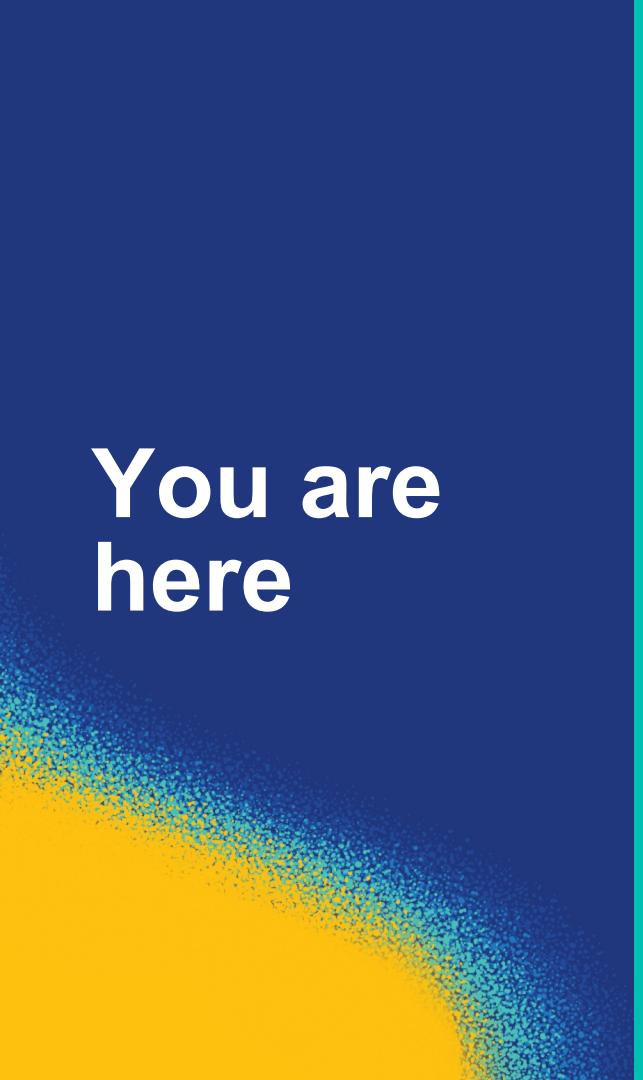


Data Streams & Pipelines

Observing with the Elastic Stack

Pete Hampton
Data Engineer, Elastic Security



You are
here

- 01 Intro**
- 02 Data Engineering**
- 03 Case Study 1 – Batch**
- 04 Case Study 2 – Streaming**
- 05 Advanced Usage**

Who am I?

-  **Name:** Pete Hampton
-  **Live:** Belfast, N. Ireland
-  **PhD:** Artificial Intelligence
-  **Work:** Data processing & Distributed systems

Who are you?

- 😎 Developer / DevOps / Data Scientist
- 😎 Builds / Operates pipelines
- 😊 Elastic Stack rookie

Mission

The goals for this presentation

- Discuss modern challenges in Data Engineering
- Elastic Stack - Get started / enhance your usage
- Share anecdotes, use cases, and recommendations

Intro

Pipelines? 😕 Streams? 😕

Pipelines

Definition for this presentation

- Typically, **batch / scheduled** based
- Can be slow(er) to deliver insights
- Moving data around on demand
- Think Airflow, Luigi, & other workflow/schedule systems

Streams

Definition for this presentation

- Not talking about low latency ($<100\mu\text{s}$) or hard real-time ($<20\text{ms}$)
- Think fast data & *near* real-time
- Processing of data as it received (**Polled or Pushed**)
- Think tech *Flink, Kafka, Spark Streaming, RabbitMQ*, etc

Data Engineering

Systems, Services & Infrastructure



“Data” engineers design and build pipelines that transform and transport data into a format wherein, by the time it reaches the Data Scientists or other end users, it is in a highly usable state.

[https://quanhub.com/what-is-data-engineering/](https://quanthub.com/what-is-data-engineering/)

Data Engineering

- *Find and make data usable for products, systems and stakeholders*
- Use Cases
 - Big Data Ingest
 - Business Intelligence
 - Contextual product enhancements, eg Recommender Systems
- Rapidly evolving
 - Monoliths -> SOA -> Microservices -> Cloud Native -> Low Code -> No Code
 - New tools and programming languages come and go
 - Industry adopting streaming approaches more and more

Data Engineering

Day to day struggles

- Taming Large / Complex data infrastructures and integrations
- Few or no engineers understand the entire data plane
- Development slows down as things get more complex
- Data processing can continue after it no longer yields value

Integration Points

Systems and Services that connect

- Databases / Data Warehouses
- Cloud Providers
- Edge Services
- 3rd Party APIs
- And much more

What makes it so complex?

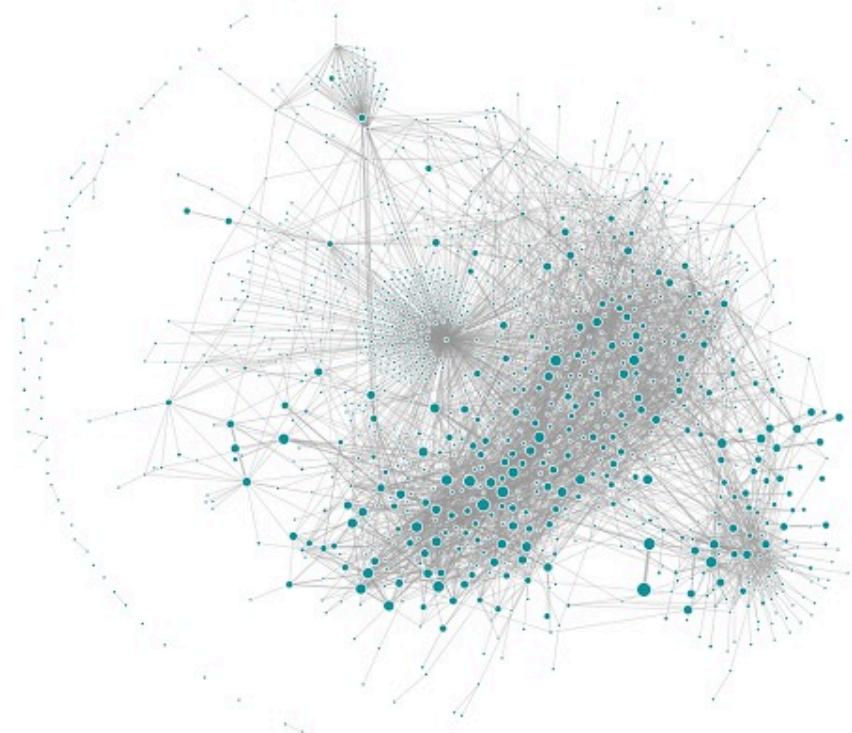
- Too many tools, clouds, SaSS
- Expertise comes and goes
- No or scattered system of record
- Synchronization is a tricky business
- Scaling Batch / ETL / Lambda architectures are difficult

Uber Service Map

~2018

- 2k+ critical microservices
- *“engineers had to work through around 50 services across 12 different teams in order to investigate the root cause of the problem.”*

Ref: <https://eng.uber.com/microservice-architecture/>





Elasticsearch

```
~ → curl localhost:9200 | jq
% Total    % Received % Xferd  Average Speed   Time     Time     Time  Current
                                         Dload  Upload   Total   Spent   Left  Speed
100  552  100  552    0      0  107k      0 --:--:-- --:--:-- --:--:-- 134k
{
  "name": "Pete's-MacBook-Pro.local",
  "cluster_name": "elasticsearch",
  "cluster_uuid": "pHen4Y7oRo00-sU0zGbkgzg",
  "version": {
    "number": "7.14.0",
    "build_flavor": "default",
    "build_type": "tar",
    "build_hash": "dd5a0a2acaa2045ff9624f3729fc8a6f40835aa1",
    "build_date": "2021-07-29T20:49:32.864135063Z",
    "build_snapshot": false,
    "lucene_version": "8.9.0",
    "minimum_wire_compatibility_version": "6.8.0",
    "minimum_index_compatibility_version": "6.0.0-beta1"
  },
  "tagline": "You Know, for Search"
}

~ → █
```

Elasticsearch

Handy components for Data Engineers

- Data Ingest Pipelines
- Data Streams
- Index Lifecycle Management (ILM)
- Data Tiers
- Painless scripts
- ...SQL, EQL

Stack Management

Topologies for pipeline observability

- Official Terraform provider available for Cloud deployments
- If you aren't using *Elasticsearch as a Service*, consider ECK/ECE
- A couple of effective topologies (hardware profiles)

Instances

Health ▾ Instance configuration ▾ Data tier ▾

Zone europe-west2-a

Instance #0

● Healthy · v7.14.1 · 1 GB RAM · GCP.APM.N2.68X32X45

Instance #0

● Healthy · v7.14.1 · 8 GB RAM ·
GCP.ES.DATAHOT.N2.68X32X45 · data_hot · data_content ·
master · coordinating · ingest

Disk allocation
47 MB / 360 GB 0%

JVM memory pressure
Normal 2%

Zone europe-west2-b

Tiebreaker #2

● Healthy · v7.14.1 · 1 GB RAM ·
GCP.ES.MASTER.N2.68X32X45 · master eligible

Disk allocation
0 GB / 45 GB 0%

JVM memory pressure
Normal 9%

Zone europe-west2-c

Instance #1

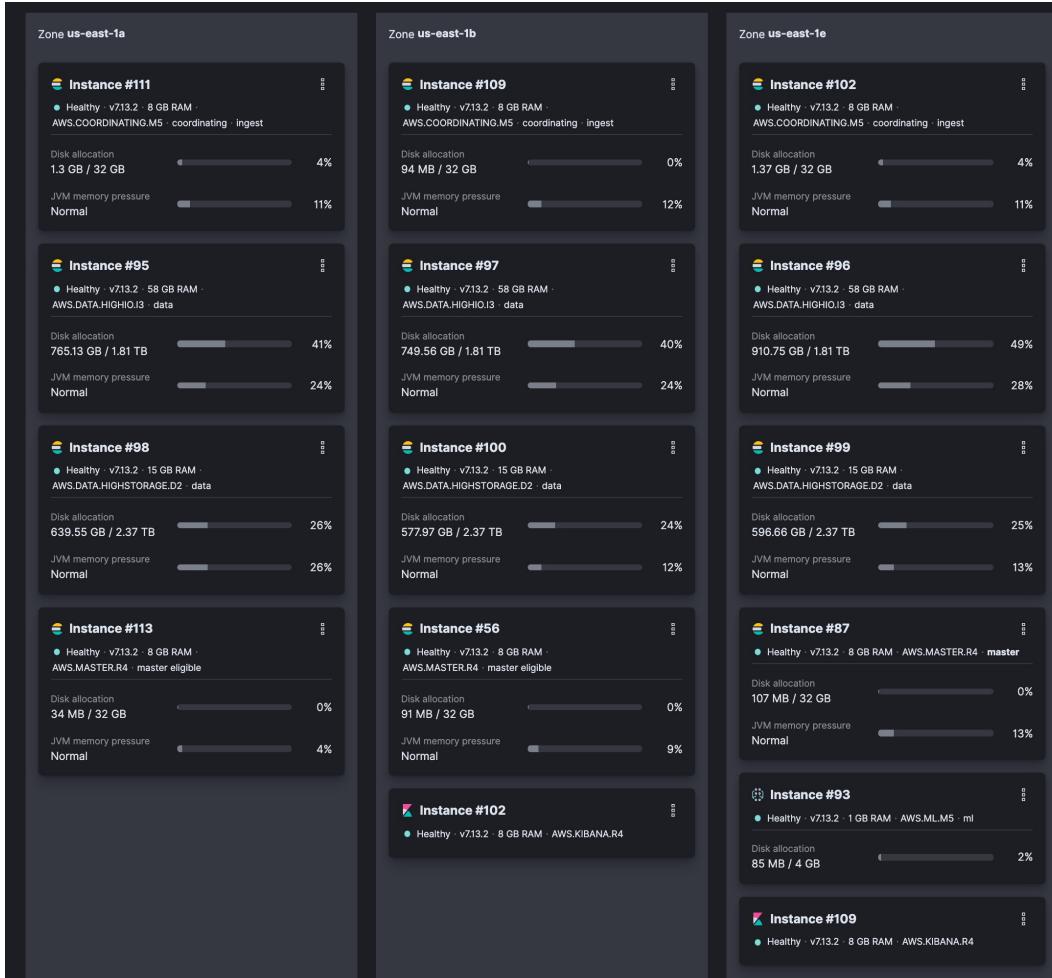
● Healthy · v7.14.1 · 8 GB RAM ·
GCP.ES.DATAHOT.N2.68X32X45 · data_hot · data_content ·
master eligible · coordinating · ingest

Disk allocation
49 MB / 360 GB 0%

JVM memory pressure
Normal 2%

Instance #0

● Healthy · v7.14.1 · 8 GB RAM · GCP.KIBANA.N2.68X32X45

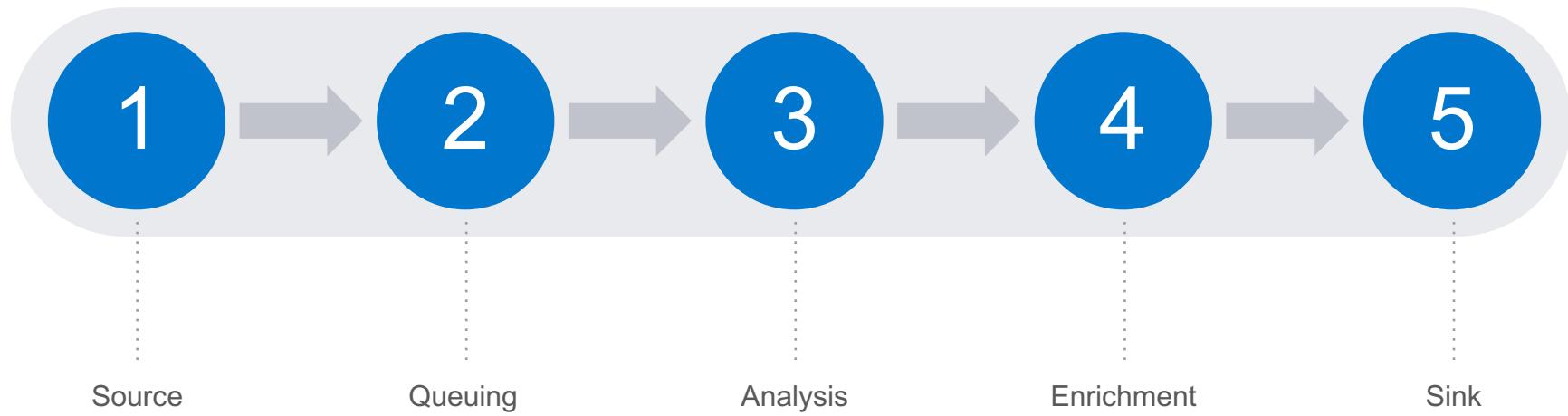


Case Study 1

Batch

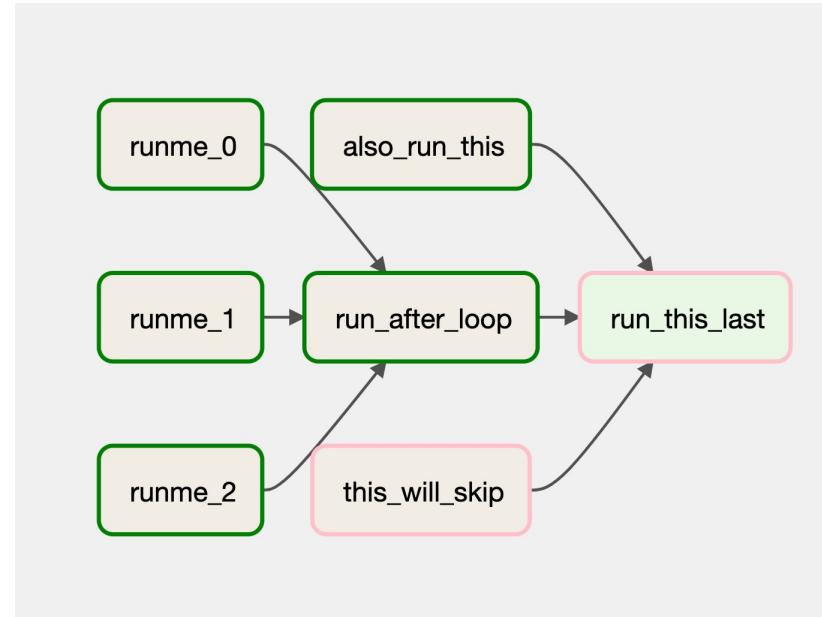
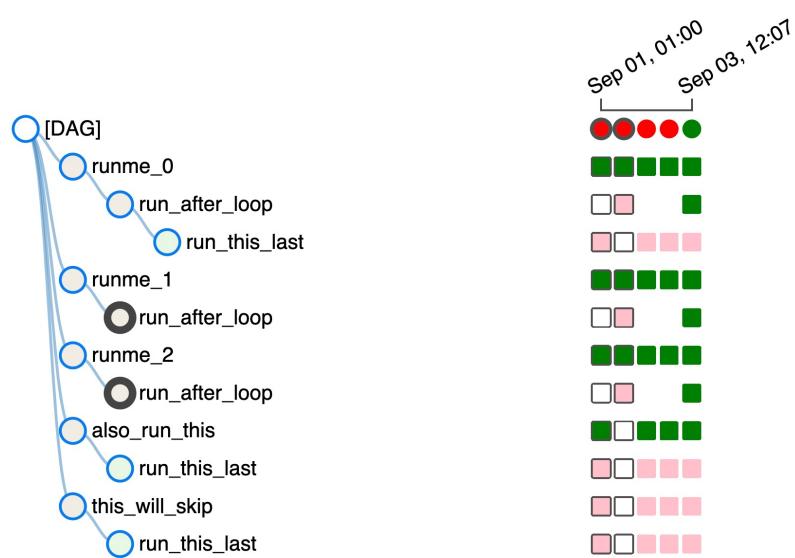
Common Pattern

Ingestion to Storage



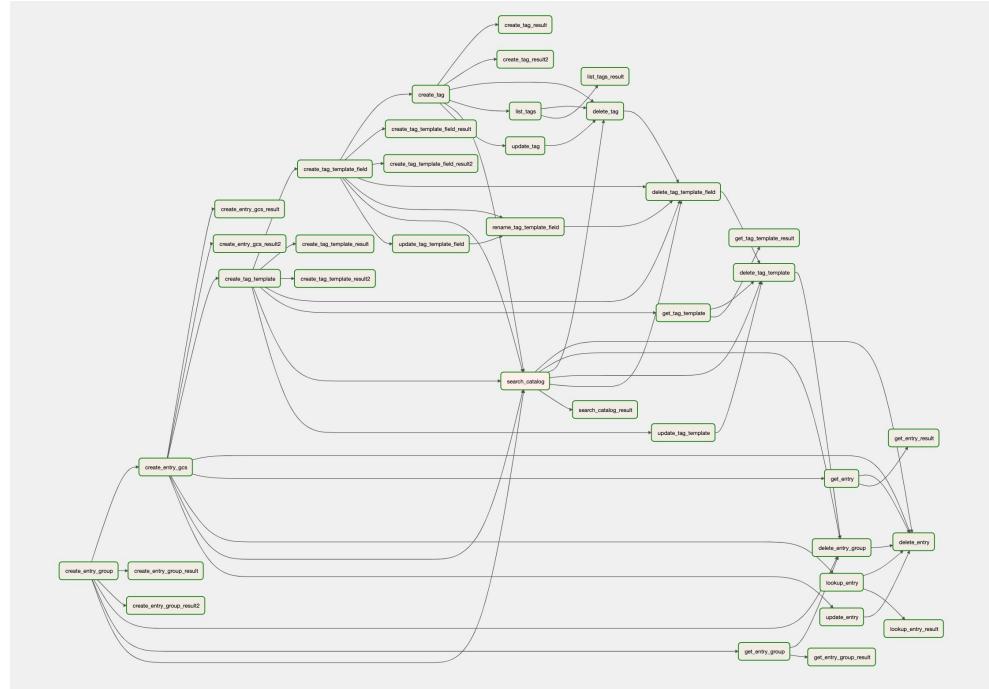
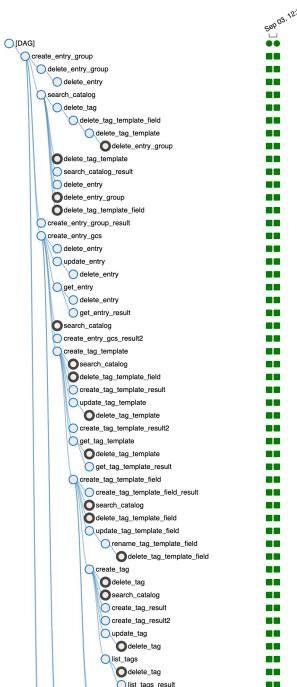
Batch Architecture

Simple



Batch Example

Not so simple



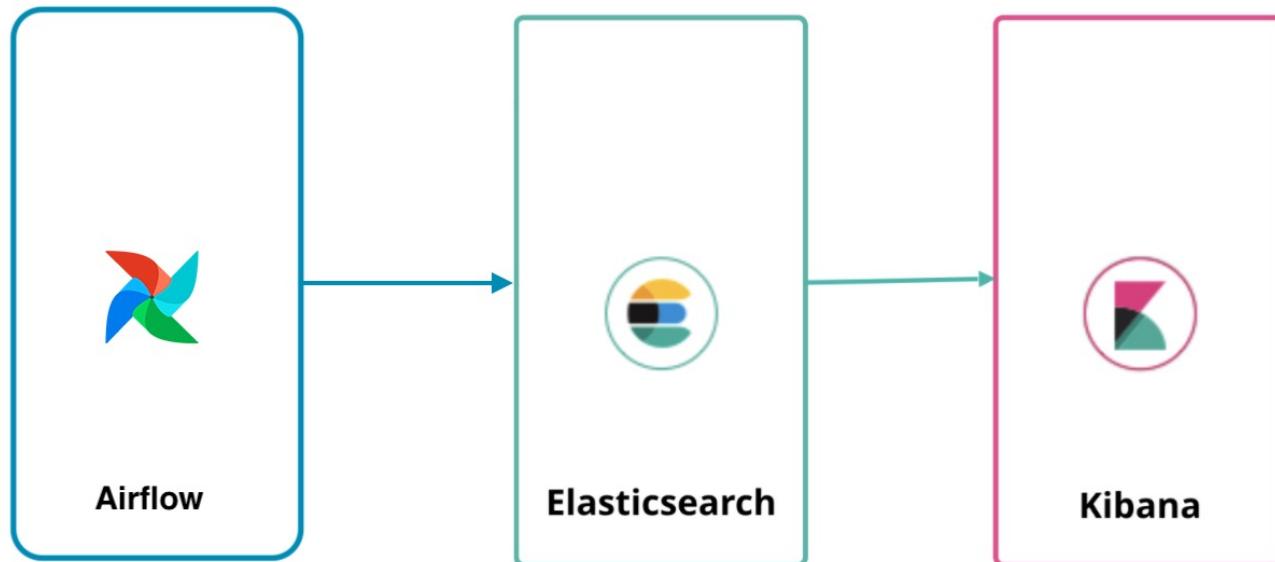
Logs?

First line of defense

- Collect them if you can
- How to ship though? 🤔

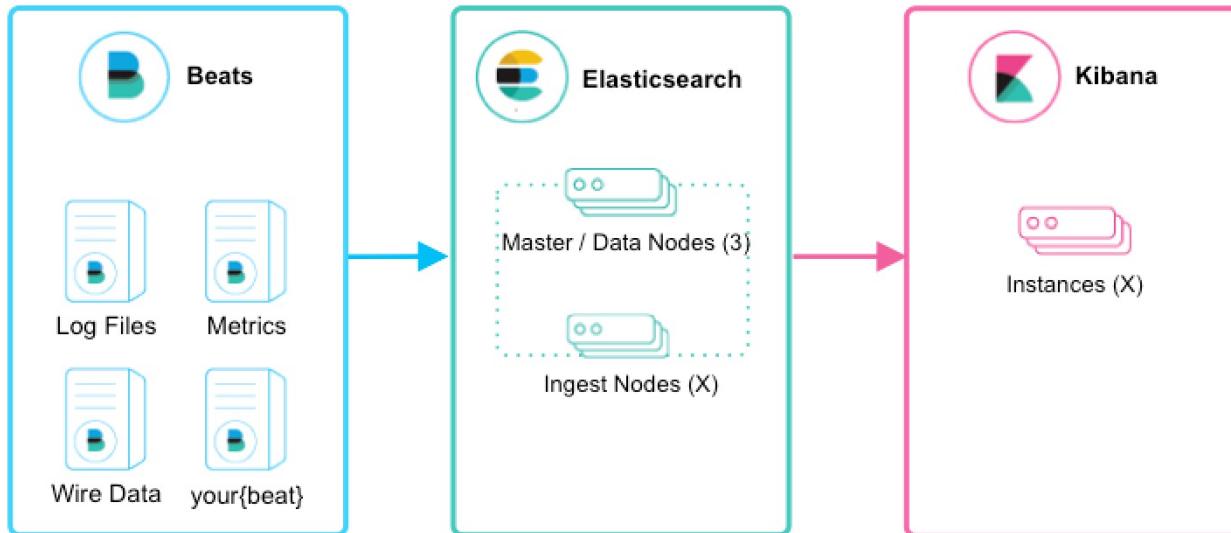
Basic Setup

Orchestrator to Elasticsearch



Basic Setup

Beats to Elasticsearch



Logs

Get the most out of your logs

- More **isn't** always better
- Structure logs and **make consistent** across services
- Build in **contextualization**
- Develop a **schema**, or adopt an existing one (<https://github.com/elastic/ecs>)

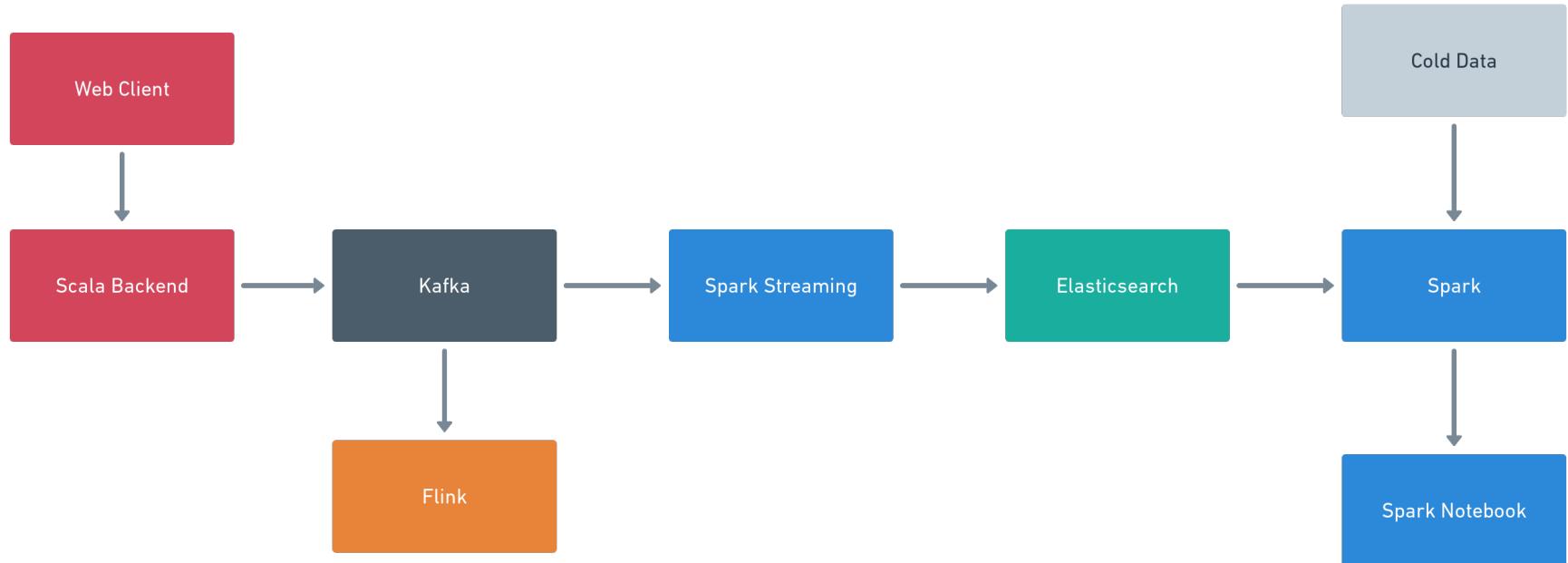
```
abstract class Event(  
    exception: Exception? = null,  
    throwable: Throwable? = null,  
    val level: LogLevel = LogLevel.DEBUG,  
    val message: String? = null,  
    val reason: String? = null,  
    val statusCode: String? = null,  
    val warning: String? = null  
) {  
    val application: String? = applicationLocal.get()  
    val caller: String = Thread.currentThread().name  
    val correlationId: String? = contextualLocal.get()  
    val event: String = this.javaClass.simpleName  
    val exceptionMessage: String? = exception?.message  
    val thread: String = Thread.currentThread().name  
    val time: Long = System.currentTimeMillis()  
}
```

```
class BadRequestReceived(message: String) : Event(  
    level = LogLevel.ERROR,  
    message = message  
)  
  
class CleaningUpAsyncWorker(name: String) : Event(  
    level = LogLevel.INFO,  
    message = "cleaning up async worker $name"  
)
```

Case Study 2

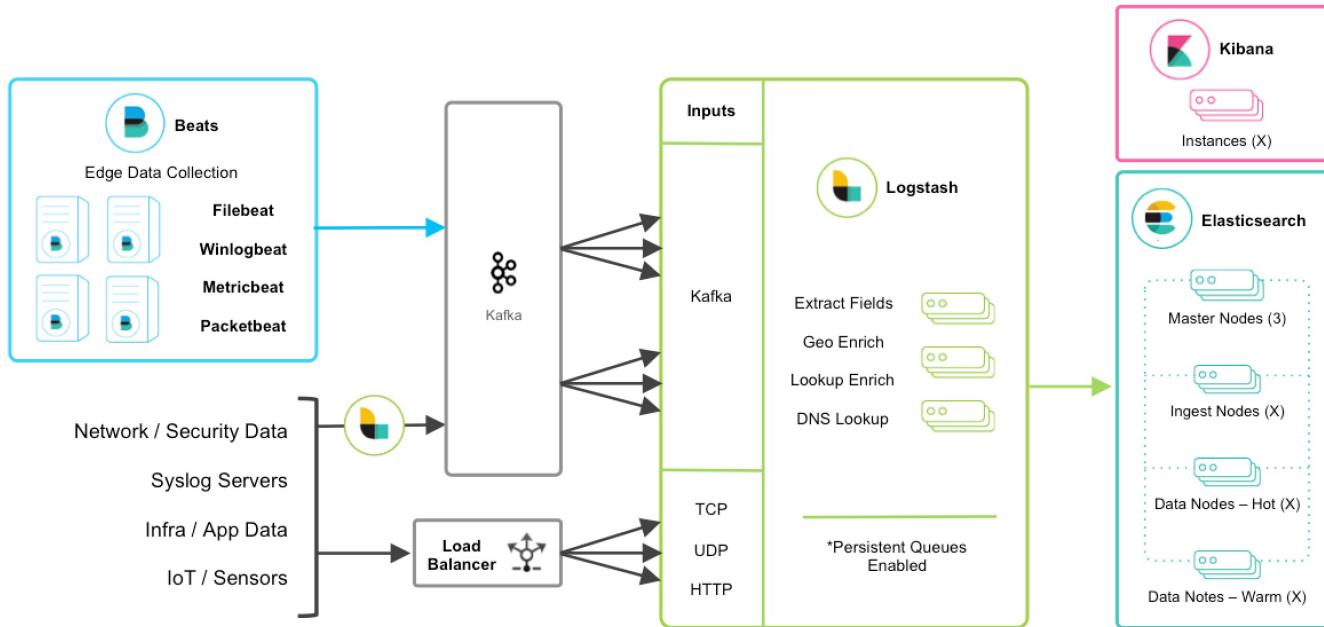
Streaming

Example Streaming Architecture



Advanced Ingestion Setup

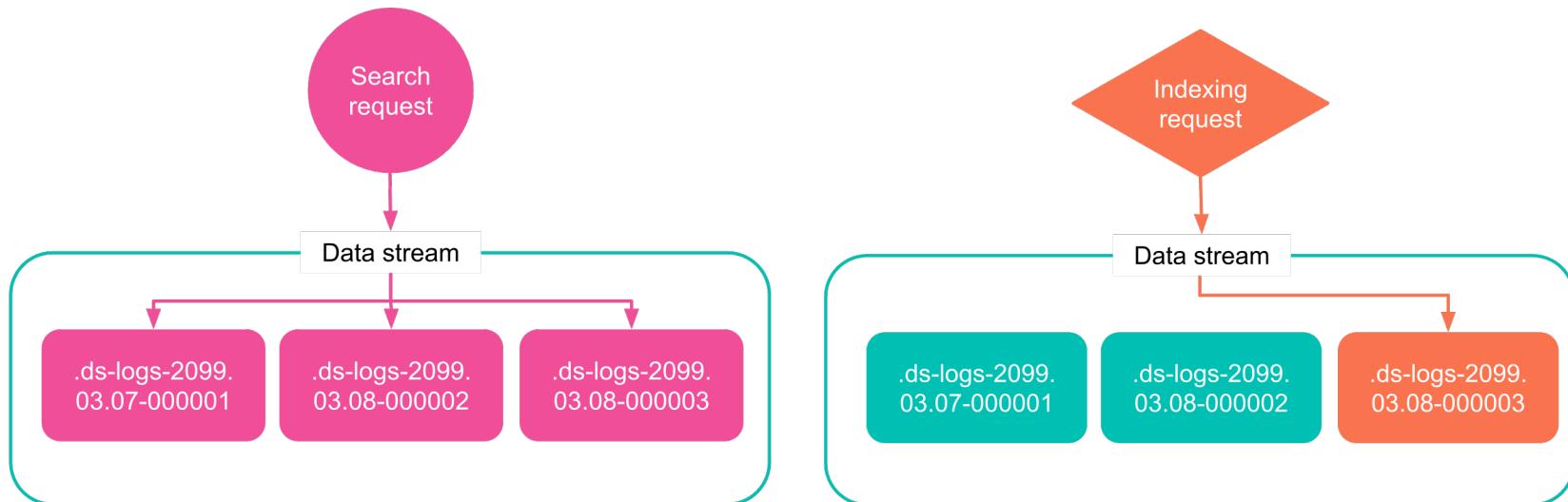
Additional On / Off Stack infrastructure



Elasticsearch Data Streams

- Append only
- Great for logs, metrics and other continuously generated data
- Can use Index Lifecycle Management (ILM) to automate index management
- ILM can help reduce costs significantly

Elasticsearch Data Streams



Homebrew Stream Processing

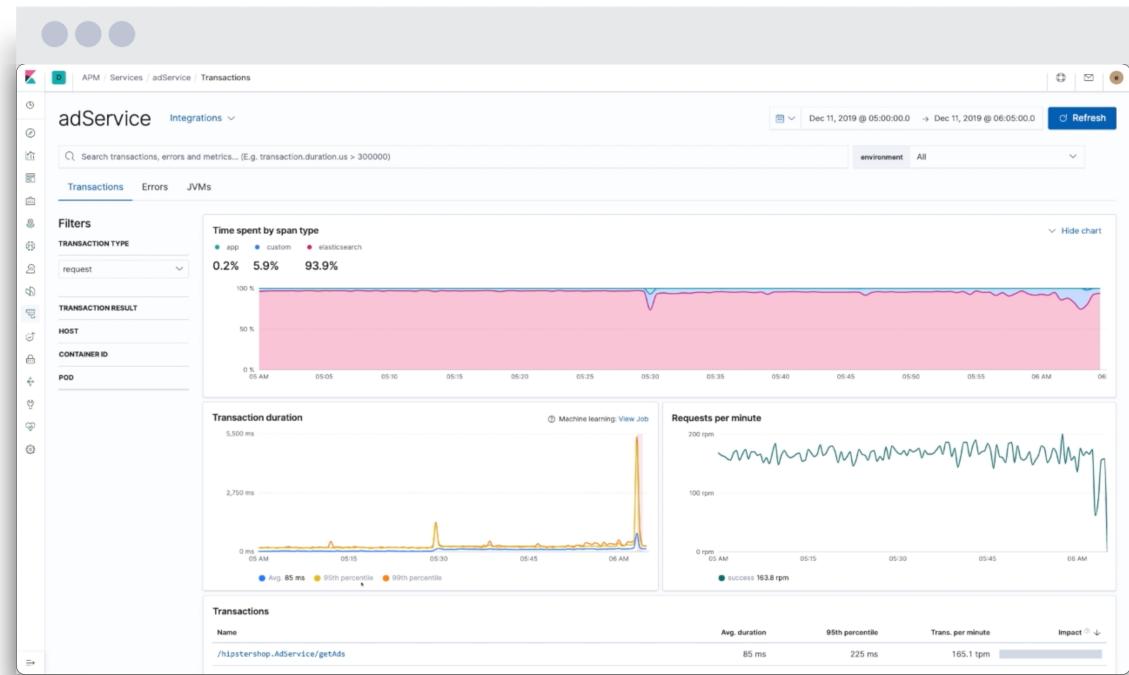
- Some teams don't wish to use *Big Data* tooling
- Bring your own code / tools

Advanced Stack Usage

APM & Machine Learning



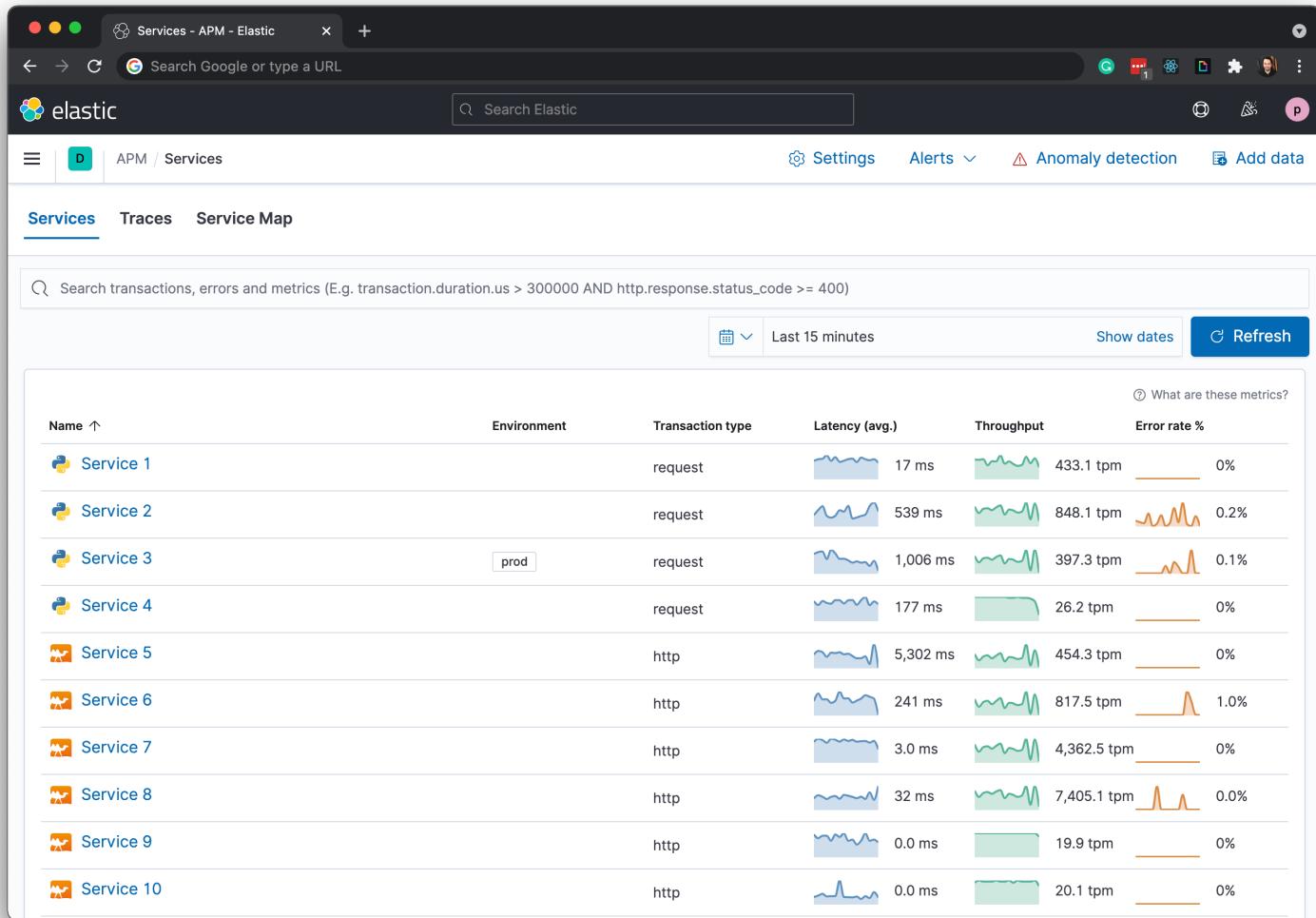
Elastic APM

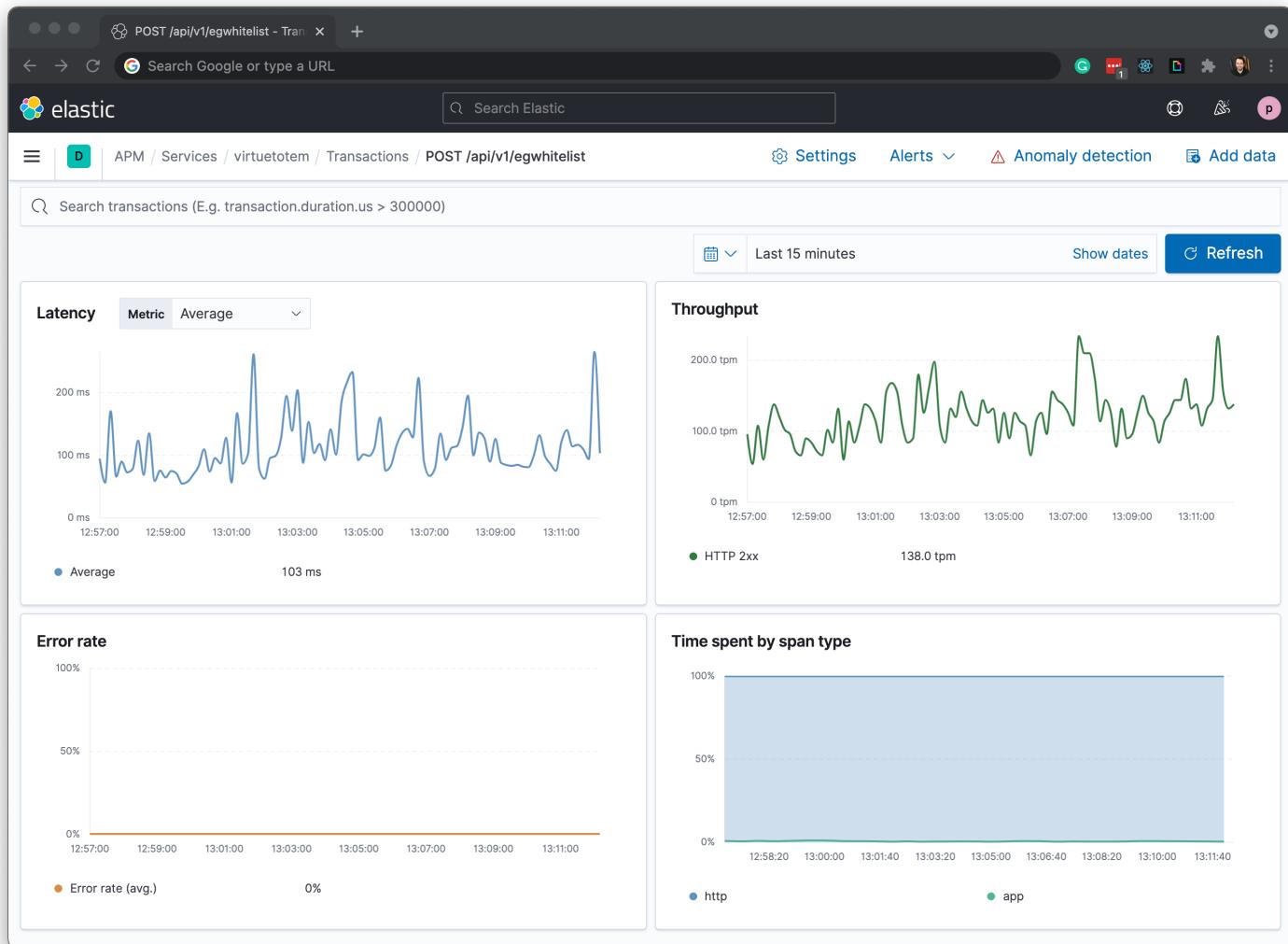


Elastic APM

Profile pipeline and stream processing components

- Check service and infrastructure uptime and health
- Service Maps & Distributed tracing
- Monitor TLS Cert expiry
- Alert when things go wrong (& create cases!)
- Data Engineering friendly
 - Supports integrations such as Kafka, RabbitMQ, ...
 - Programming language support – Java, Python, Go, ...
 - Community support for OCaml, Erlang and others





POST /api/v1/egwhitelist - Trans x +

Search Google or type a URL

elastic Search Elastic

APM / Services / virtuetotem / Transactions / POST /api/v1/egwhitelist

Settings Alerts Anomaly detection Add data

Latency distribution

Latency Range (ms)	Transactions
0 - 100	~1000
100 - 200	~200
200 - 1900	< 100

Trace sample

15 minutes ago | 45 ms (100% of trace) | POST http://virtuetotem.prod.endgamecloud.com/api/v1/egwhitelist?ori... | 200 OK | Python Requests (2.17)

Investigate View full trace

Timeline Metadata Logs

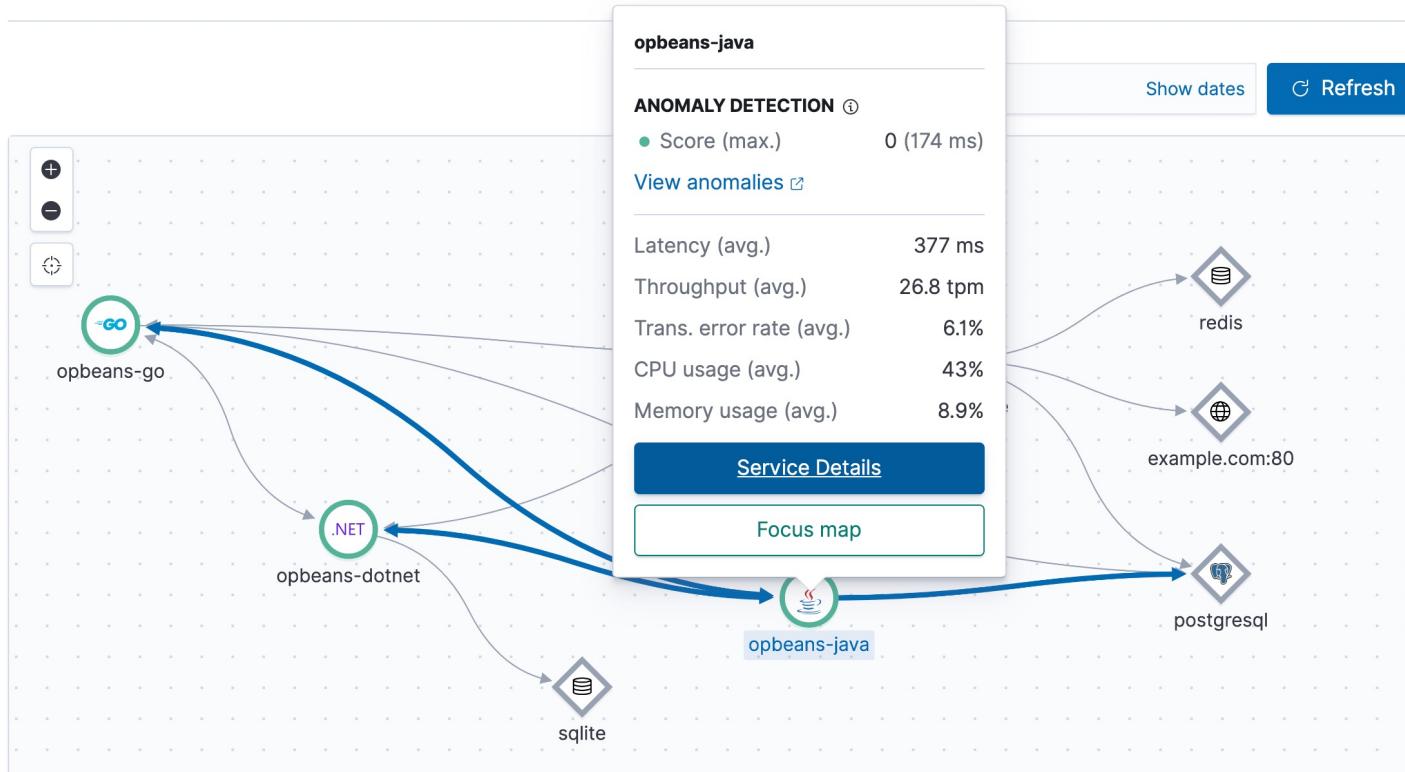
Type: virtuetotem http

HTTP 2xx POST /api/v1/egwhitelist 45 ms

GET cloud.security.elastic.co 44 ms async

Service Map

Environment All



Machine Learning

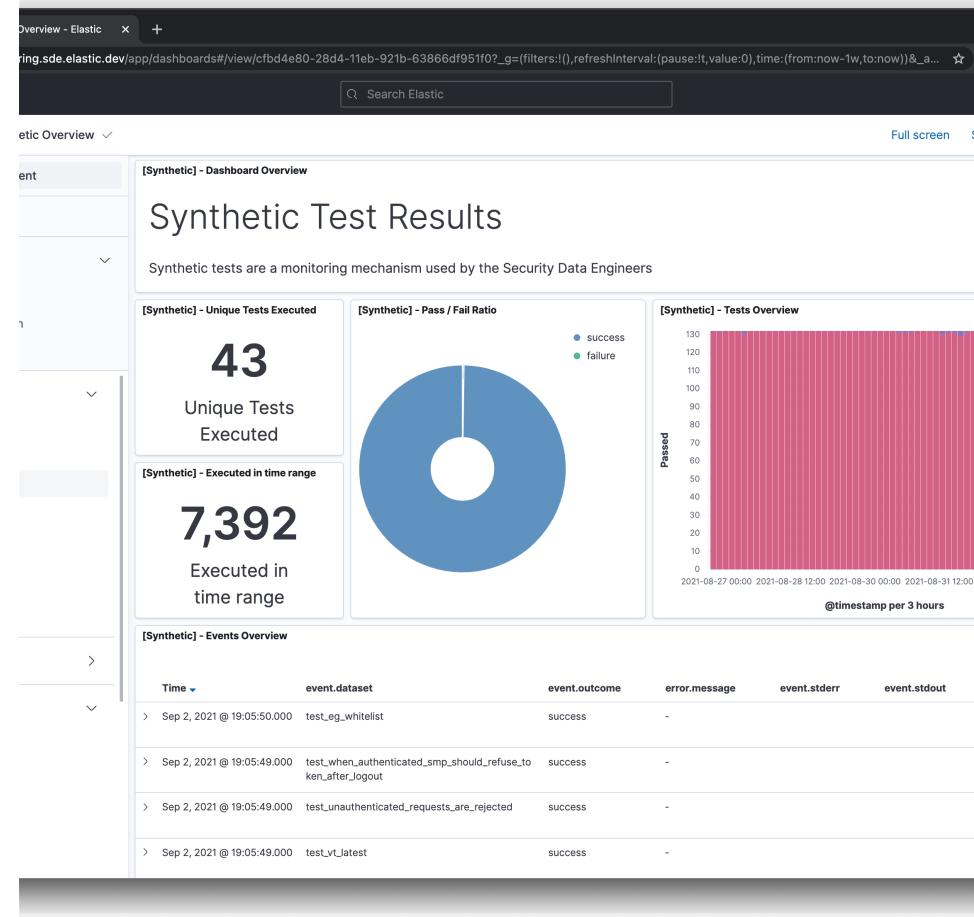
Unsupervised analysis of APM Data

- Highlights anomalies
- Deep insight into performance of individual service + trace
- Spot issues before they happen

Capturing Tests

Monitoring Synthetic events

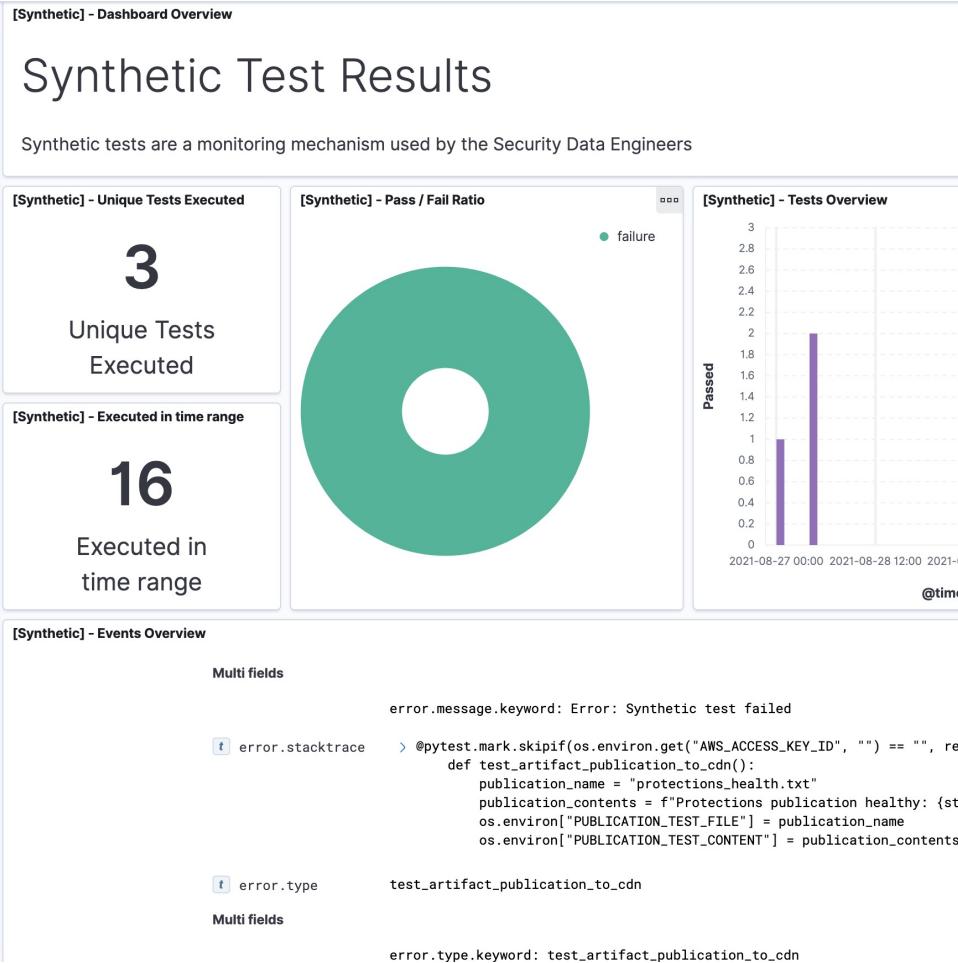
- Run tests on a temporal cadence
- Insert test output into Elasticsearch
- Create simple visualizations and dashboards



Components Fail

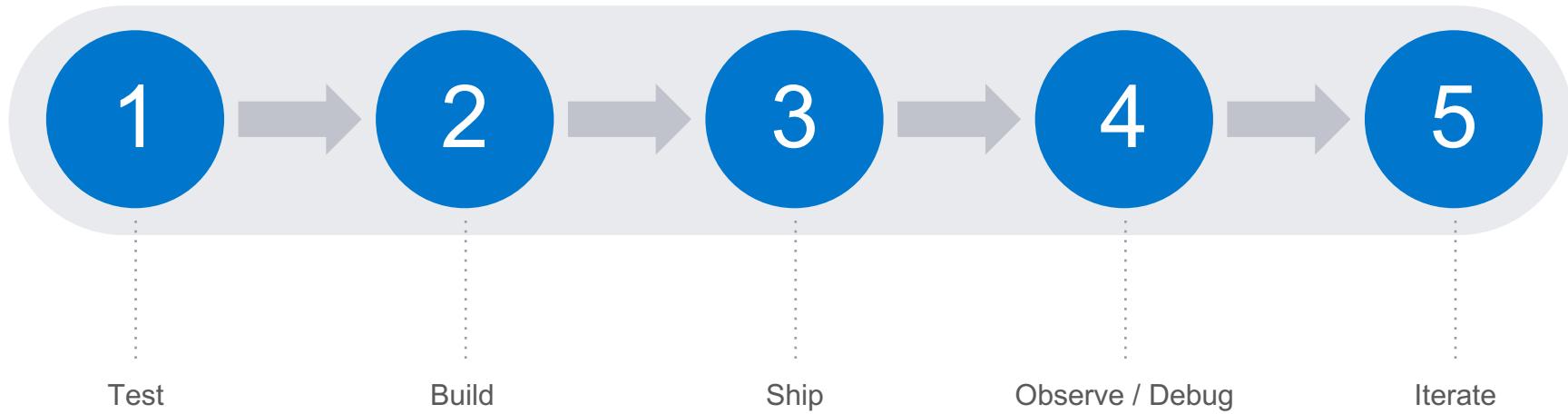
Context matters

- Determine trends in services not meeting our SLAs



Summary

Feedback loops are key



Takeaways

- Elastic Stack helps operate data workloads
- Understand data integration and processing
- Use APM & ML for enhanced observability
- Lots of great integrations available
(Public Cloud, JMS, JDBC, Vault, Zookeeper, ...)

Thanks!

<https://github.com/pjhampton>