

Developers

Realtime data pipelines with
Apache Pulsar™ and Apache Cassandra

Data in motion with Pulsar Functions



Director of Developer Relations



- Trainer
- Public Speaker
- Developers Support
- Developer Applications
- Developer Tooling

- Creator of ff4j (ff4j.org)
- Maintainer for 8 years+

- Happy developer for 14 years
- Spring Petclinic Reactive & Starters
- Implementing APIs for 8 years



Cédrick Lunven

Streaming Developer Advocate



ddieruf



ddieruf



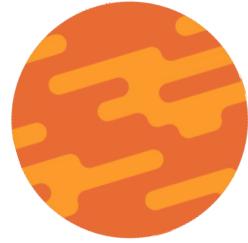
@dierufdavid



- Developer/Architect
- Specialize in cloud-native
- Kubernetes
- Application modernization
- CI/CD



David Dieruf



S



Cedrick
Lunven



David
Dieruf



Rags
Srinivas



Artem
Chebotko



Stefano
Lottini



Aleksandr
Volochnev



Aaron
Ploetz



S



Jack
Fryer



Kirsten
Hunter



Gary
Harvey



Mary
Grygleski



Ryan
Welford



David
Gilardi

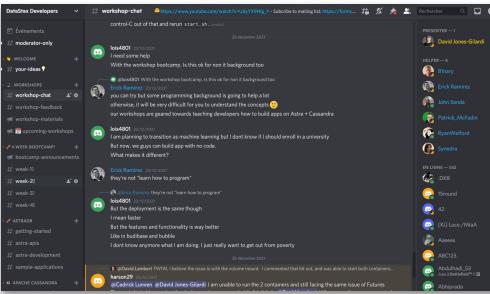


DataStax Developers Crew

Livestream: youtube.com/DataStaxDevs

Questions: <https://dtsx.io/discord>

Agenda



YouTube
(with nighbot)

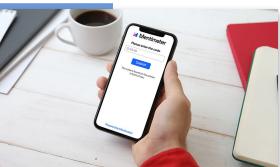
Discord
(#workshop-chat)



!discord

Games and quizzes: menti.com

How much experience do you have with the Spring Framework ?



DataStax Developers



Mentimeter

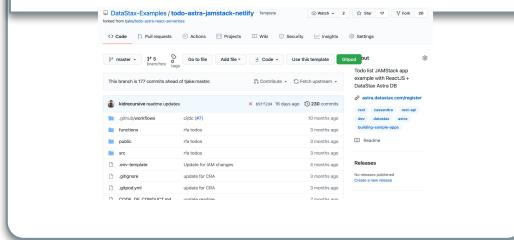


DataStax Developers

Live Sessions

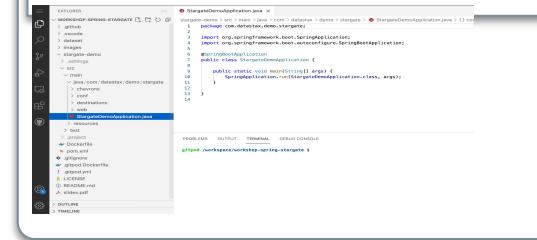
Nothing to install !

Source code + exercises + slides



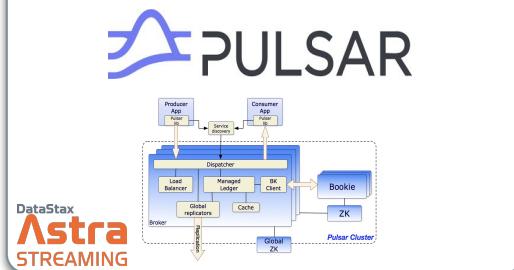
 GitHub

IDE

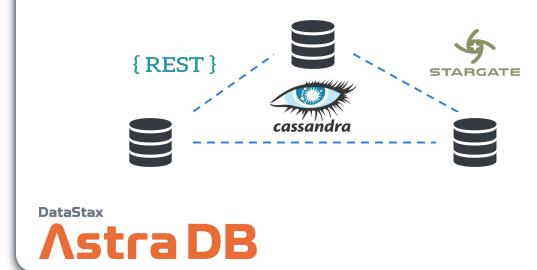


 Gitpod

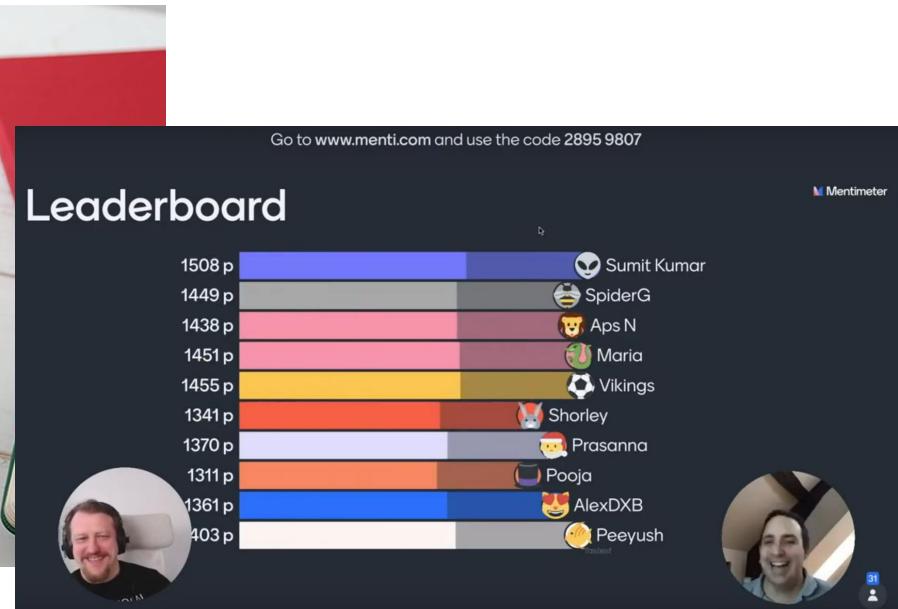
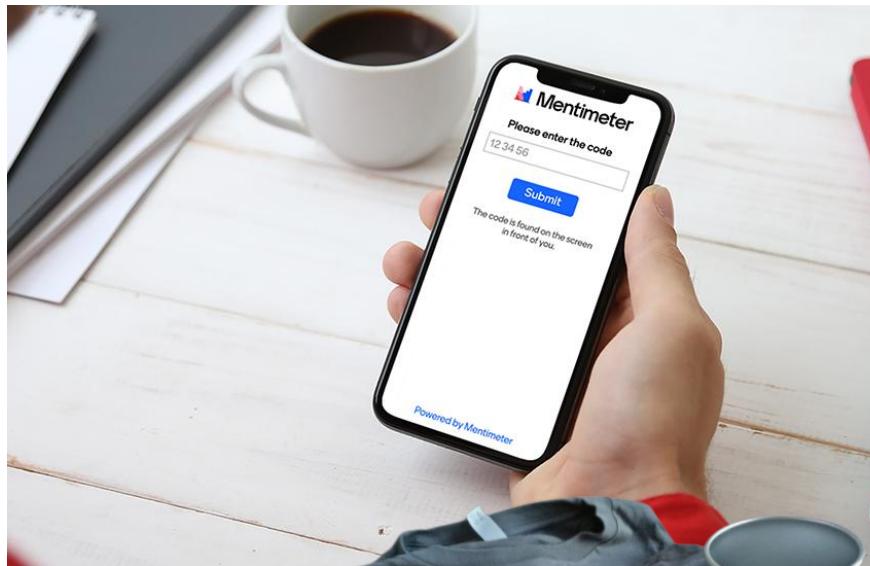
Messaging system (based on Pulsar)



Database + GraphQL + PlayGround



Hands-On Housekeeping



**menti . com ⇒ enter code
Don't answer in YT chat
Look at phone (not at YT)**

Quiz on "Menti" !

01



Introduction to
Apache Pulsar™

02



DataScience in
Event Streaming

03



Pulsar Functions

04



Introduction to
Apache Cassandra™

05



Pulsar Sinks

06



What's next?
Quiz, Homework, Next week



Agenda

01



Introduction to
Apache Pulsar™

02



DataScience in
Event Streaming

03



Pulsar Functions

04



Introduction to
Apache Cassandra™

05



Pulsar Sinks

06



What's next?
Quiz, Homework, Next week



Agenda

Event Streaming == Message Streaming

- Watch for events with “the system” or application
- Publish messages and receive events
- Make decisions on data in real time
- Ingest high frequency of messages with very low latency and consume at a different rate



Event Streaming == Message Streaming



Open source

Created by Yahoo

Contributed to the Apache Software Foundation 2016

Top-level project 2018

Cloud-native design

Cluster based

Multi-tenant

Simple client APIs (Java, C#, Python, Go, Node, ...)

Separate compute and storage!

Guaranteed message delivery

If a message successfully reaches a Pulsar broker, it will be delivered to its intended target.

Light-weight serverless functions framework

Create complex processing logic within a Pulsar cluster (aka: data pipeline)

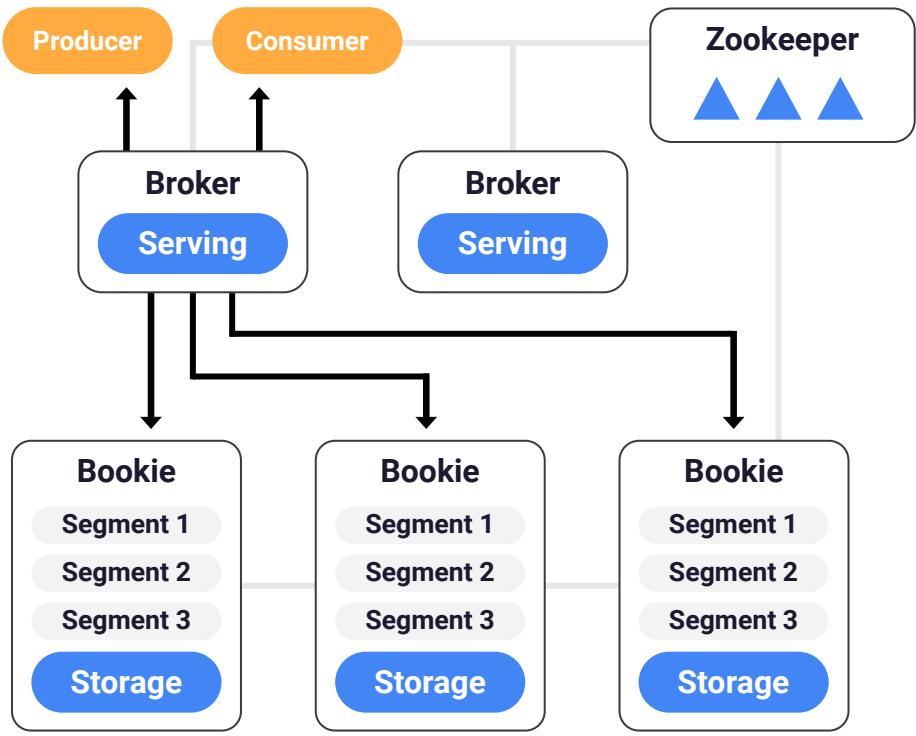
Tiered storage offloads

Offload data from hot/warm storage to cold/long-term storage when the data is aging out



Apache Pulsar™

- Distributed, tiered architecture
- Separated compute from storage
- Zookeeper holds metadata for the cluster
- Stateless Broker handles producers and consumers
- Storage is handled by Apache Bookkeeper



Producer

Client application sending messages to topic managed by Broker

Consumer

Client application reading messages from a topic managed by Broker

Broker

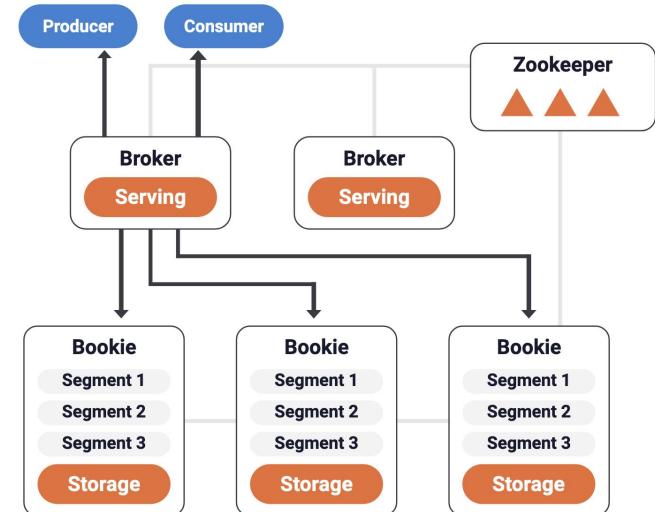
A stateless process that handles incoming message, message dispatching, communicates with the Pulsar configuration store, and stores messages in BookKeeper instances

BookKeeper

Persistent message store

ZooKeeper

Holds cluster metadata, handles coordination tasks between Pulsar clusters





Pulsar-as-a-Service

Streaming-as-a-Service built on Apache Pulsar



No Operations

Eliminate the overhead to install, operate, and scale Pulsar



Powerful Tools and APIs

Leverage the same tools used to interact with Pulsar on prem



Cloud Native

Built to run on any cloud



Zero Lock-in

Leverage Pulsar's built in integration with existing developer tools



Start for Free

Free monthly credits to help you get started quickly



DataStax Astra: Streaming Made Easy in the Cloud



Raw messages are produced: some are reviews. Of these, some are of restaurants. Must be cleaned, filtered, etc.

publish

schemaless
Pulsar topic with mixed "raw" contents

Business Architecture

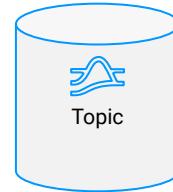
Review Injector

```
{  
    type:"hotel",  
    rating: 4.5,  
    comment:"--"  
}
```



python™

Publisher



Consumer

Output on terminal



Logical Architecture

DataStax



Lab 1

Producer & Consumer

[https://github.com/datastaxdevs/
workshop-pulsarfunctions-data-in-motion](https://github.com/datastaxdevs/workshop-pulsarfunctions-data-in-motion)



01



Introduction to
Apache Pulsar™

02



DataScience in
Event Streaming

03



Pulsar Functions

04



Introduction to
Apache Cassandra™

05



Pulsar Sinks

06



What's next?
Quiz, Homework, Next week



01

Fraud Detection

Needed to ingest high-speed writes of customer event traffic for real-time fraud detection and analysis.
Geo-replication must have little to no latency.



02

Secure Social Media, Protect Customer Privacy

Identify out-of-the-ordinary patterns to prevent malicious attacks on digital and physical assets from unauthorized applications and individuals.

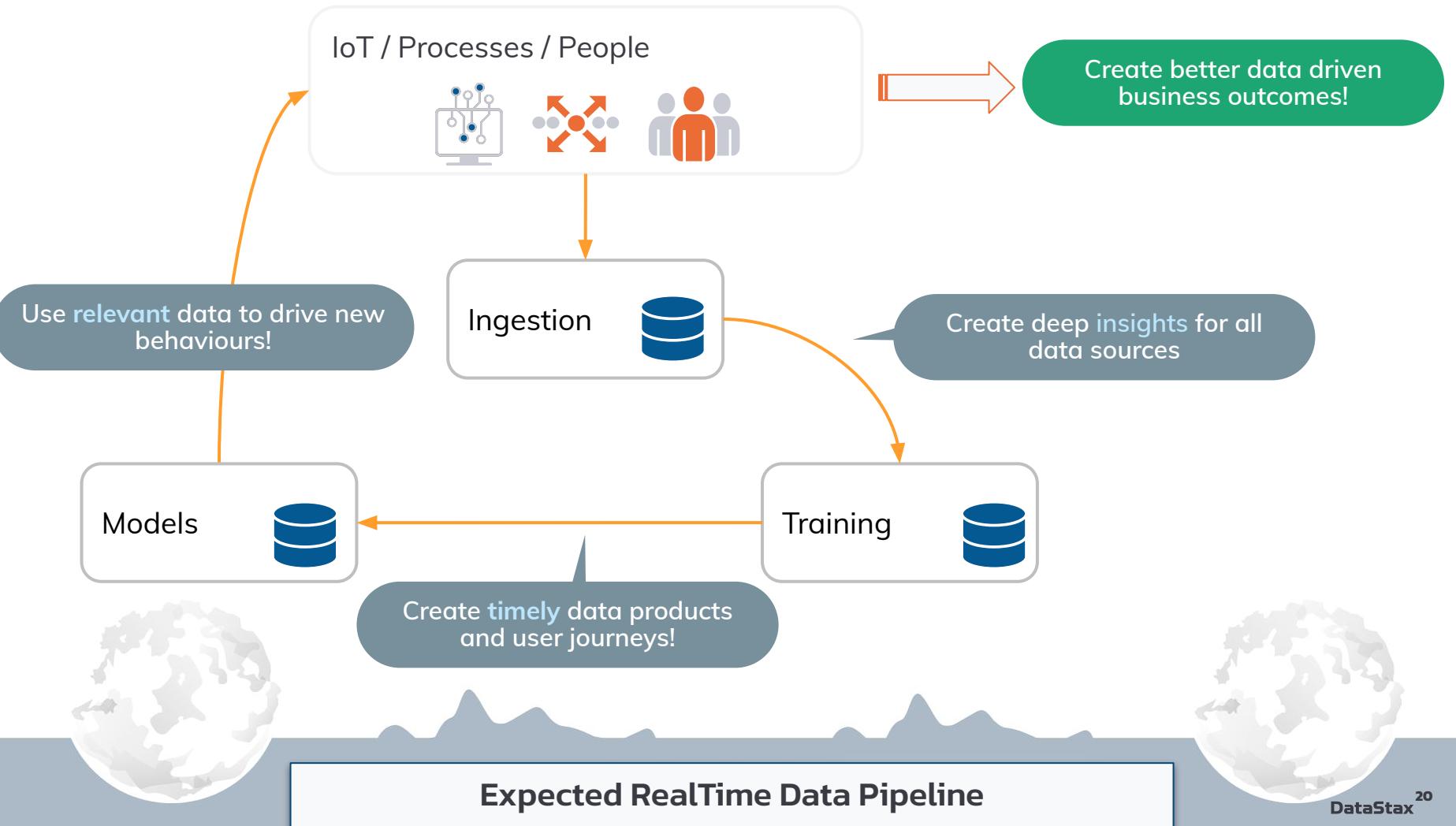
DataScience with events

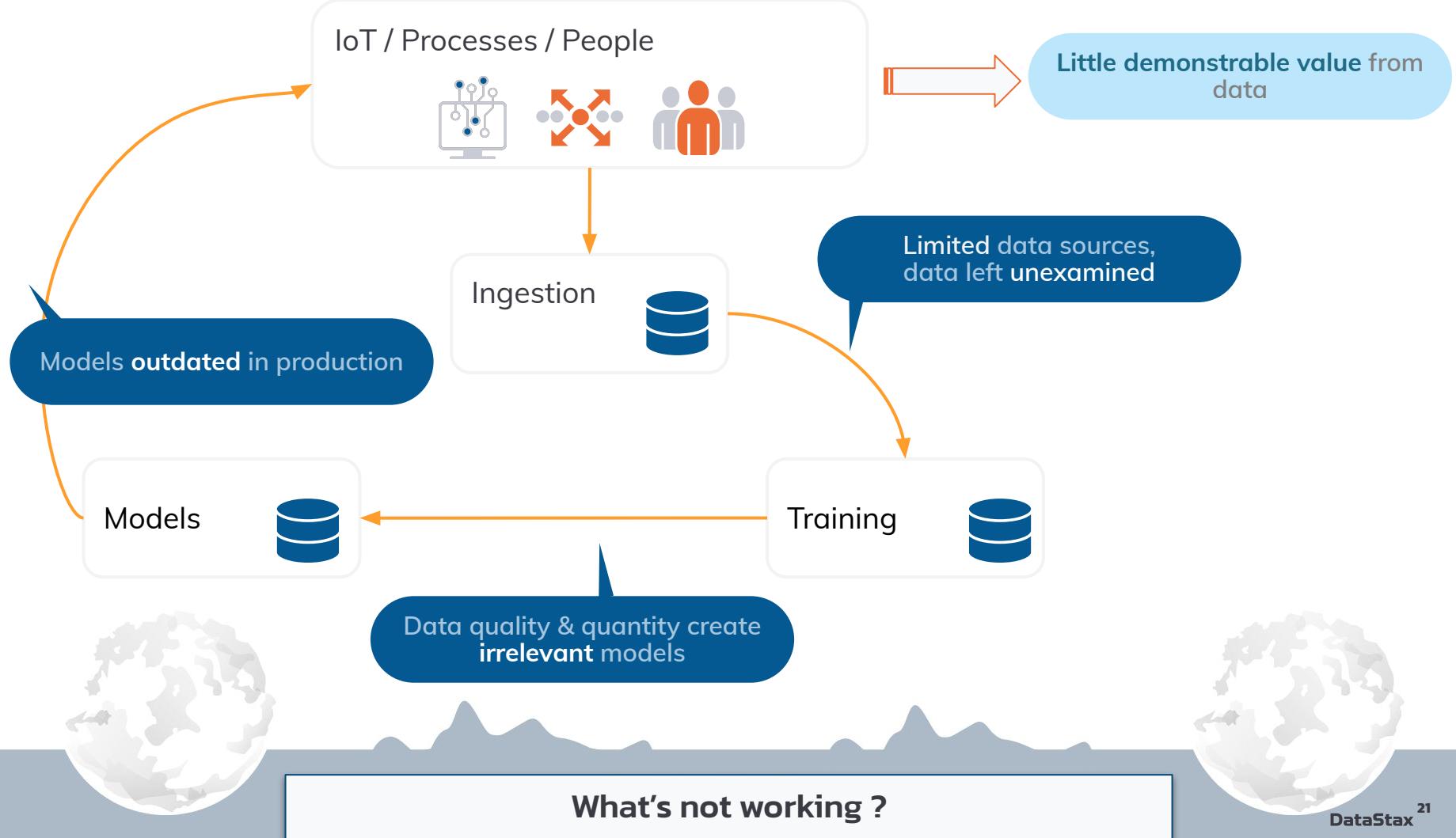
03

IoT Data Ingestion and Classification

Take in high speed data with very little latency, while processing at a different [slower] speed internally.







IoT / Processes / People



Significant business outcomes achieved!

Publish time-sensitive models - faster

Ingestion



Remove complexity of pipelines & lakes

Models



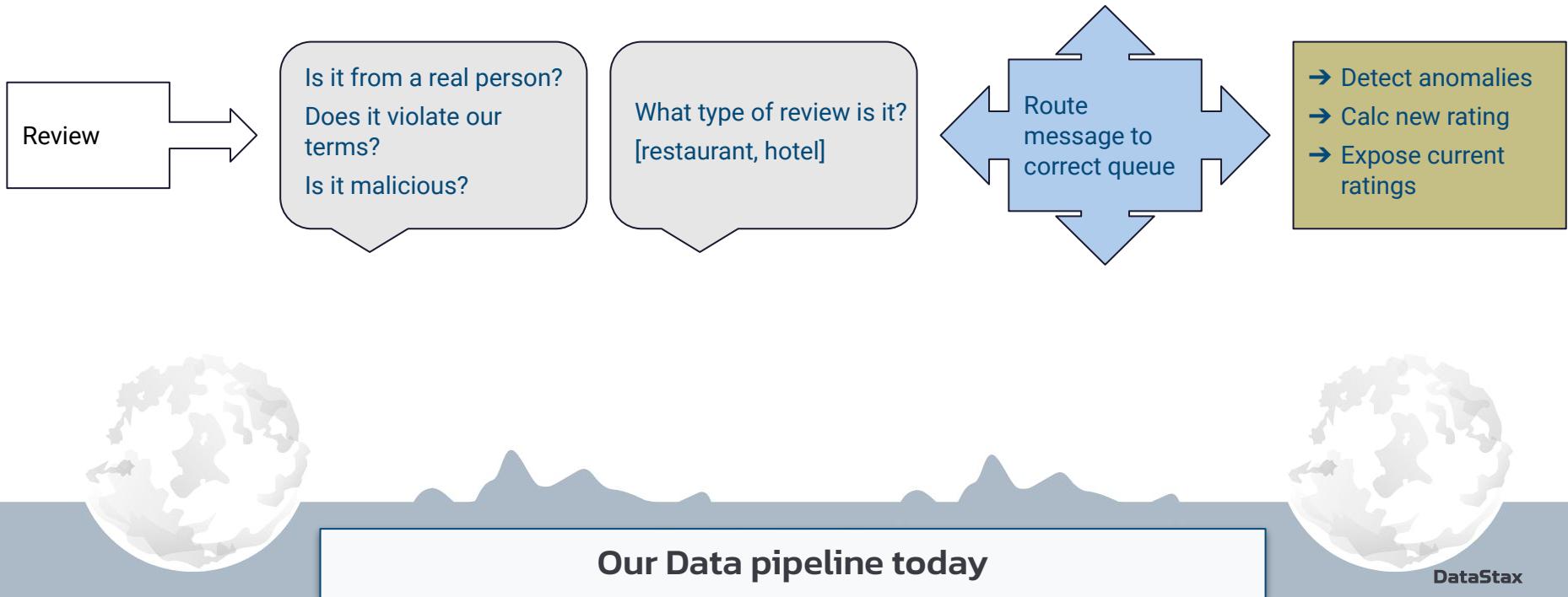
Training



Deeper analysis of data sets to enrich the models - faster



Cassandra and Pulsar to the rescue



01



Introduction to
Apache Pulsar™

02



DataScience in
Event Streaming

03



Pulsar Functions

04



Introduction to
Apache Cassandra™

05



Pulsar Sinks

06



What's next?
Quiz, Homework, Next week

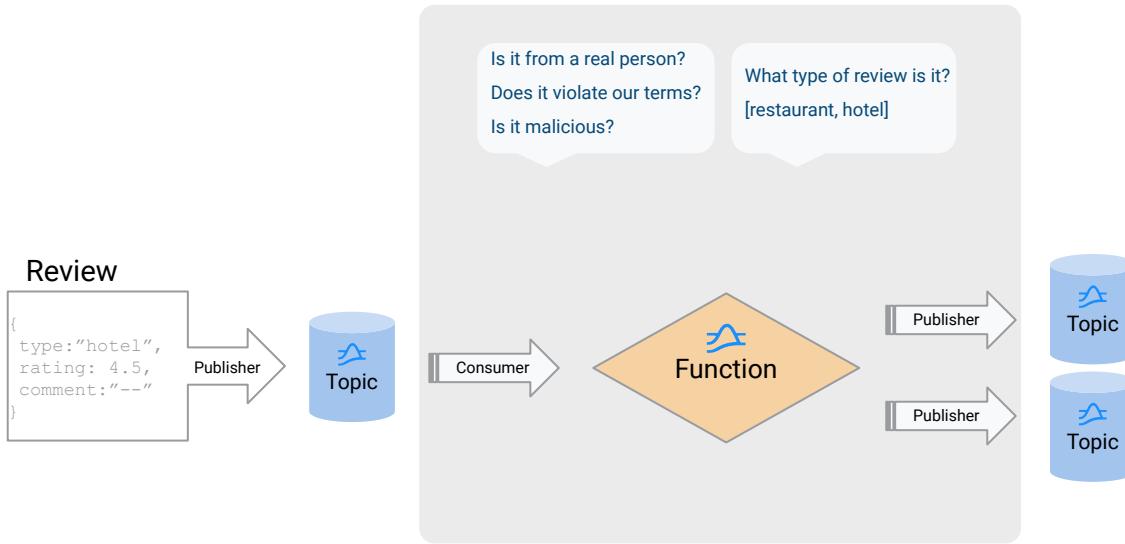


Agenda



Pulsar Functions





Architecture Overview



Lab 2

Pulsar Functions

[https://github.com/datastaxdevs/
workshop-pulsarfunctions-data-in-motion](https://github.com/datastaxdevs/workshop-pulsarfunctions-data-in-motion)



01



Introduction to
Apache Pulsar™

02



DataScience in
Event Streaming

03



Pulsar Functions

04



Introduction to
Apache Cassandra™

05



Pulsar Sinks

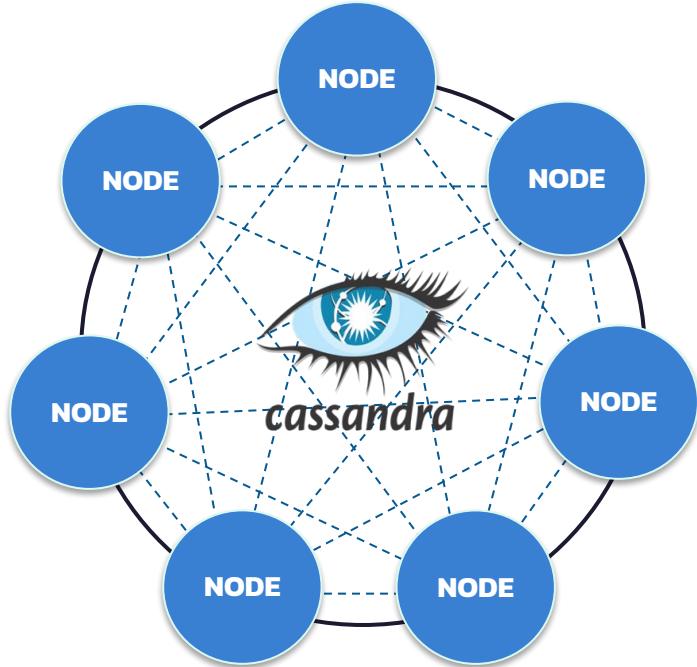
06



What's next?
Quiz, Homework, Next week



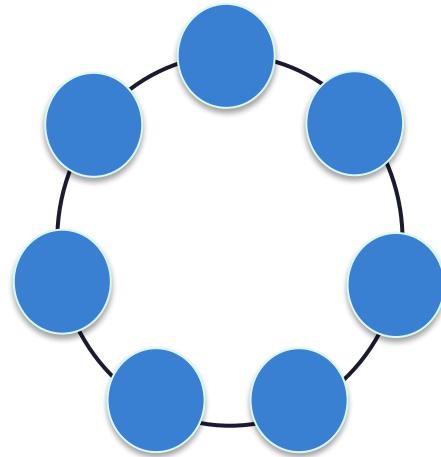
Agenda



- Big Data Ready
- Read / Write Performance
- Linear Scalability
- Highest Availability
- Self-Healing and Automation
- Geographical Distribution
- Platform Agnostic
- Vendor Independent

Apache Cassandra's Awesomeness

Partitioning over distributed architecture makes the database capable to handle data of any size: we mean petabytes scale. Need more volume? Add more nodes.

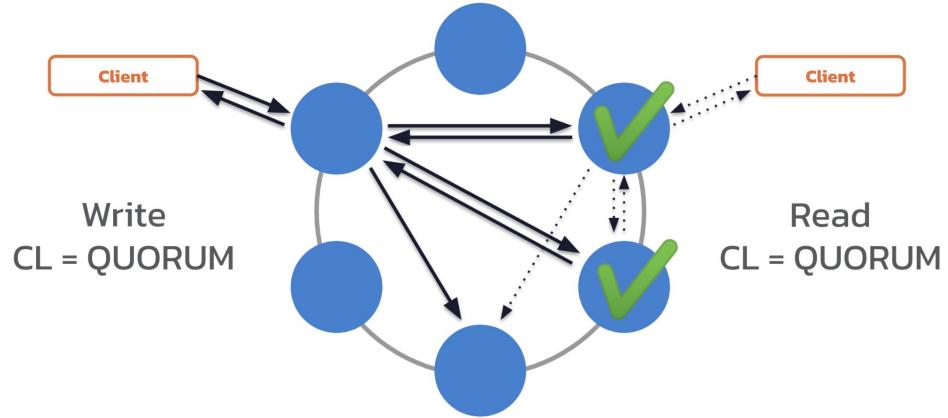


Big Data Ready



Even a single Cassandra node is very performant but a cluster consisting of multiple nodes and data centers brings throughput to the next level.

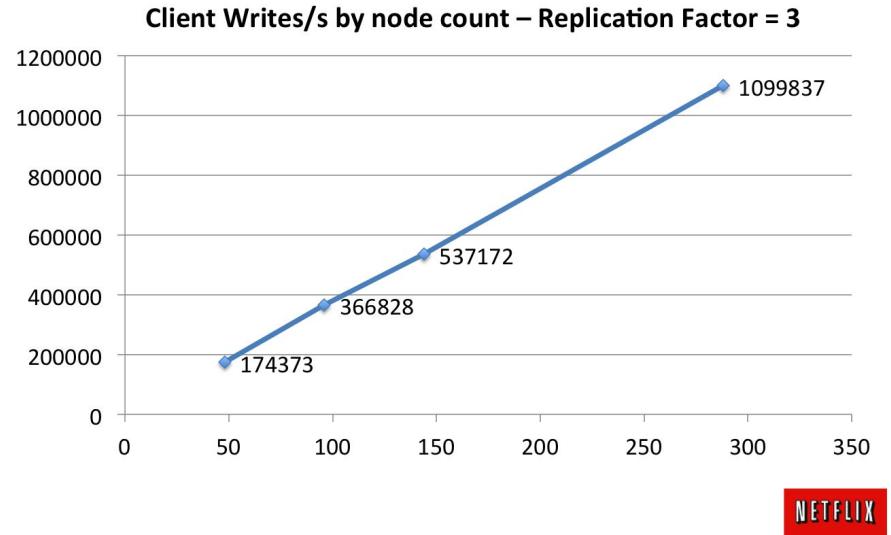
Decentralisation (**masterless architecture**) means that every node is able to deal with any request, read or write.



Read / Write Performance



- For volume or velocity, there are no limitations
- **Linear** - No overhead on new nodes, scales with your needs*

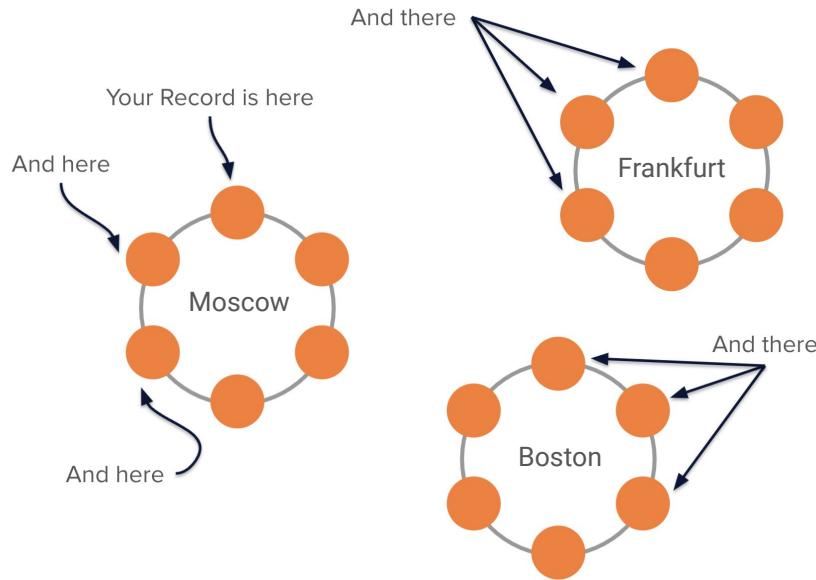


Linear Scalability



Replication, Decentralisation, and Topology-Aware Placement Strategy take care of possible downtimes:

- Multiple Live Replicas
- No Single Point of Failure
- Network topology-aware data placement
- Client-side Smart Reconnection and Strong Retry Mechanism



Highest Availability



Operations for a huge cluster can be exhausting so Apache Cassandra clusters are smart and able to scale, change data placement and recover automatically.

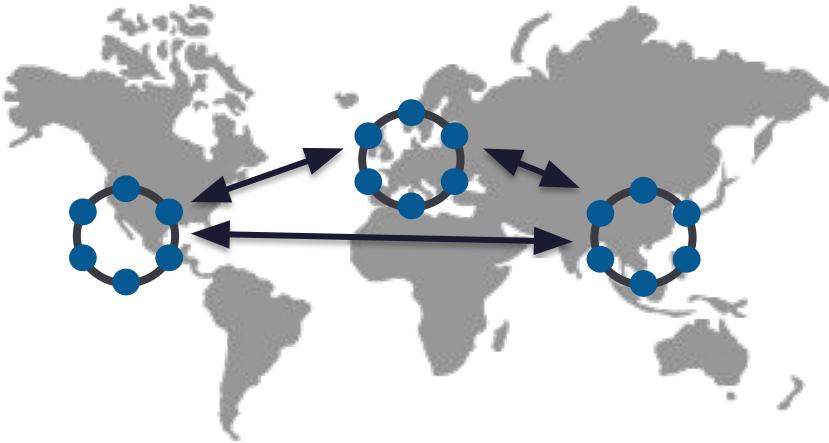


Self-Healing and Automation



Cassandra's trademark is multi-datacenter deployments, granting you an exceptional capability for disaster tolerance while keeping your data close to your clients - worldwide.

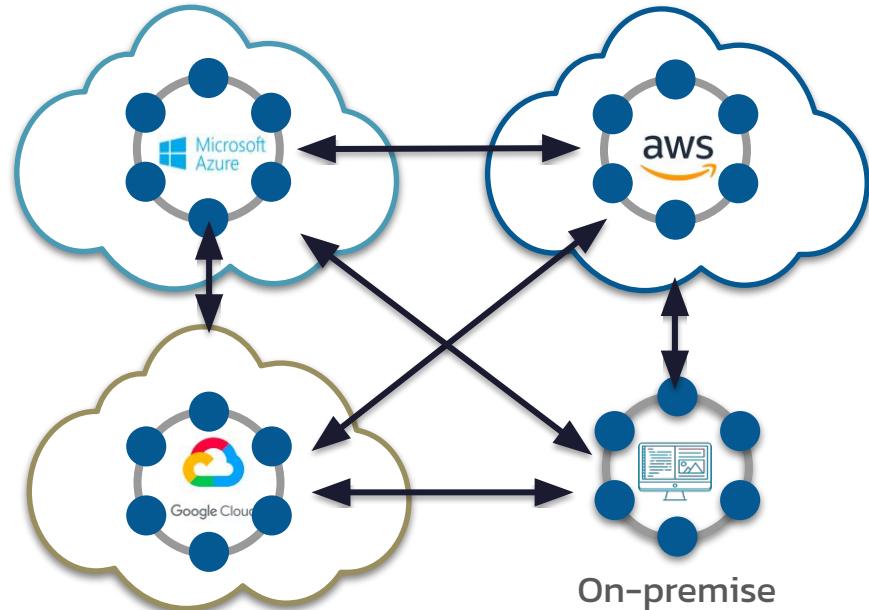
All DCs are active (available for both writes and reads)!



Geographical Distribution



Apache Cassandra is **not bound to any platform** or service provider, helping you build hybrid-cloud and multi-cloud solutions with ease.



Platform Agnostic

Cassandra doesn't belong to any of commercial vendors but controlled by a non-profit Open Source **Apache Software Foundation**, already familiar to you by *Hadoop*, *Spark*, *Kafka*, *Zookeeper*, *Maven* and many other projects.



Vendor Independent



01



Introduction to
Apache Pulsar™

02



DataScience in
Event Streaming

03



Pulsar Functions

04



Introduction to
Apache Cassandra™

05



Pulsar Sinks

06



What's next?
Quiz, Homework, Next week

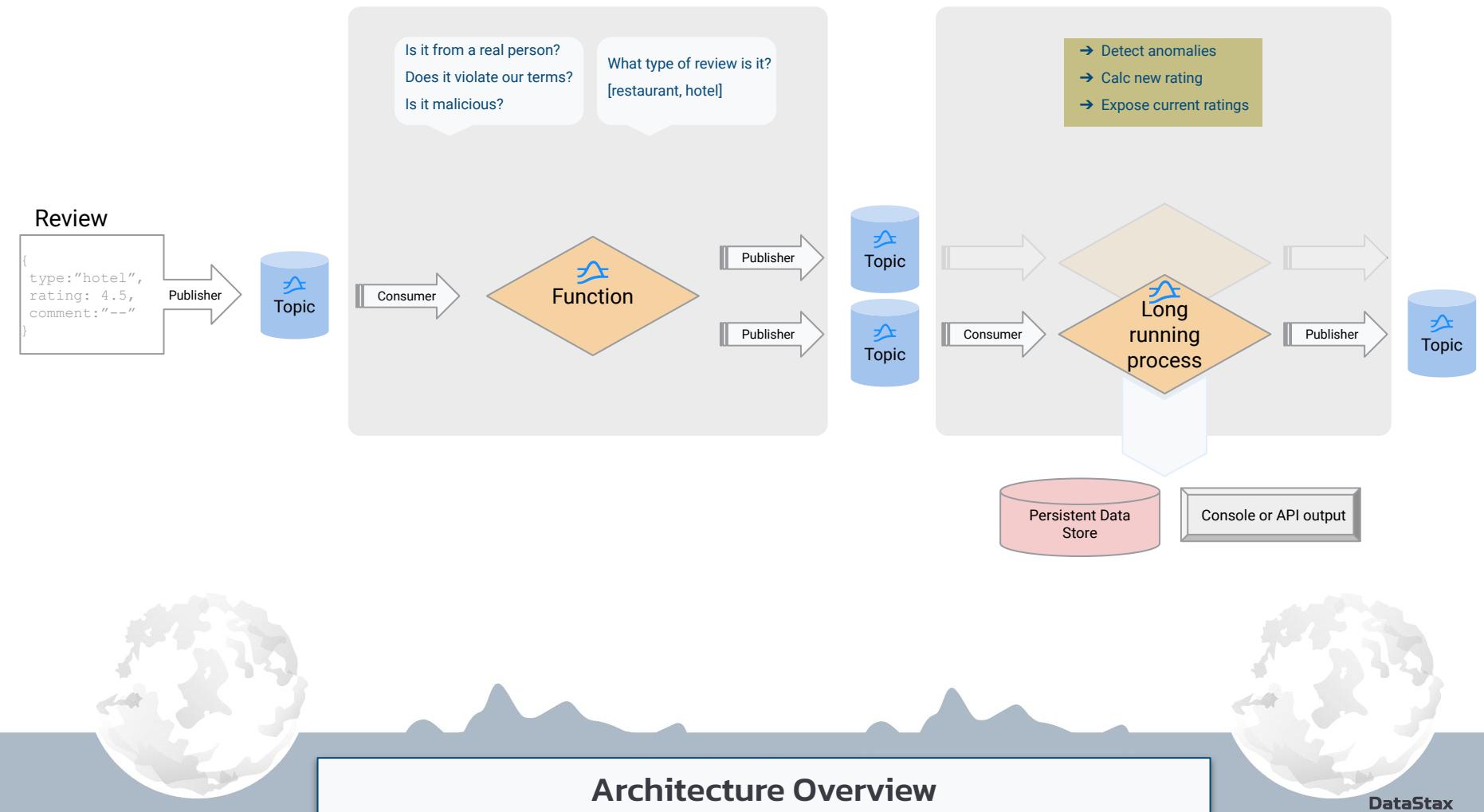


Agenda



Pulsar Sinks







Lab 3

Pulsar Sinks

[https://github.com/datastaxdevs/
workshop-pulsarfunctions-data-in-motion](https://github.com/datastaxdevs/workshop-pulsarfunctions-data-in-motion)



01



Introduction to
Apache Pulsar™

02



DataScience in
Event Streaming

03



Pulsar Functions

04



Introduction to
Apache Cassandra™

05



Pulsar Sinks

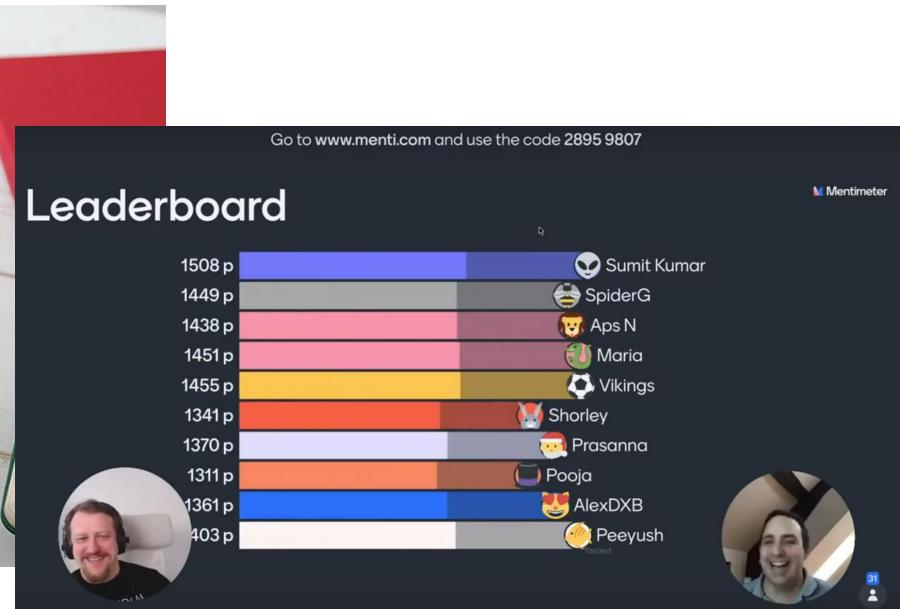
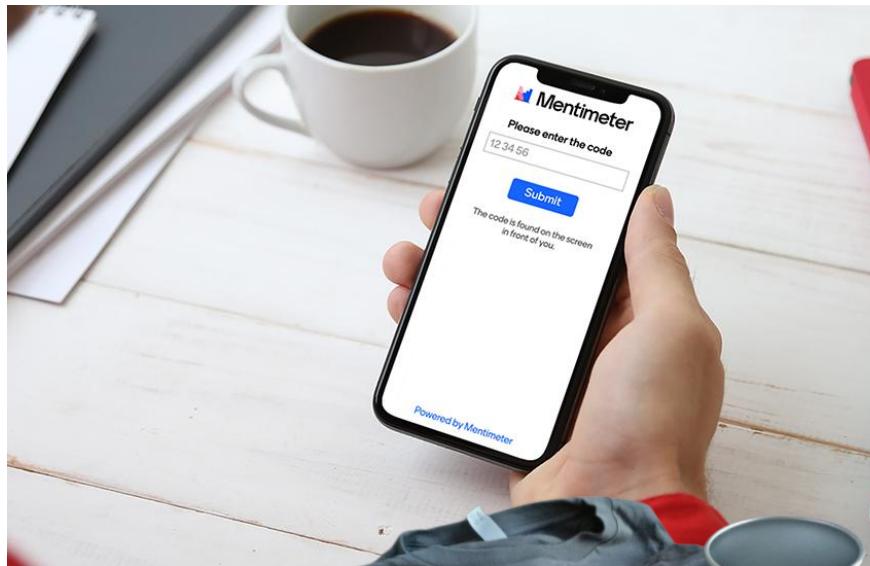
06



What's next?
Quiz, Homework, Next week



Agenda



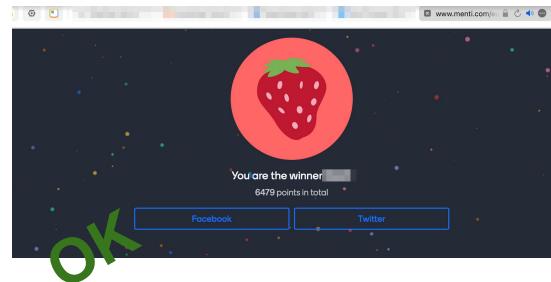
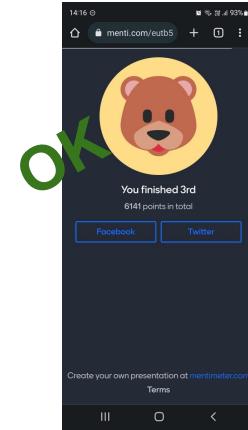
**menti . com ⇒ enter code
Don't answer in YT chat
Look at phone (not at YT)**

Quiz on "Menti" !

SWAG WINNERS



Congratulations to 1st, 2nd and 3rd place
on the Menti quiz!



To claim your prize:

Take a screenshot of your Menti screen

Fill the form at dtsx.io/workshop-swag

NO!



Swag Winners!



!discord

dtsx.io/discord



David Jones-Gilardi

012345

AaronP

BInary

Chelsea Navo

Jeremy Hanna

John Sanda

Patrick_McFadin

-samu-

6304-42J8

Aahlya

Abdurahim

abhi3pathi

Abhiis.s

Abhineet

Abirish

En LIGNE — 560

RIGGITYREKT Hier à 21:14

I have a 5 node datacenter, 4 nodes are on dse version 5.1.20, one is on dse5.0.15. I am doing some mixed version testing for a class and the one node that is 5.0.15 is coming up as an analytics workload. I don't have /etc/default/dse, instead I am using /etc/init.d/dse-cassandra. how do I make that node start in cassandra workload, not in analytics?

RIGGITYREKT Hier à 23:39

Okay I found out my issue, when I started DSE 5.0.15 it had endpointsnitch set to DseSimpleSnitch, the rest of my cluster is using PropertyFileSnitch, when I change it to PropertyFileSnitch, it still uses the simple snitch config. looking at the docs I see there is a way to go to GossipingPropertyFileSnitch, but I need the property file one. I can wipe this db, do anything with this node to get this done. how do I fix this? @here

19 novembre 2021

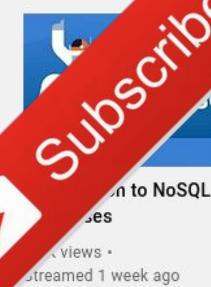
@RIGGITYREKT Okay I found out my issue, when I started DSE 5.0.15 it had endpointsnitch set to DseSimpleSnitch, the rest... Erick Ramirez Aujourd'hui à 02:19 mixed versions isn't supported and you're guaranteed to run into weird issues that will cause further problems down the track

@RIGGITYREKT I have a 5 node datacenter, 4 nodes are on dse version 5.1.20, one is on dse5.0.15. I am doing some mixed v... Cedrick Lunven Aujourd'hui à 09:01 When you start a node you have parameters -k for analytics, -g for graph and -s for search. To remove analytics check and remove -k

Envoyer un message dans #workshop-chat

Datastax Developers Discord (18k+)

Subscribe



Subscribe



How to create an Authentication Token in...
37 views • 4 weeks ago

How to use the Data Loader in Astra DB
62 views • 4 weeks ago

Astra DB Sample App Gallery
36 views • 4 weeks ago

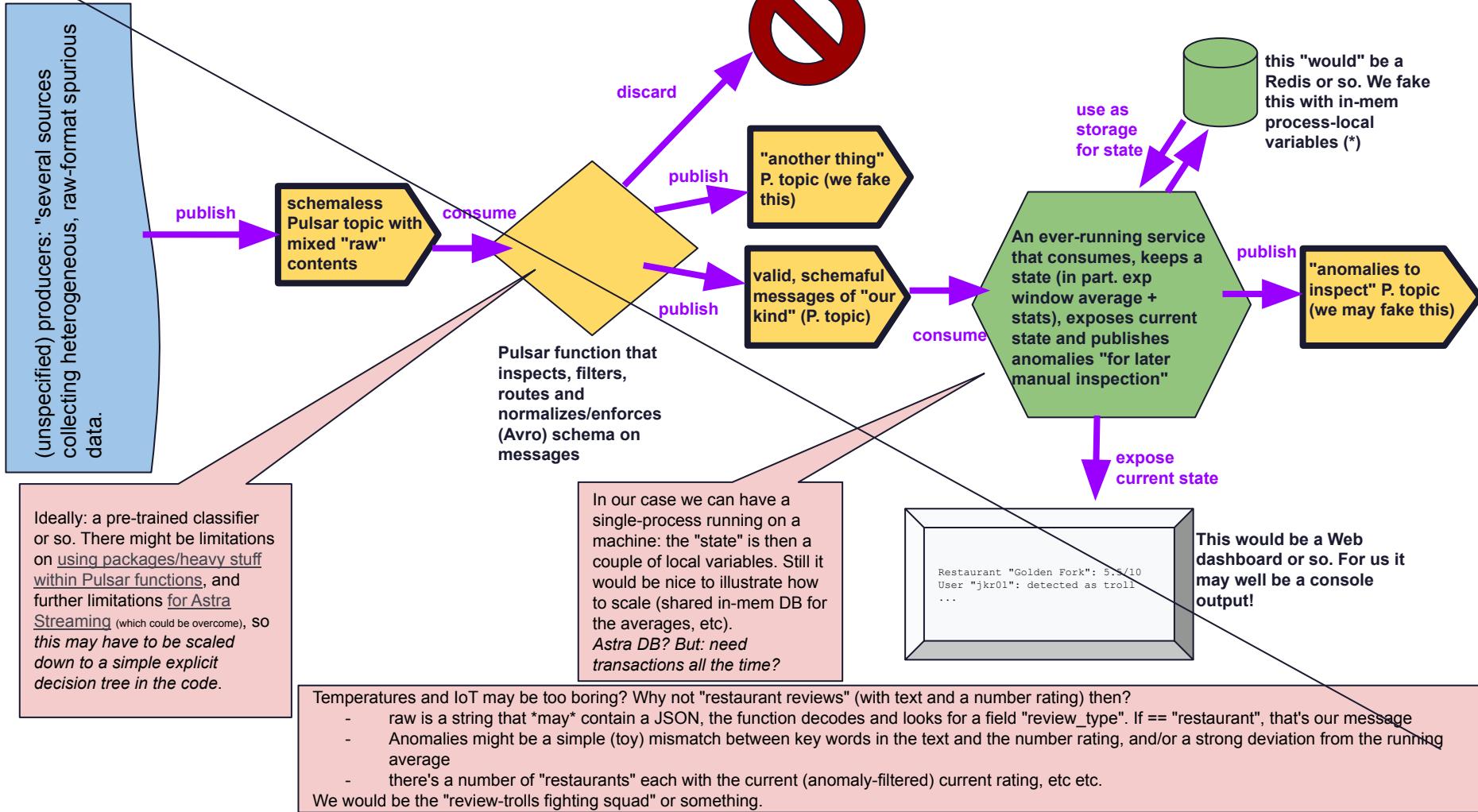
How to use Secure Connect in Astra DB
42 views • 4 weeks ago

Cassandra Day India: CL Room (Workshops)
2.4K views • Streamed 4 weeks ago

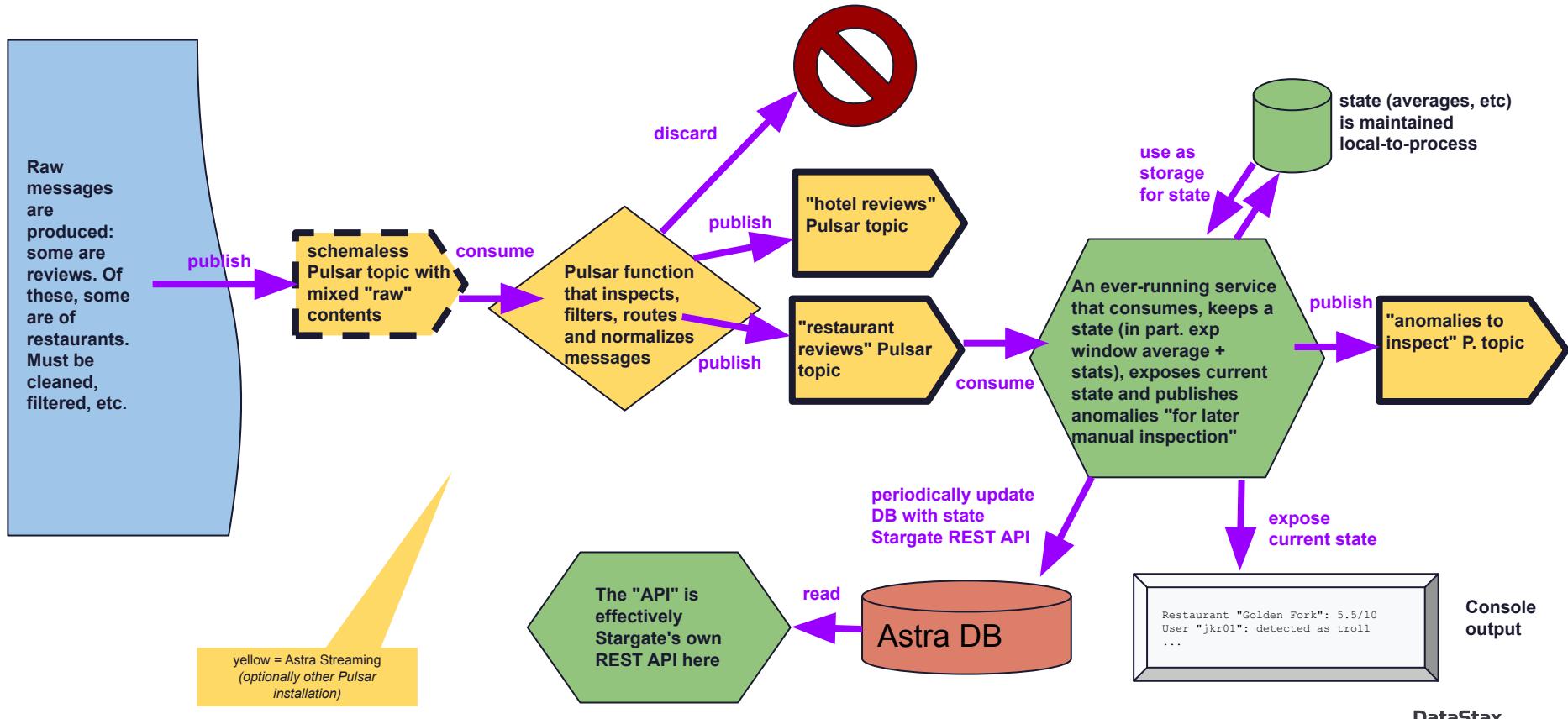
Cassandra Day India: RF Room (Talks)
1.3K views • Streamed 1 month ago

Thank You!

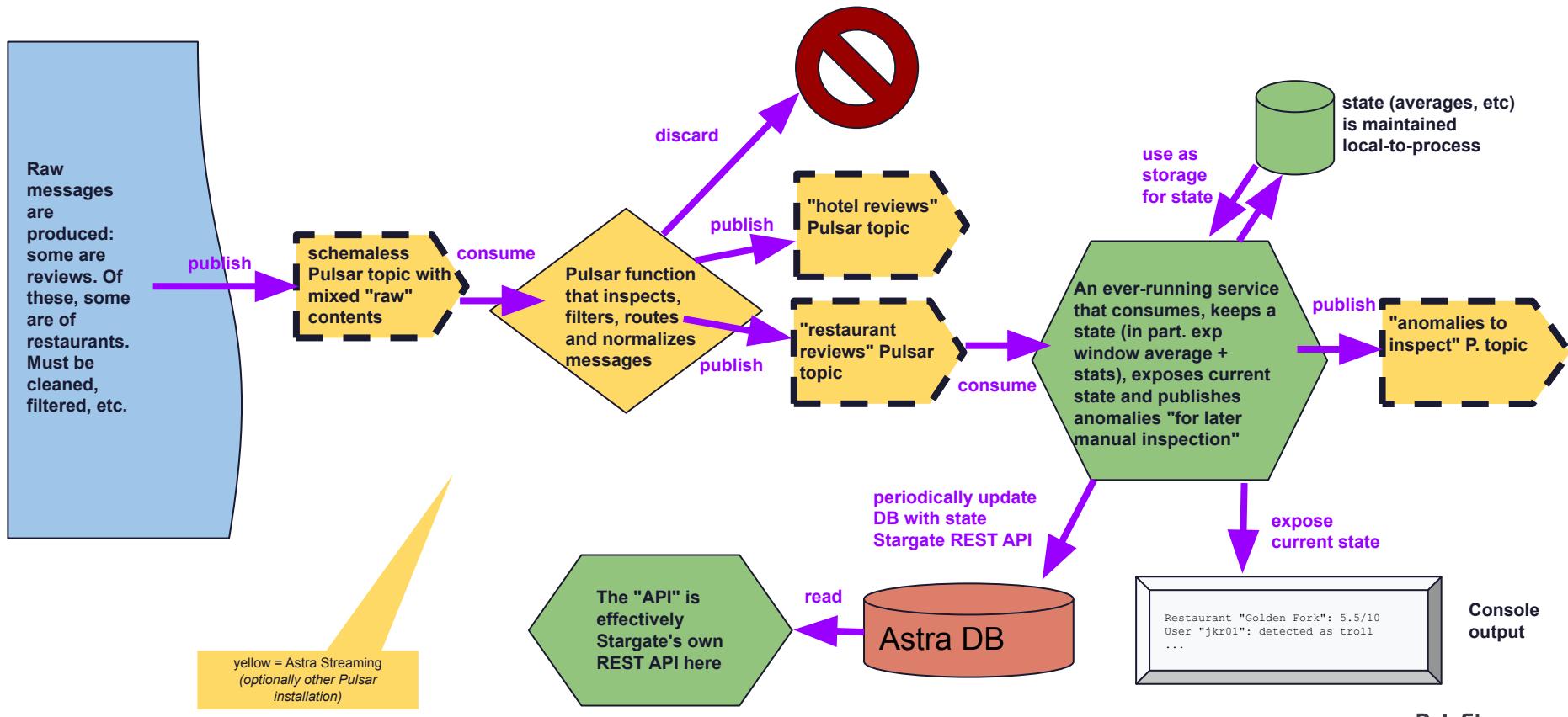




Goal Architecture



Current Architecture



yellow = Astra Streaming
(optionally other Pulsar
installation)

The "API" is effectively Stargate's own REST API here

Astra DB

Restaurant "Golden Fork": 5.5/10
User "jkr01": detected as troll

Console output

DataStax

Data Layer: Architecture

