NETFLIX



Beaming Flink to the Cloud @ Netflix

Monal Daxini

Engineering Manager Stream Processing, Real Time Data Infrastructure



World's Leading Internet Streaming Service

(Global launch Jan 6, 2016)



Netflix Service Scale

- 83+ Million Members, 190+ Countries
- **1000**+ device types
- 35% of downstream Internet traffic

Netflix Service Scale - Daily viewing hours

125,000,000,000+

Whoa!



1,000,000,000,000+

1.4 PB

That's a huge number!

Event Scale

Peak

- 1T unique events ingested/day
- **16M** / sec
- **43GB** / sec
- 10MB / message

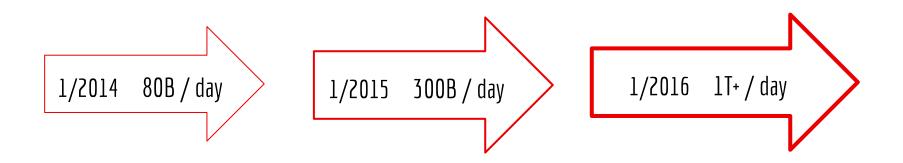
Daily Averages

- 1T+ events processed
- 600B unique events ingested
- 1.4 PB / day
- 4K / event

Keystone Scale

99.99% + Availability / Four 9s

Keystone Events Trend



Evolution

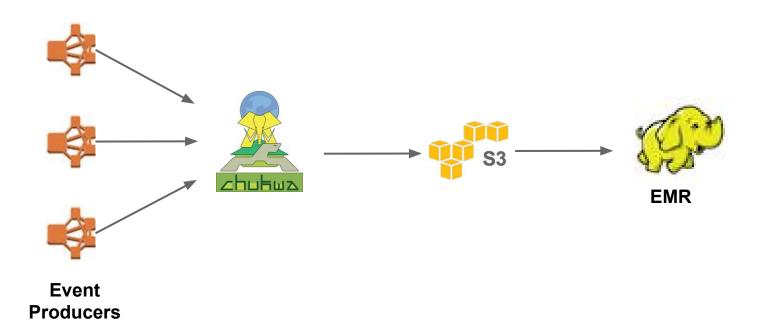


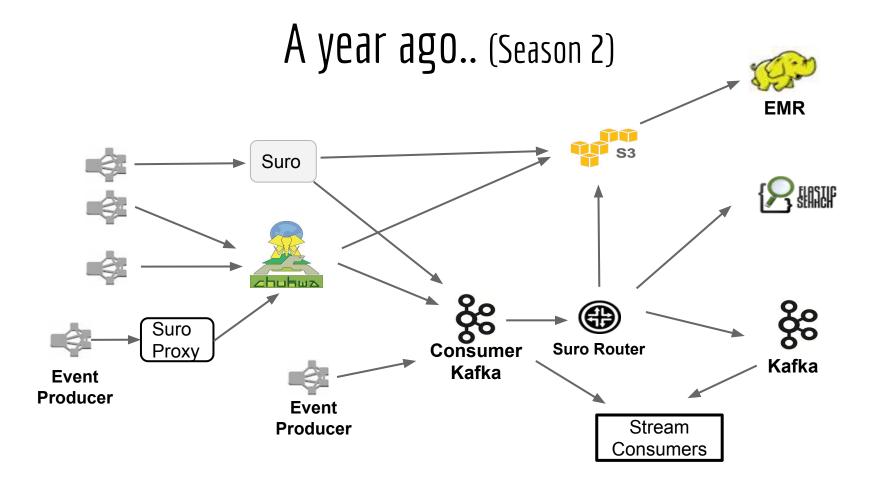
SPaaS

Goal - Migrate 1.3 PB of event data to a new Pipeline in flight, while ensuring data diff < 0.1%



Few years ago (Season 1)

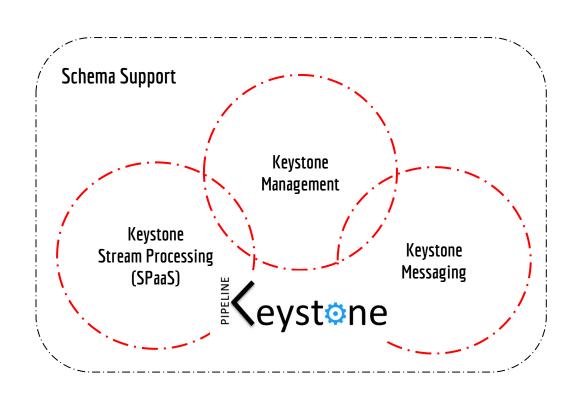




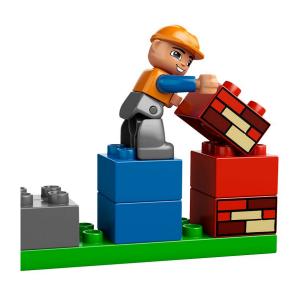
Season 3



Keystone 100% in AWS

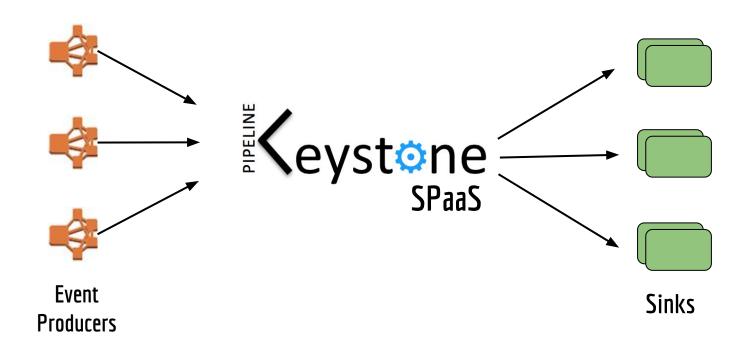


Our Philosophy

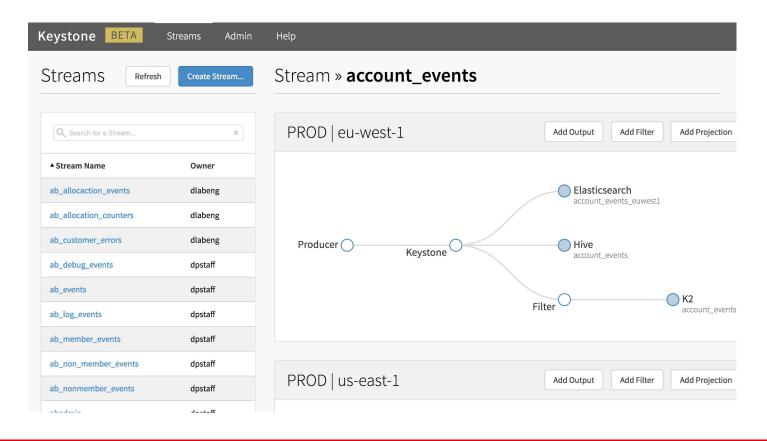


Create Duplo® Blocks: Let reusability drive new value

Q4 2015



Keystone Management

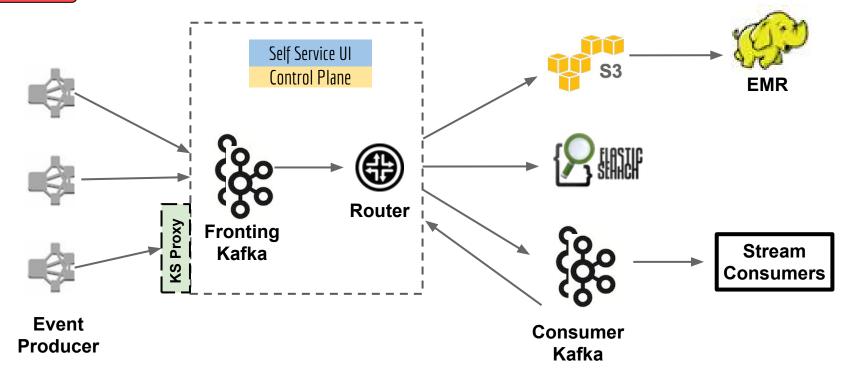


Per Stream Auto Dashboard

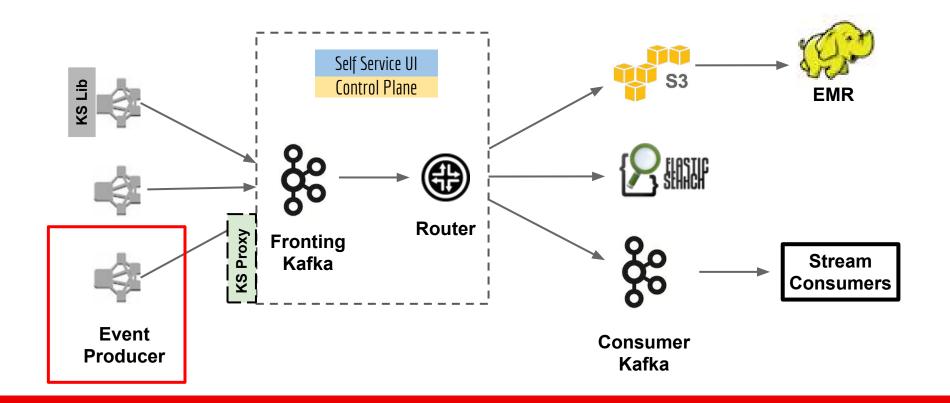


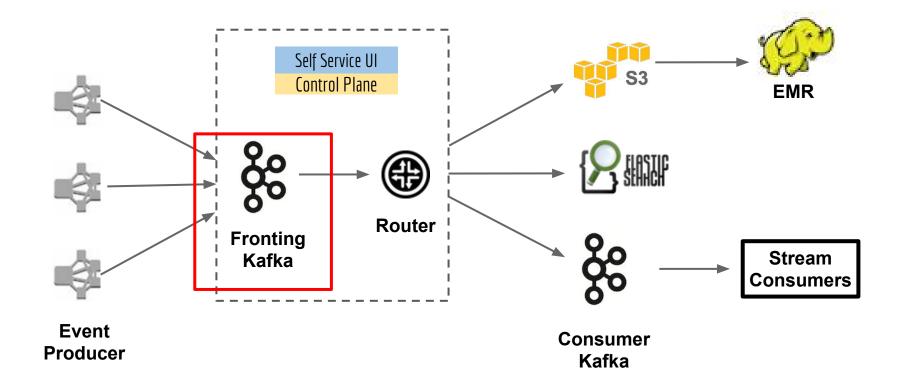
100% in AWS 24 x 7 Region failover

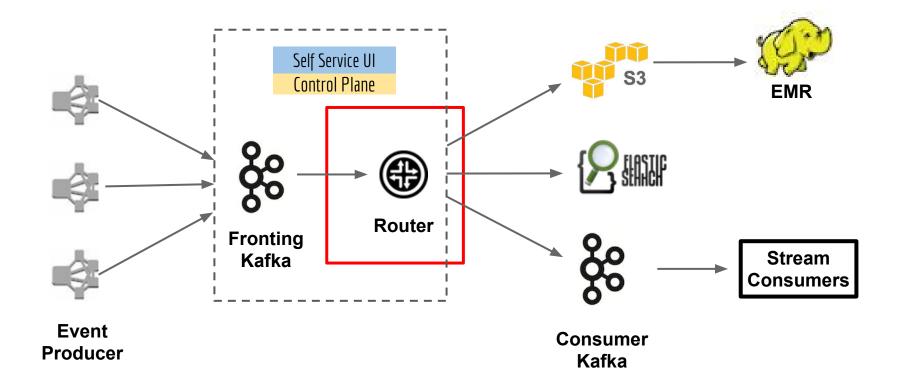
Keystone SPaaS

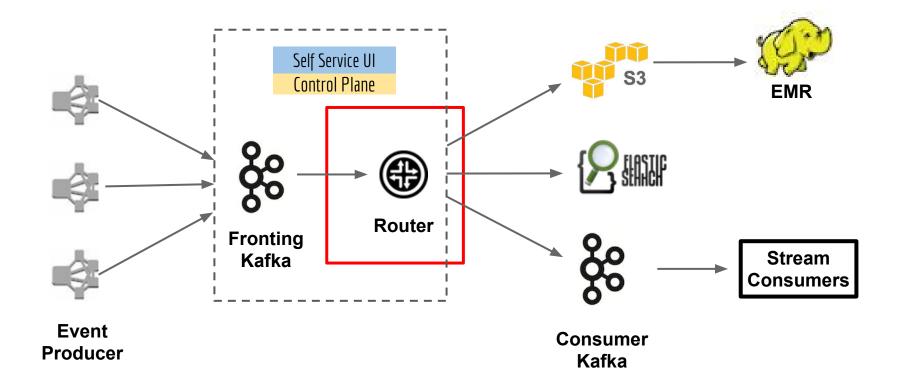


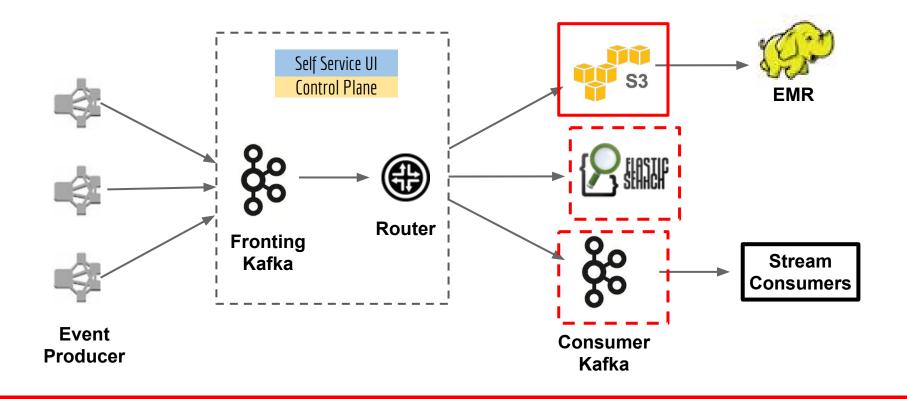
Event flow Keystone Pipeline As a Service

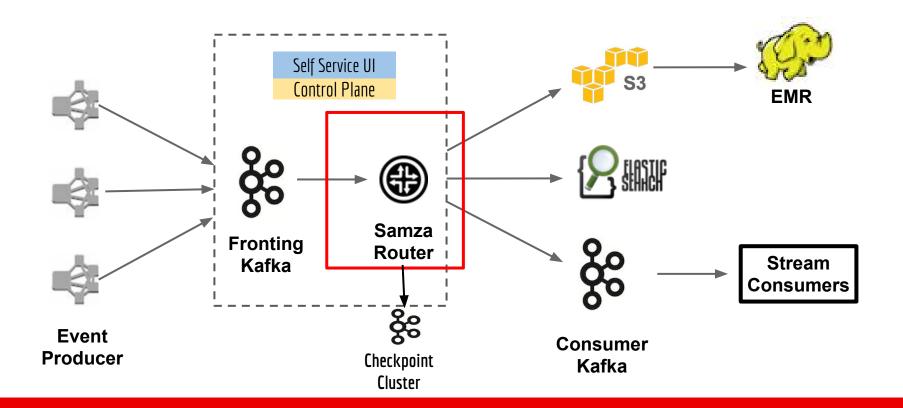










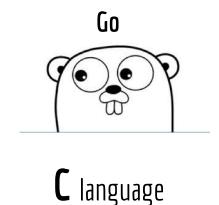


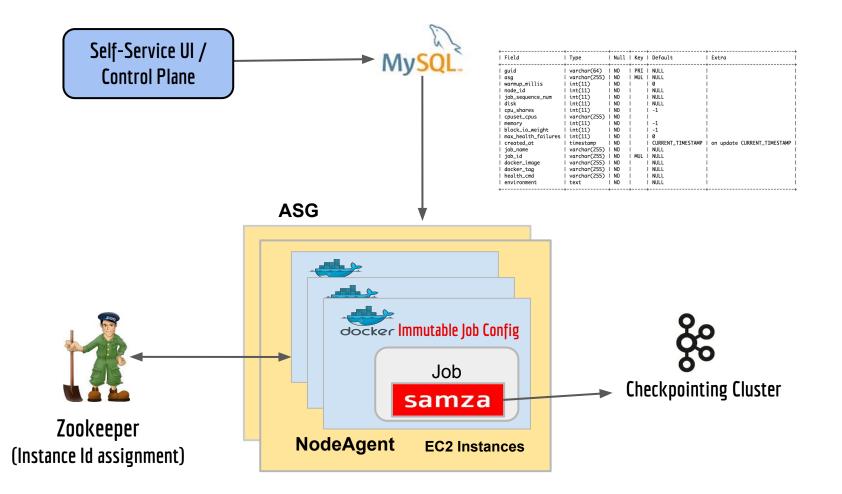
Details..

- Massively parallel use-case
 - Per element processing declarative filtering & projection
- Stateless except Kafka offset checkpointing state

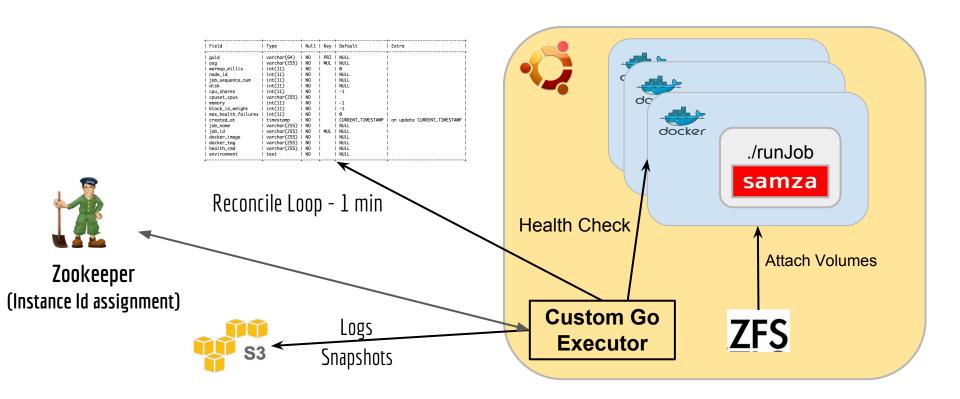
Routing Infrastructure







What's running on the host?



Yes! You inferred right!

No Mesos & No Yarn

Keystone Salient Features

Keystone Scale

- 4000+ Kafka brokers
- 14000+ Samza Jobs
 - Running in docker containers
 - On 1600+ nodes

Keystone Salient Features 100% in AWS

- At-least-once delivery semantics*
- Multi-Tenant
- Self Service
- Scalable
- Fault Tolerant

Keystone Salient Features

- Event payload is immutable
- Inject event metadata
 - o guid, timestamp, host, app
- Custom extensible wire protocol
- Kafka producer wrapper

Keystone Salient Features

- Kafka Cluster failover
 - Kafka Kong
- Routing regional / failover
- Scales based on historical traffic
 - Externally driven

Season 4 Plot & Pilot...



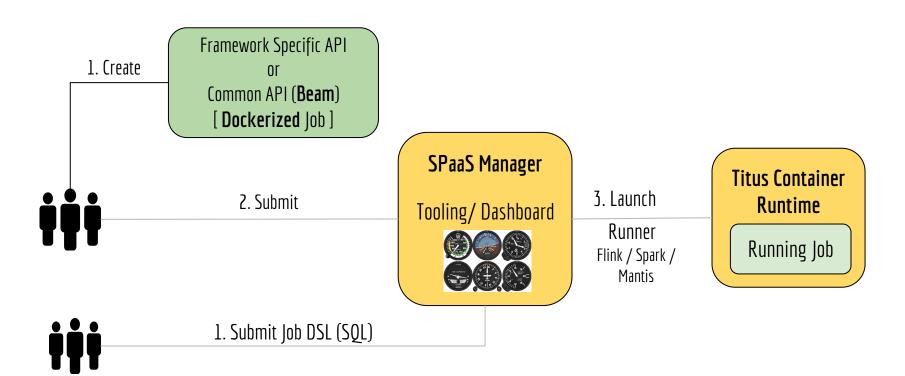
Stream **P**rocessing **A**s a **S**ervice (SPaaS)

(more capable)

SPaaS Vision (plot)

- Multi-tenant support for stateful stream processing apps
- Autoscaling managed infrastructure
- Support for schemas
- Self Service Tooling

SPaaS Architecture (plot)



SPaaS - "Beam Me Up, Scotty!"

Why Apache Beam?

- Portable API layer to build sophisticated data processing apps
 - Support multiple execution engines
- Unified model API over bounded and unbounded data sources
- Millwheel, FlumeJava, <u>Dataflow model</u> lineage

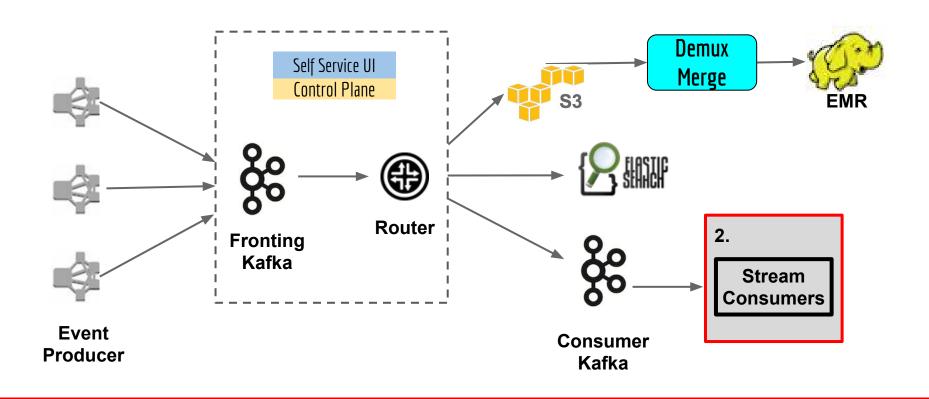
SPaaS - Pilot

Iterative build out: then

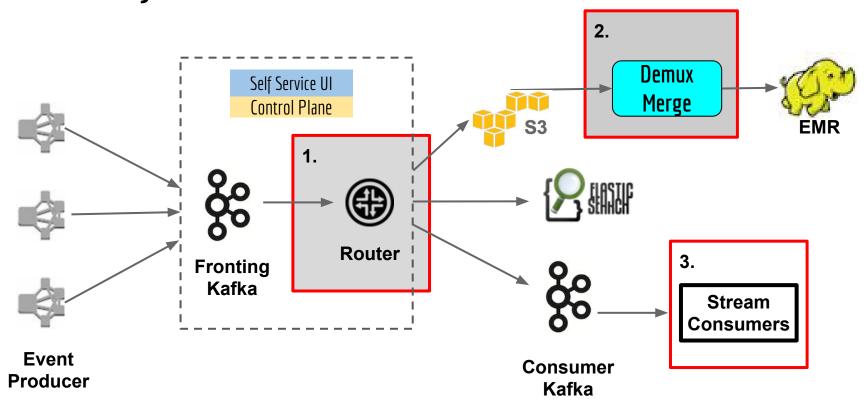


- First Flink on Titus in VPC, AWS
 - **Titus** is a cloud runtime platform for container based jobs
- Next Apache Beam and Flink runner

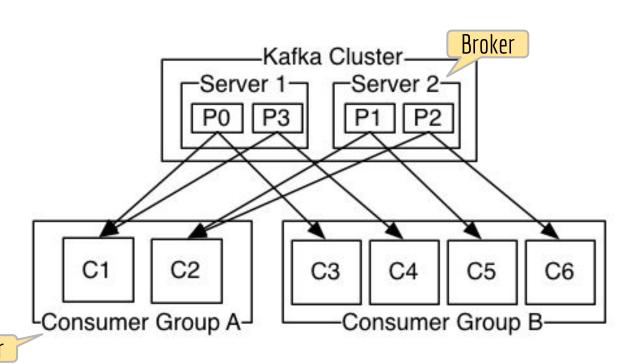
Keystone SPaaS-Flink Pilot Use Cases



Keystone SPaaS-Flink Pilot Use Cases

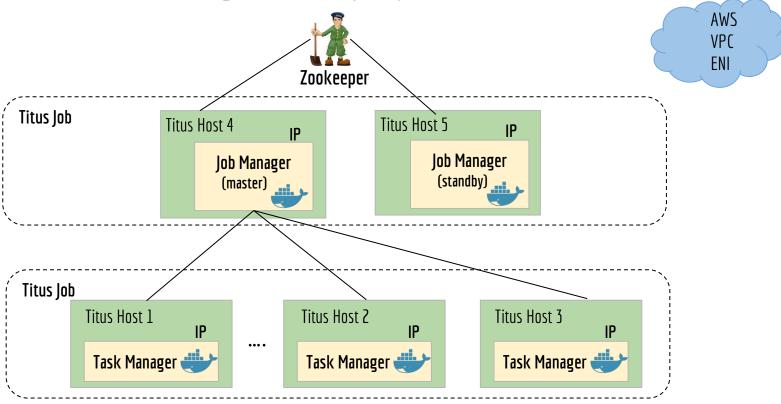


Router - Massively parallel use case

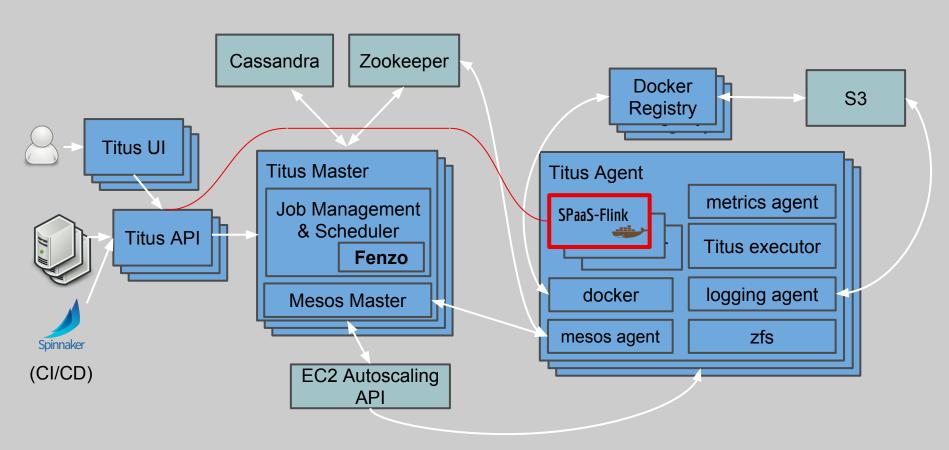


Router

Flink (1.2) Program Deployment (prod shadow)



Titus High Level Architecture



Titus Current Tasks Job Status Cluster

Tasks Submit Job

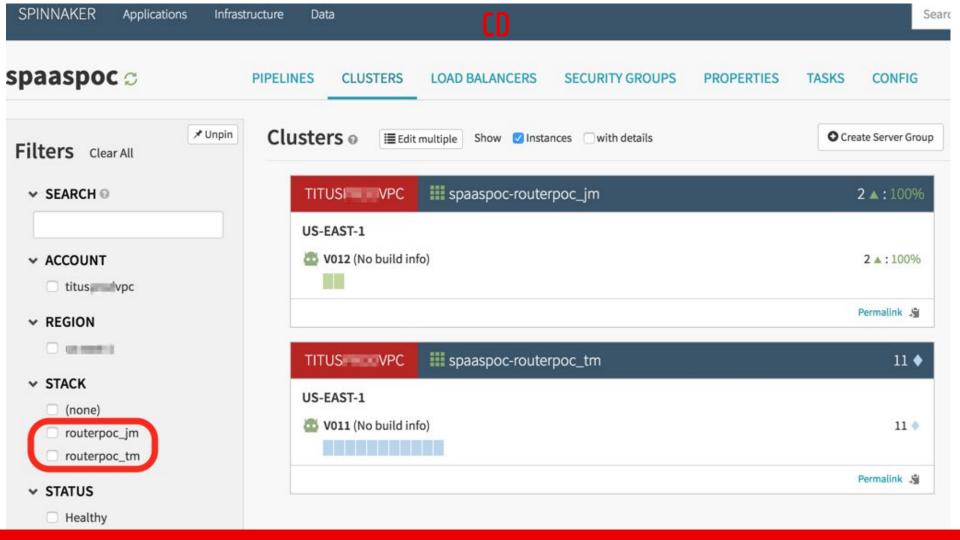
Titus UI

Running

Queued

Dispatched

Show 25 e	ntries					Search:	spaas-router
ID A	Image	⊕ Command	\$	SubmittedAt	Jobld	*	Actions
Titus-240138- worker-0-2	spaas-routerpoc	/apps/spaas/spaasbin/prep- jobmanager.sh	1	2016-09-07T21:43:59.916Z	Titus-240138	3	×
Titus-240138- worker-1-3	spaas-routerpoc	/apps/spaas/spaasbin/prep- jobmanager.sh		2016-09-07T21:43:59.916Z	Titus-240138	:	×
Titus-240139- worker-0-2	spaas-routerpoc	/apps/spaas/spaasbin/prep- taskmanager.sh		2016-09-07T21:44:07.641Z	Titus-240139		×
Titus-240139- worker-1-3	spaas-routerpoc	/apps/spaas/spaasbin/prep- taskmanager.sh		2016-09-07T21:44:07.641Z	Titus-240139		×
Titus-240139- worker-10-12	spaas-routerpoc	/apps/spaas/spaasbin/prep- taskmanager.sh		2016-09-07T21:44:07.641Z	Titus-240139		×
Titus-240139- worker-2-4	spaas-routerpoc	/apps/spaas/spaasbin/prep- taskmanager.sh		2016-09-07T21:44:07.641Z	Titus-240139		×
Titus-240139- worker-3-5	spaas-routerpoc	/apps/spaas/spaasbin/prep- taskmanager.sh		2016-09-07T21:44:07.641Z	Titus-240139		×
Titus-240139- worker-4-6	spaas-routerpoc	/apps/spaas/spaasbin/prep- taskmanager.sh		2016-09-07T21:44:07.641Z	Titus-240139	;	×
Titus-240139- worker-5-7	spaas-routerpoc	/apps/spaas/spaasbin/prep- taskmanager.sh		2016-09-07T21:44:07.641Z	Titus-240139	;	×
Titus-240139- worker-6-8	spaas-routerpoc	/apps/spaas/spaasbin/prep- taskmanager.sh		2016-09-07T21:44:07.641Z	Titus-240139	;	×
Titus-240139- worker-7-9	spaas-routerpoc	/apps/spaas/spaasbin/prep- taskmanager.sh		2016-09-07T21:44:07.641Z	Titus-240139		×
Titus-240139- worker-8-10	spaas-routerpoc	/apps/spaas/spaasbin/prep- taskmanager.sh		2016-09-07T21:44:07.641Z	Titus-240139		×
Titus-240139- worker-9-13	spaas-routerpoc	/apps/spaas/spaasbin/prep- taskmanager.sh		2016-09-07T21:44:07.641Z	Titus-240139		ĸ



Dashboard

Last hour

End Minus 5 mins \$ mins

Shift None

Step Auto

Time Zone US/Pacific \$

Logarithmic

Flink Router POC

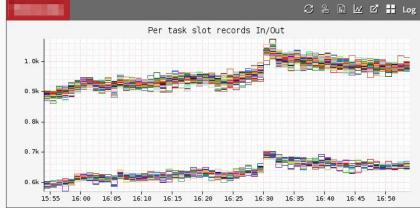
C Refresh All

Auto Refresh

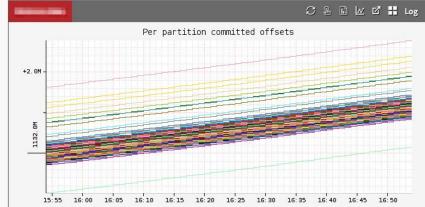
Show Legend

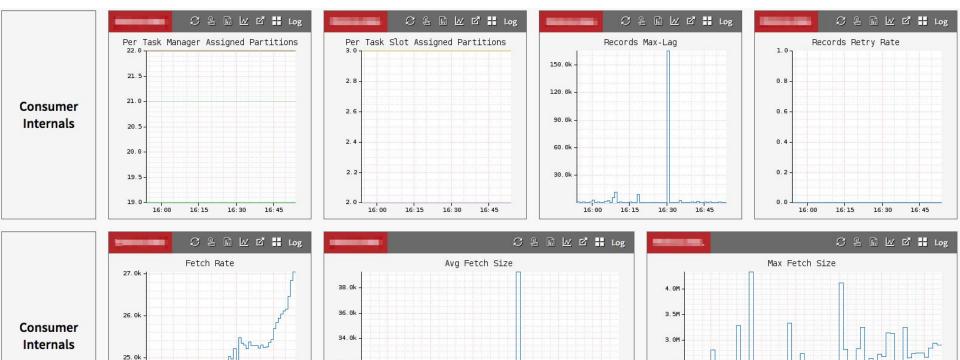












16:30

16:40

16:50

2.5M

2. 0M

1.5M-

16:00

16:10

16:20

16:30

16:40

16:50

32.0k

30.0k

28. 0k

16:00

16:10

16:20

24. 0k

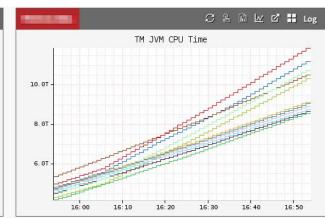
16:15

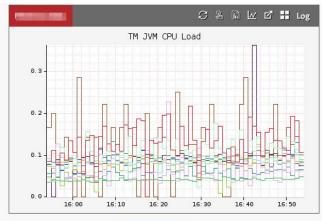
16:00

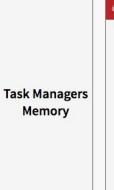
16:45

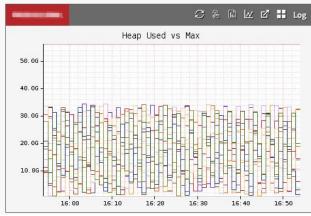
16:30

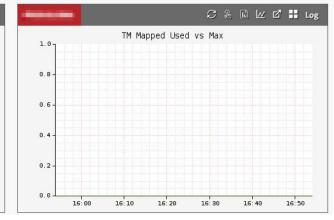


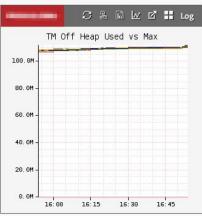


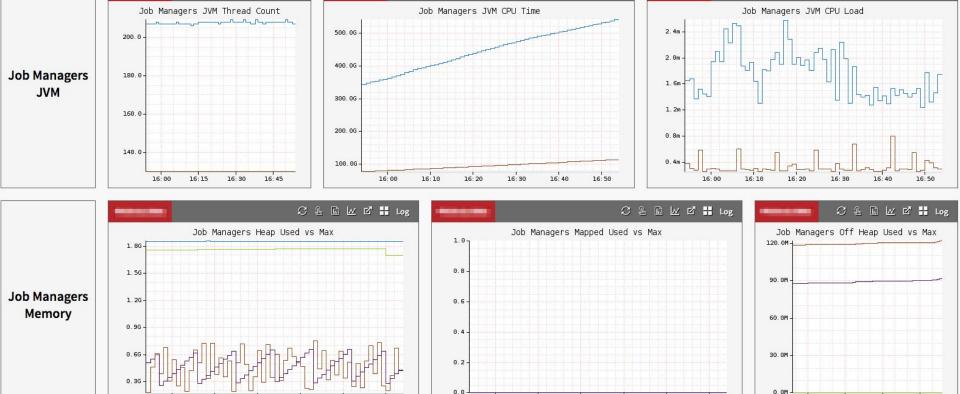












16:00

16:10

16:30

16:40

16:50

16:00

16:15

16:30

16:45

16:50

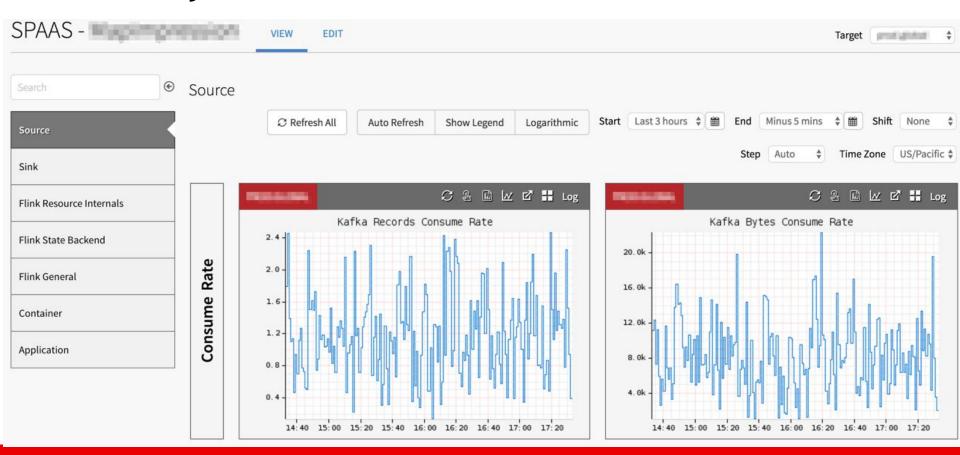
16:40

C & W W B Log

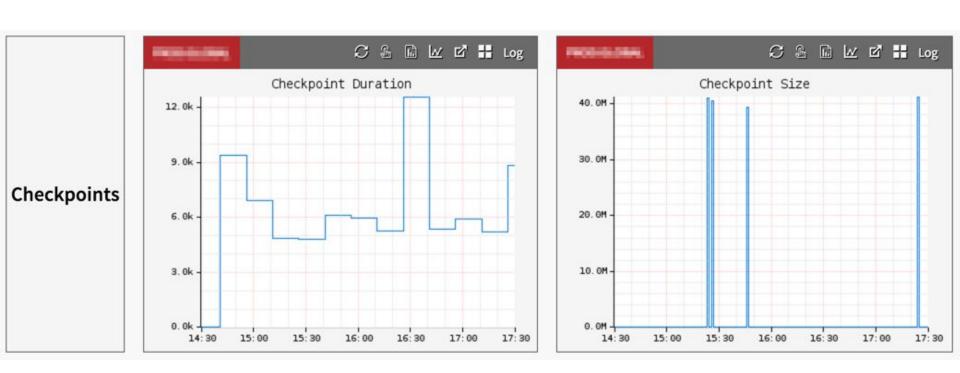
C & L W I H Log

16:10

Keystone SPaaS-Flink Pilot Use Case - 2



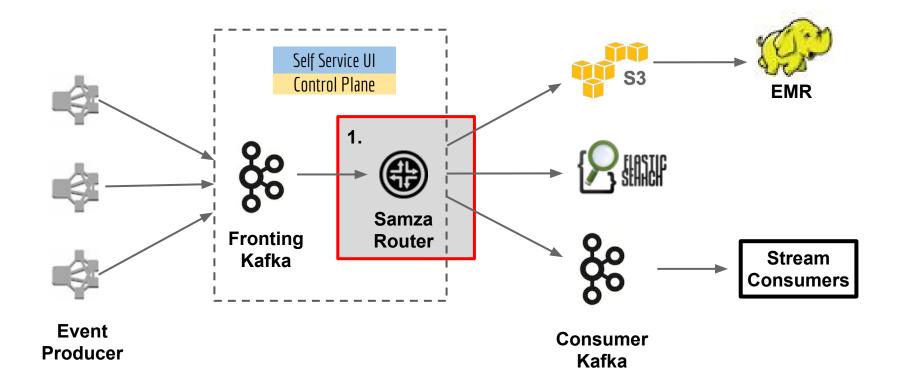
Keystone SPaaS-Flink Pilot Use Case - 2



Flink Router perf test (YMMV)

- Note
 - The tests were performed on a specific use case,
 - running in a specific environment, and with
 - with **one specific event stream**, and setup.

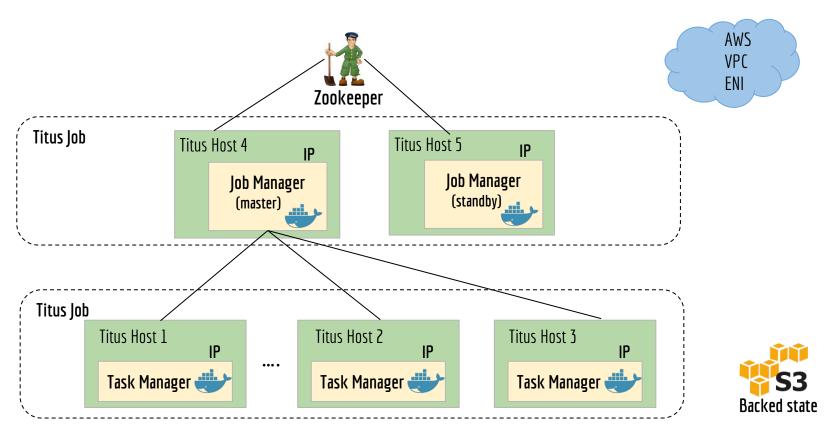
Keystone



Details..

- Different runtimes for Flink & Samza routers
- Massively parallel use-case
 - Per element processing
- Focused on net outcomes

Flink (1.2) Router

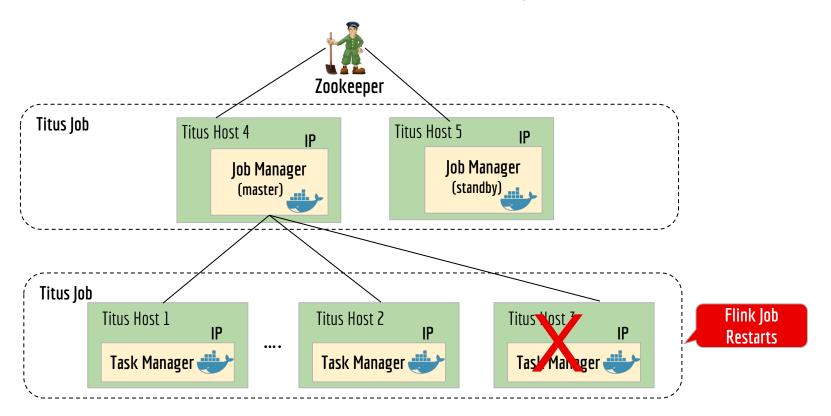


Flink Router outcome (YMMV)

- Cost ≈17% savings
- Memory utilization ≅16% better
- Cpu utilization ≈ 40% better
- Network utilization ≅ 10% better
- Msg. throughput $\approx 1\%$ (avg) 4% (peak) better

Awaiting Flink Features

Fine Grained Recovery FLIP-1



Awaiting Flink Features (now)

- Checkpoints and savepoints
 - FLINK-4484 Unify, Persistent checkpoints, periodic savepoints
 - Compatible across Flink version upgrades
 - Inspection tool for debugging

Awaiting Flink Features (now)

- Atomic savepoint and stop (pause) the job
- Dynamic Scaling DONE
 - Resume from savepoint with different parallelism
 - Job elasticity
- FLINK-4545 Adjust TaskManager network buffer when scaling

Awaiting Flink Features (now)

- Cluster Runtime and Elasticity (FLIP 6)
- Extending Window Function Metadata (FLIP-2)
- Large State support
 - Incremental checkpointing
 - hot standby

Awaiting Flink Features

- Metric tags as key-value pairs FLINK-4245
- FLIP-9 Window Trigger DSL DONE
- FLIP-11 Table API
- Side Inputs / Side Outputs (handle late data)



Ponder over

Could Stream Processing Engine enable building non-analytics Applications?

More brain food...

- Netflix Keystone Pipeline Evolution
- Netflix Kafka in Keystone Pipeline
- Samza Meetup Presentation
- <u>Titus talk</u>
- Netflix OSS