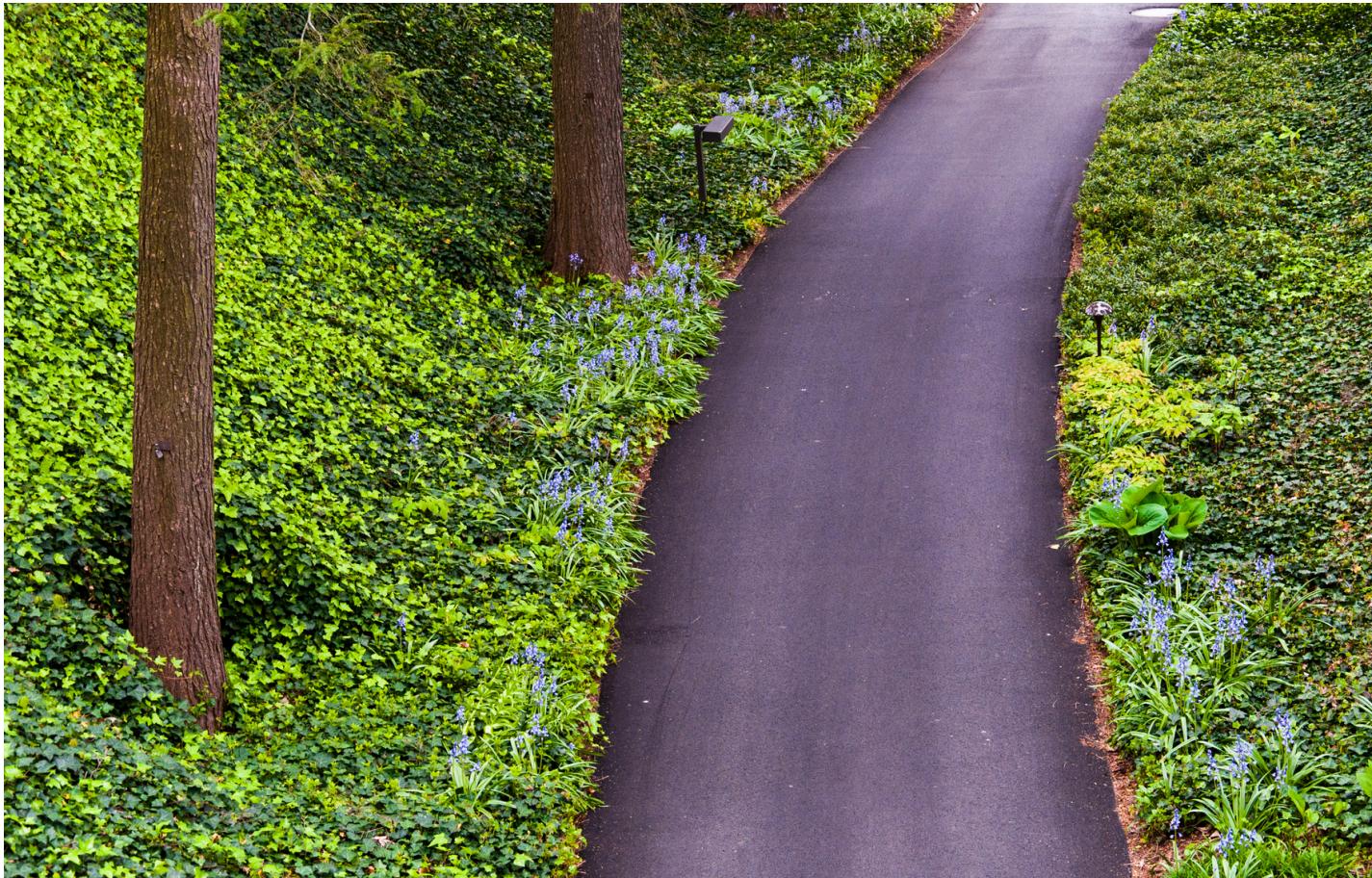






Development & Deployment Flow









Why?

Netflix Culture Deck



Google Search: “Netflix Culture Deck”

Questions?

NETFLIX

<http://jobs.netflix.com/>
kurtbrown@netflix.com