

# Building a Cross Cloud Data Protection Engine

Richard Conway  
CEO and Founder  
@azurecoder

Sandy May  
Lead Data Engineer  
@spark\_spartan

# Speaker Bio

**Richard Conway** - @azurecoder

Microsoft Azure Most Valuable Professional

Microsoft Regional Director

UK Azure User Group Co-Organizer and Co-Founder

Data Science London Co-Organizer

Worldwide technology speaker

Passionate about big data, AI and security in Microsoft Azure

**DATA+AI SUMMIT EUROPE**

**#DataTeams #DataAISummit**



# Speaker Bio

**Sandy May** - @spark\_spartan

Databricks Champion

Data Science London Co-Organizer

Tech speaker across the UK

Passionate about Apache Spark, Databricks, AI,  
Data Security and Reporting platforms in  
Microsoft Azure

# Agenda

## Richard Conway

What is a Data Protection Engine and Why do we need it? Let's also look at some architecture

## Sandy May

Building a simple Data Protection Engine, from the ground up to cover GDPR and give the business a starting point



# Data Protection Engine Overview



**DATA+AI SUMMIT EUROPE**

**@azurecoder @spark\_spartan**

**#DataTeams #DataAISummit**

# What is the problem?

- **GDPR & CCPA fines can be in billions \$ now**
  - British Airways €204m July 2019 – 500,000 effected customers
  - Highest theoretical fine \$21b based on 4% 2019 revenue
- **Off the shelf products are expensive**
  - With Slow delivery roadmaps that you can't control
  - You still must pay to run them in cloud = more \$\$\$
- **Products don't mitigate risk, you still own risk**
  - You are responsible to run products over your data
  - Some products won't even own liability for bugs in their software
  - Most don't "detect" PII within data



# Should we Build or Buy?

## Build

- Own the IP
- Prioritise the features you want
- Built for your use case
- No licence fees
- Use your core technology

## Buy

- May have track record
- Bugs fixed by vendor
- Features not thought about by business
- Service Level Agreements

# Business Needs

- Run as part of Data Pipeline and ad-hoc
- Track lineage of Data Protected
  - Use a metadata store for all transformations
- Support Pseudonymisation, Anonymisation & Generalization
  - Re-identification required from Pseudonymisation
  - Joining datasets required from Pseudonymisation
- Allow Pseudonymisation Tokens to be migrated to another solution

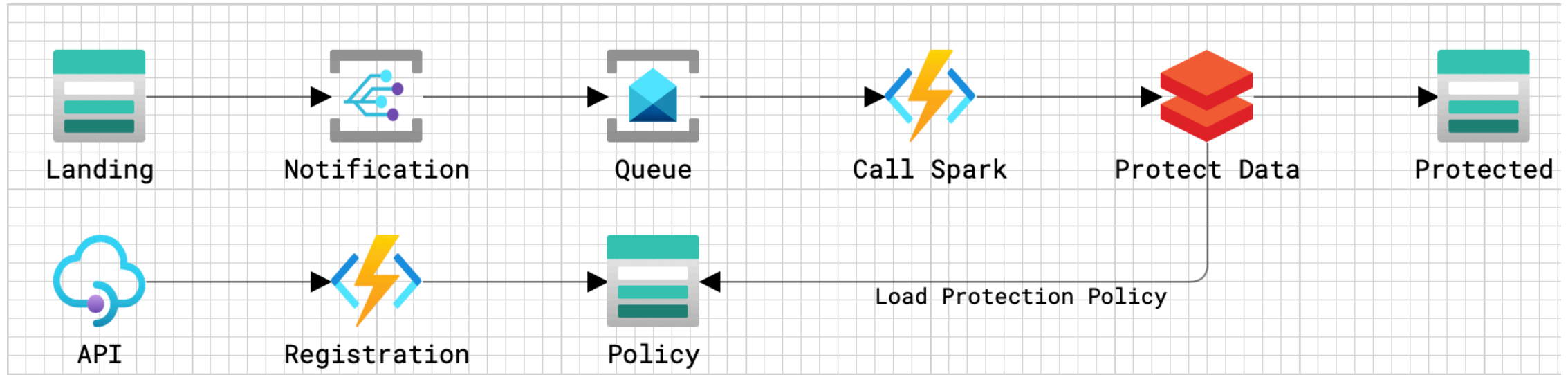


# Key Design Decisions

- Support to Run On-Premise and Cloud
- Use Native tools in Azure and AWS
- Token Vault consistency and auditability
- Single Reporting Platform
- Metadata driven

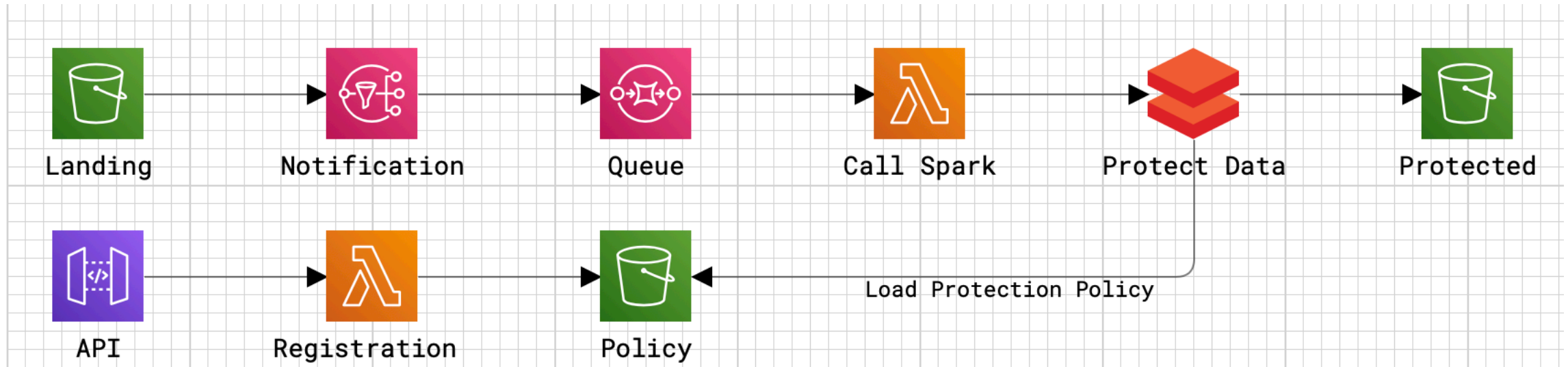


# Architecture - Azure





# Architecture - AWS



# Config Driven Design

```
c:\Users> rusarg > Desktop > DtreePMML.xml
1 <?xml version="1.0" encoding="UTF-8" standalone="yes"?>
2 <PMML xmlns="http://www.dmg.org/PMML-4_2" xmlns:data="http://jpmml.org/jpmml-model/InlineTable" version="4.2">
3 <Header copyright="Copyright (c) 2020 ivnard" description="Decision Tree Model">
4 <Application name="JPMML-SparkML" version="1.5.4"/>
5 <Timestamp>2020-03-21T13:30:39Z</Timestamp>
6 </Header>
7 <DataDictionary>
8 <DataField name="BAD" optype="categorical" dataType="integer">
9 <Value value="0"/>
10 <Value value="1"/>
11 </DataField>
12 <DataField name="MORTDUE" optype="continuous" dataType="double"/>
13 <DataField name="YOJ" optype="continuous" dataType="double"/>
14 <DataField name="DEROG" optype="continuous" dataType="double"/>
15 <DataField name="DELINQ" optype="continuous" dataType="double"/>
16 <DataField name="CLAGE" optype="continuous" dataType="double"/>
17 <DataField name="NINQ" optype="continuous" dataType="double"/>
18 <DataField name="CLNO" optype="continuous" dataType="double"/>
19 <DataField name="DEBTINC" optype="continuous" dataType="double"/>
20 </DataDictionary>
21 <TreeModel functionName="classification" splitCharacteristic="binarySplit">
22 <MiningSchema>
23 <MiningField name="BAD" usageType="target"/>
24 <MiningField name="MORTDUE"/>
25 <MiningField name="YOJ"/>
26 <MiningField name="DEROG"/>
27 <MiningField name="DELINQ"/>
28 <MiningField name="CLAGE"/>
```

# Let's Build it!



**DATA+AI SUMMIT EUROPE**

**@azurecoder @spark\_spartan**

**#DataTeams #DataAISummit**

# Summing Up



**DATA+AI SUMMIT EUROPE**

**@azurecoder @spark\_spartan**

**#DataTeams #DataAISummit**



# Future Work

- Machine Learning PII Detection
- K-Anonymity
- Batching Service
  - Databricks cluster only runs one job with 3-5 minute spin up time
  - Delta Lake ensures ACID transaction on requests originating from a single cluster
- De-centralize the solution
  - Allow individual data teams to control their own data protection and pay for their usage
  - Maintain a central reporting solution for business
  - Consideration needs to be given to joins across tokenized data

# Conclusions

- Building can be quick and secure
- Prioritise your own business needs
- Can be used as a stop gap while you create a service for an off the shelf product
- No false promises of protection, you control all



**databricks**



Thanks for listening!  
Questions?



# Feedback

Your feedback is important to us.

Don't forget to rate  
and review the sessions.

