# Building Cloud Infrastructure

Aaron Davidson

databricks

# Who am I?

- Early Databricks engineer (4 years)
- Apache Spark committer & PMC member
- Worked on a lot of things @ DB
- Most recently, cloud infrastructure
  - Helping eng produce efficient, secure, and reliable software.

databricks

# What is Databricks?

- Big Data & Machine Learning in the Cloud
  - Yes - our customers are *data scientists* and *data engineers*
- Thinking about getting into self-driving cars
- Yes, we have some Go and Rust code, but prefer FP

databricks

# Databricks Product

- People love Spark, but:
    - How do I get and maintain a Spark cluster?
    - How do I configure that cluster?
    - How do I run jobs reliably and periodically?
    - How do I interface with Spark?

**Operations**

**Usability**

databricks

# Databricks Product

- People love Spark, but:
    - How do I get and maintain a Spark cluster?
    - How do I configure that cluster?
    - How do I run jobs reliably and periodically?
    - How do I interface with Spark?

    **Operations**

    **Usability**

- Enter Databricks…

databricks

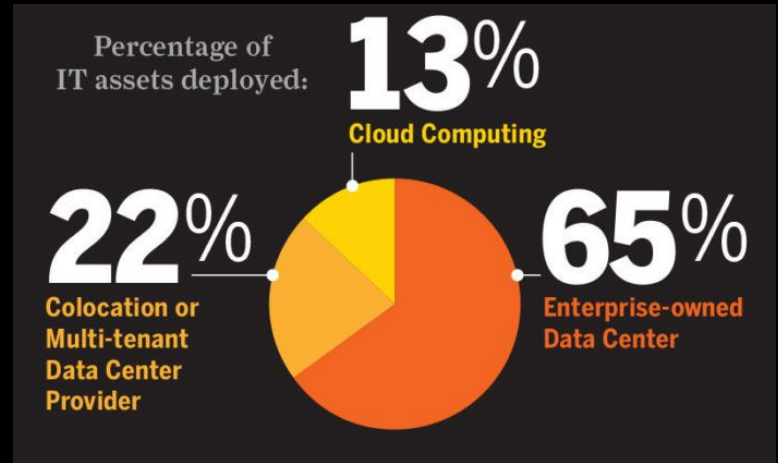# Databricks Product

- People love Spark, but:
    - How do I get and maintain a Spark cluster?
    - How do I configure that cluster?
    - How do I run jobs reliably and periodically?
    - How do I interface with Spark?

    **Operations**

    **Usability**

- Enter Databricks…
- What hardware do we have?

# What does it mean to be a Cloud Company?

- Most money is **still** in on-premise, but trend is towards Cloud.
- "Enterprise:" Financial institutions, government, health care, etc.
- Berkeley & probably Stanford, too



Percentage of IT assets deployed:

**13%** Cloud Computing

**22%** Colocation or Multi-tenant Data Center Provider

**65%** Enterprise-owned Data Center

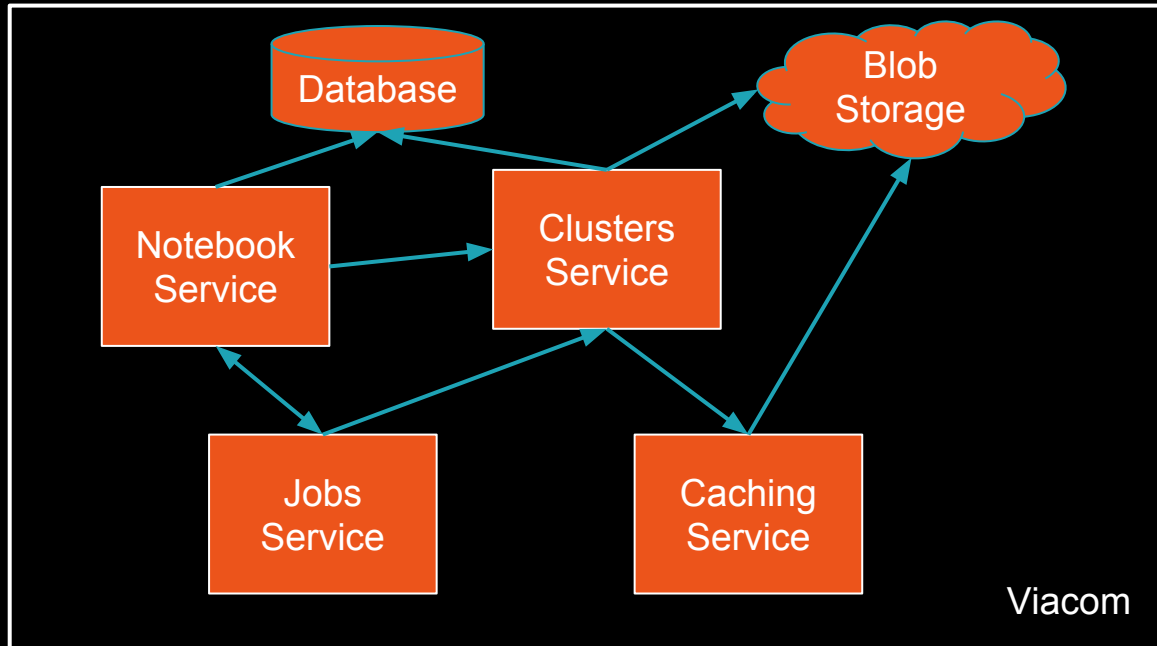databricks

# What does it mean to be a Cloud Company?

- Infrastructure in the Cloud (vs on-prem infrastructure):
    - **Infrastructure is dynamic** -- provisioning new hardware in O(minutes) rather than O(months).
    - No operations team, but high-level primitives provided instead.
        - Storage (DBs, blob storage), networking (routing/firewalls), etc
- Running Software as a Service (vs on-prem appliance) means:
    - We operate the product on behalf of our customers.
    - Often, the software we run is **multitenant**.
    - **Update often** -- deliver features and fixes faster than 3/6/12 months

databricks

# In this talk

- We'll use a real-life motivating example from Databricks to talk about building a **cloud service**.
- Focus on three major aspects:
    - Scaling out a multitenant service
    - Updating services safely
    - Deploying the infrastructure to run our service.

databricks

# Databricks Community Edition

- In The Beginning, Databricks provided a single-tenant product
- Easier:
  - Security
  - Isolation
  - Selling
- But:
  - Costly
  - Failures



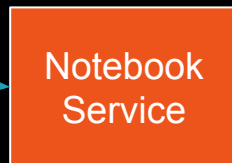databricks

# Databricks Community Edition

- We wanted to make a free, multitenant version
- Use-cases: people playing around with Spark, training/classes, MOOCs (now: all new customers)
- Problems:
  - How do we scale our single-tenant services out?
  - How do we update when there is constant usage?
  - How do we maintain this larger, more dynamic infrastructure?

# The Notebook Service

- Collaborative notebook UI
  - Users mainly edit their own notebooks, but sometimes want to collaborate
  - Collaboration requires merging changes from multiple users in real-time.
- Originally: ~10 concurrent users.
- Now: Training of 500 people -- or a 50,000-person MOOC!
- How do we scale this service out?

databricks

# The Notebook Service

# Service Replication

# Replication: Logical Stickiness

# Replication: Logical Stickiness

# Replication: Logical Stickiness

# Replication: Stateless

# Replication: Stateless

# Replication: Stateless

# Replication: Stateless



How do we deal?
- Push logic into database
- Take fine-grained locks

# Replication: User/Session Stickiness

# Replication: User/Session Stickiness

# Replication: User/Session Stickiness

# Replication: User/Session Stickiness

# Replication: User/Session Stickiness



TCP-sticky load balancer
Easy to find -- probably default!

HTTP-sticky load balancer
Cookie-based -- a bit more
complicated, but also common

Database

Notebook Service 1

Notebook Service 2

Notebook Service 3

{char: "s"}
{char: "e"}
{char: "l"}

Load Balancer

read ntbk

databricks

# Replication: User/Session Stickiness



TCP-sticky load balancer
Easy to find -- probably default!

HTTP-sticky load balancer
Cookie-based -- a bit more complicated, but also common

| Pros | Cons |
| --- | --- |
| + Easy to find<br>+ Built-in fault recovery | - Only supports single-flow/user locality<br>- Failures may be harder to reason about |

# Service replication: How to decide?

- Review:
  - Stateless replication: Simplest
    - Simplest ("best") replication model, hardest to program against
  - Session/user stickiness
    - Particularly common replication model -- well-supported by tooling
  - Logical/tenant stickiness
    - Most complicated ("worst") replication model, easiest to program against
- Considerations:
  - Higher is better, but have to start thinking from beginning.
  - If not, then the last will be the only option (that's exactly what we did for notebooks!)

databricks

# Service replication: How to implement?

- VM-level: Cloud providers have TCP & HTTP load balancers:
  - Static or scalable pool of machines registered with a port & protocol.
  - Health checking mechanism to remove machines from routable pool.
- Container-level: YMMV; Kubernetes also provides TCP- and HTTP-level load balancing, between containers.

# Service replication: How to implement?

- Tenant-stickiness?

- Need a consistent, highly-available leader election store
  - ZooKeeper, consul, etcd (Googlers: Chubby)

- Need an HTTP load balancer
  - Probably nginx or go -- not recommended to build your own, in JVM

# Recap: Databricks Community Edition

- We wanted to make a free, multitenant version
- Use-cases: people playing around with Spark, training/classes, MOOCs (now: all new customers)
- Problems:
    - ✓ How do we scale our single-tenant services out?
    - – How do we update when there is constant usage?
    - – How do we maintain this larger, more dynamic infrastructure?

databricks

# Service updates

- Can leverage our earlier work in service replication to perform updates without downtime.
- Update strategies:
    - The ol' off 'n' on
    - Blue-green
    - Rolling
    - Traffic control

databricks

# Service updates: Blue/green

# Service updates: Blue/green



databricks

# Service updates: Blue/green

| Pros | Cons |
|---|---|
| + Easy to implement<br>+ Can work with single replica | - Unused infra<br>- All-or-nothing -- bugs exposed immediately |

Version 2

Notebook Service 2

Notebook Service 1

Notebook Service 3

Load Balancer

Load Balancer/DNS

databricks

# Service updates: Rolling update

# Service updates: Rolling update



databricks

# Service updates: Rolling update

# Service updates: Rolling update



Notebook Service V2

Notebook Service V2

Notebook Service V2

Load Balancer

databricks

# Service updates: Rolling update



| Notebook Service V2 |
| Notebook Service V2 | Notebook Service V2 |
| Load Balancer |

| Pros | Cons |
| --- | --- |
| + Gradual roll out | - Coarse-grained |
| + All infra used | |

databricks

# Service updates: Traffic control

# Service updates: Traffic control

# Service updates: Traffic control

# Service updates: Traffic control

# Service updates: Traffic control

| Pros | Cons |
|------|------|
| + Google-scale quality control<br>+ Simple extension: shadowing traffic | - Requires complicated load balancer |

**Gaining traction:**
Envoy & Istio starting to add support

# Update strategy: How to decide?

- **-** Review:
  - – Blue/green
    - – Useful for stateful applications
    - – Useful for acceptance testing
    - – Complicated roll-out procedure
  - – Rolling update
    - – Most common -- simple roll-out procedure
  - – Traffic control
    - – Best-in-class -- requires complicated load balancer
- **-** Considerations:
  - – Design with at least one updates strategy in mind and you can keep downtime minimal, even for unreplicated services.

# Update strategy: How to implement?

- VM-level: Cloud providers have (auto)scaling groups.
    - Create a new group for the new version.
    - For blue-green, switch DNS when tested.
    - For rolling update, have load balancer use both groups and increase/decrease replicas.
    - Netflix does this -- see Spinnaker
- Container-level: Kubernetes provides first-class support for rolling updates within one cluster, other stuff is as manual as VM case.

databricks

# Recap: Databricks Community Edition

- We wanted to make a free, multitenant version
- Use-cases: people playing around with Spark, training/classes, MOOCs (now: all new customers)
- Problems:
    - ✓ How do we scale our single-tenant services out?
    - ✓ How do we update when there is constant usage?
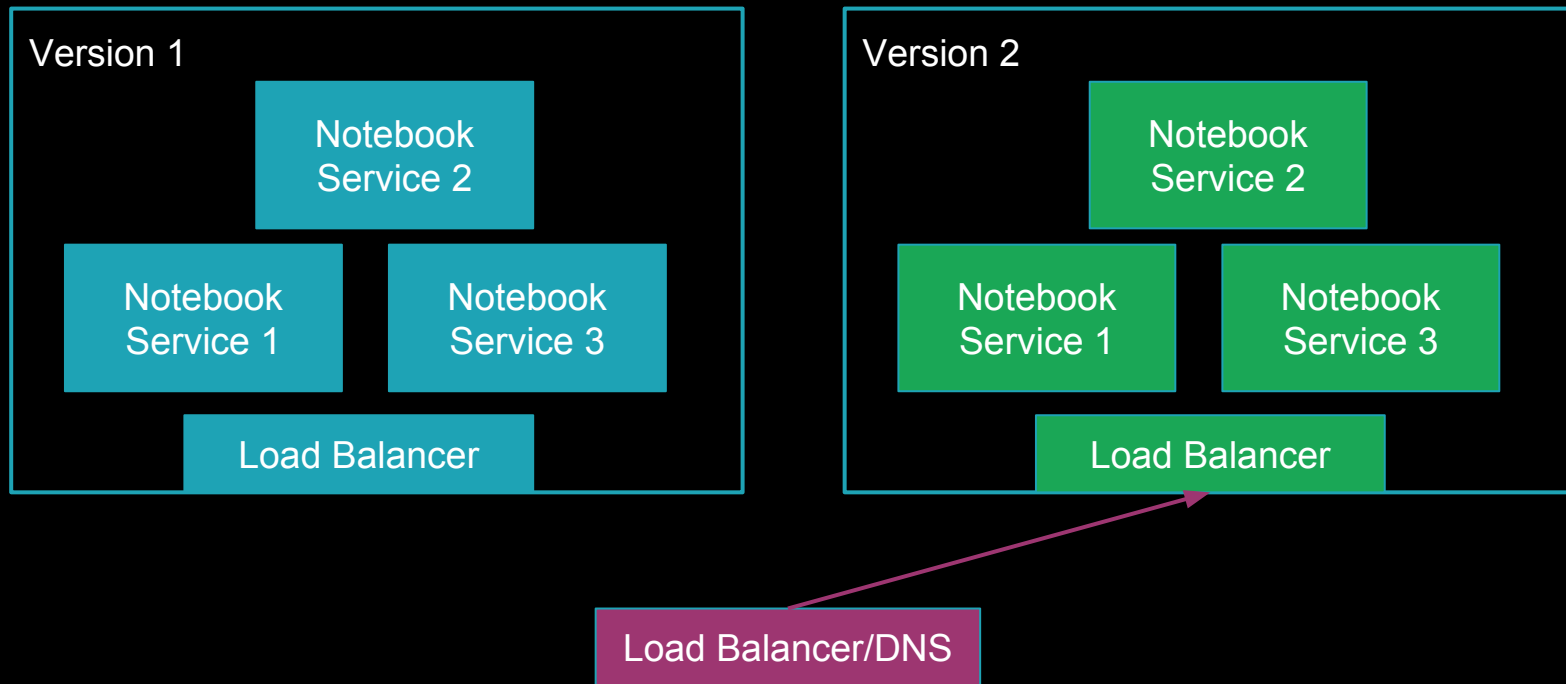    - – How do we maintain this larger, more dynamic infrastructure?

databricks

# Infrastructure as Code

- I want to provision 3 VMs for my Notebook Service.
- On-prem: Ask ops team for 3 machines, wait 1-3 months
- Cloud: **Launch Instance**

- Scenarios:
  - Scale out to 5 VMs.
  - VM crashes, need to replace it.
  - Change VM parameter (e.g., instance size)
  - Replicate environment to a new region.
  - Create a testing environment.
  - Security breach! Tear it all down and recreate everything.

databricks

# Infrastructure as Imperative Code

```
def createInfra():
  for i in range(3):
    ec2.createInstance(
      name = s"NotebookService-$i",
      type = "m4.xlarge")
```

- Scenarios:
  - ✖ Scale out to 5 VMs.
  - ✖ VM crashes, need to replace it.
  - ✖ Change VM parameter (e.g., instance size)
    Replicate environment to a new region.
    Create a testing environment.
    Security breach! Tear it all down and recreate everything.

databricks

# Infrastructure as Imperative Code

```
def createInfra(region):
  for i in range(3):
    ec2.createInstance(
      name = s"NotebookService-$i",
      type = "m4.xlarge",
      region = region)
```

- Scenarios:
  - ✖ Scale out to 5 VMs.
  - ✖ VM crashes, need to replace it.
  - ✖ Change VM parameter (e.g., instance size)
  - ✓ Replicate environment to a new region.
    Create a testing environment.
    Security breach! Tear it all down and recreate everything.

# Infrastructure as Imperative Code

```
def createInfra(region, accountId):
  for i in range(3):
    ec2.createInstance(
      name = s"NotebookService-$i",
      type = "m4.xlarge",
      region = region,
      accountId = accountId)
```

- Scenarios:
  - ✖ Scale out to 5 VMs.
  - ✖ VM crashes, need to replace it.
  - ✖ Change VM parameter (e.g., instance size)
  - ✓ Replicate environment to a new region.
  - ✓ Create a testing environment.
  - ✓ Security breach! Tear it all down and recreate everything.

databricks

# Infrastructure as Imperative Code

```
def createInfra(region, accountId, oldCount, newCount):
  for i in range(oldCount, newCount):
    ec2.createInstance(
      name = s"NotebookService-$i",
      type = "m4.xlarge",
      region = region,
      accountId = accountId)
```

- Scenarios:
  - ✓ Scale out to 5 VMs.
  - ✖ VM crashes, need to replace it.
  - ✖ Change VM parameter (e.g., instance size).
  - ✓ Replicate environment to a new region.
  - ✓ Create a testing environment.
  - ✓ Security breach! Tear it all down and recreate everything.

databricks

# Infrastructure as Imperative Code

- Scenarios:
  - ✓ Scale out to 5 VMs.
  - ✖ VM crashes, need to replace it.
  - ✖ Change VM parameter (e.g., instance size).
  - ✓ Replicate environment to a new region.
  - ✓ Create a testing environment.
  - ✓ Security breach! Tear it all down and recreate everything.

- Problems:
  - Specific: Each scenario needs new code, new parameters. Not necessarily shared between use-cases, either (e.g., create a database)
  - Stateful: Correctness requires either maintaining state, writing state resolution logic, or having a human enter the state.
  - Fallible: Did you spot the incorrect error handling?

databricks

# Infrastructure as Declarative Code

```
[{ kind: "EC2::Instance",
   type: "m4.xlarge",
   name: "NotebookService-0",
   region: "oregon",
   accountId: 1234567,
}, … ]
```

Declarative Deployer

Notebook Service-0

Notebook Service-1

Notebook Service-2

- Scenarios:
  - ✓ Scale out to 5 VMs.
  - ✓ VM crashes, need to replace it.
  - ✓ Change VM parameter (e.g., instance size).
  - ✓ Replicate environment to a new region.
  - ✓ Create a testing environment.
  - ✓ Security breach! Tear it all down and recreate everything.

databricks

# Infrastructure as Declarative Code

- Scenarios:
  - ✓ Scale out to 5 VMs.
  - ✓ VM crashes, need to replace it.
  - ✓ Change VM parameter (e.g., instance size).
  - ✓ Replicate environment to a new region.
  - ✓ Create a testing environment.
  - ✓ Security breach! Tear it all down and recreate everything.
- Benefits: State, API, and error handling are all managed for us
  - Difficult to manage large, dynamic infrastructure due to duplication. (One solution here is to introduce a layer of templating)
  - Needs an implementation of "Declarative Deployer"
    - All cloud providers have a native way of doing this (e.g., CloudFormation)
    - Terraform is a cloud semi-agnostic tool
    - Quilt?

databricks

# Recap: Databricks Community Edition

- We wanted to make a free, multitenant version
- Use-cases: people playing around with Spark, training/classes, MOOCs (now: all new customers)
- Problems:
  - ✓ How do we scale our single-tenant services out?
  - ✓ How do we update when there is constant usage?
  - ✓ How do we maintain this larger, more dynamic infrastructure?

databricks

# Summary

- Cloud infrastructure is dynamic
  - Replicate multitenant services for scale-out
  - Automate deployment (imperatively or declaratively)
  - Leverage cloud provider abstractions (VMs, load balancers, databases)

- Software as a Service allows us to move quickly
  - Deliver updates on weekly cadence rather than 3/6/12-monthly
  - Reduce friction of use by taking over operational burden
  - Just make sure your updates aren't breaking things *too* often!

databricks

# Thank you!

We're hiring -- come intern with us!

Aaron Davidson - aaron@databricks.com

Try Community Edition:
https://databricks.com/try-databricks

# Appendix: Container Engines (Kubernetes)

databricks

# What problem are we trying to solve?

- I want to run my code on a remote server.
- How do I get my code there?
    - What about my code's dependencies (e.g., library A)?
    - What about my code's system dependencies (e.g., curl or ntp)?
- How do I know what's going on?
    - Logging?
    - SSHing into the machine?
- How do I update my code? How do I roll back?

databricks

# World V1: Ansible and "bare-metal"

- I want to run my code on a remote server.
- How do I get my code there?
    - Script which copies my JAR and any dependent jars.
    - Script also can install dependencies on target host.
- How do I know what's going on?
    - SSH in and find out.
- How do I update my code? How do I roll back?
    - Rerun script (how to undo dependencies?)
- Problems:
    - Script is not very general! New one per service.
    - Have to manually place services on hosts (what about node failure?)

databricks

# World V2: Ansible and Docker

- I want to run my code on a remote server.
- How do I get my code there?
    - I build a Docker container which contains all my dependencies!
    - I run a script which starts that script.
- How do I know what's going on?
    - SSH in and find out.
- How do I update my code? How do I roll back?
    - Rerun script -- dependencies inside container so can roll back.
- Problems:
    - Script is now pretty general, service-specific stuff is in container.
    - Still have to manually place services on hosts (node failures)

databricks

# World V3: Kubernetes (w/Docker)

- I want to run my code on a remote server.
- How do I get my code there?
    - I build a Docker container which contains all my dependencies!
    - I ask Kubernetes to find somewhere to put that container.
- How do I know what's going on?
    - I ask Kubernetes for logs or to SSH into the container directly.
- How do I update my code? How do I roll back?
    - I ask Kubernetes to do a rolling update.
- Problems:
    - Kubernetes replaces my custom script entirely
    - Kubernetes deals with placement of containers within a cluster, and with node failure.

databricks

# Other Kubernetes Features

- In addition to managing containers, Kubernetes helps with:
  - Exposing services to the outside world via Load Balancers
  - Maintaining a fixed set of replicas of a node.
  - Health checking and restarting services (provided service-specific health checks).
  - Managing network-attached storage.
  - Providing cross-cloud abstractions.
  - (And more!)
- Similar systems: DC/OS, Docker Swarm, Google's Borg

databricks

# Container Engines: Unsolved Problems

- Solid, authn/authz inter-service networking
    - Envoy & istio approach problem from proxy layer
    - Calico approaches problem from network layer (BGP!)
- Geo-replicated (multi-cluster) services
- Easy-to-use logical stickiness abstraction (e.g., notebooks)

# Appendix: Terraform
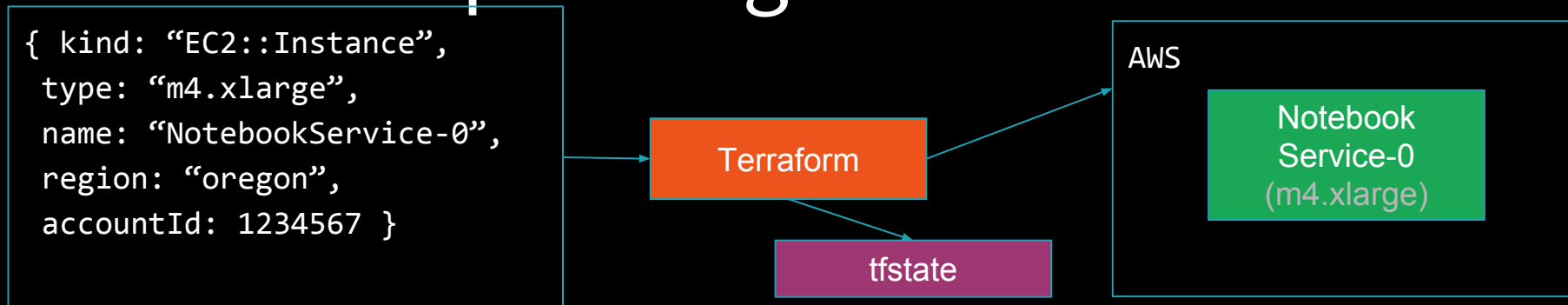
# Terraform Operating Model

```
{ kind: "EC2::Instance",
  type: "m4.xlarge",
  name: "NotebookService-0",
  region: "oregon",
  accountId: 1234567 }
```

Terraform

AWS

- Input: Template, state file, and cloud resources
- Output: Plan of how to converge state

databricks

# Terraform Operating Model

```
{ kind: "EC2::Instance",
  type: "m4.xlarge",
  name: "NotebookService-0",
  region: "oregon",
  accountId: 1234567 }
```

Terraform

tfstate

AWS

Notebook
Service-0
(m4.xlarge)

- Input: Template, state file, and cloud resources
- Output: Plan of how to converge state

# Terraform Operating Model

```
{ kind: "EC2::Instance",
 type: "m4.2xlarge",
 name: "NotebookService-0",
 region: "oregon",
 accountId: 1234567 }
```

Terraform

tfstate

AWS

Notebook
Service-0
(m4.xlarge)

- Different properties require different change procedures.
  - Changing EC2 VM instance size requires tearing down and recreating.
  - Changing RDS database instance size requires just restarting.

databricks

# Terraform Operating Model

```
{ kind: "EC2::Instance",
  type: "m4.xlarge",
  name: "NotebookService-100",
  region: "oregon",
  accountId: 1234567 }
```

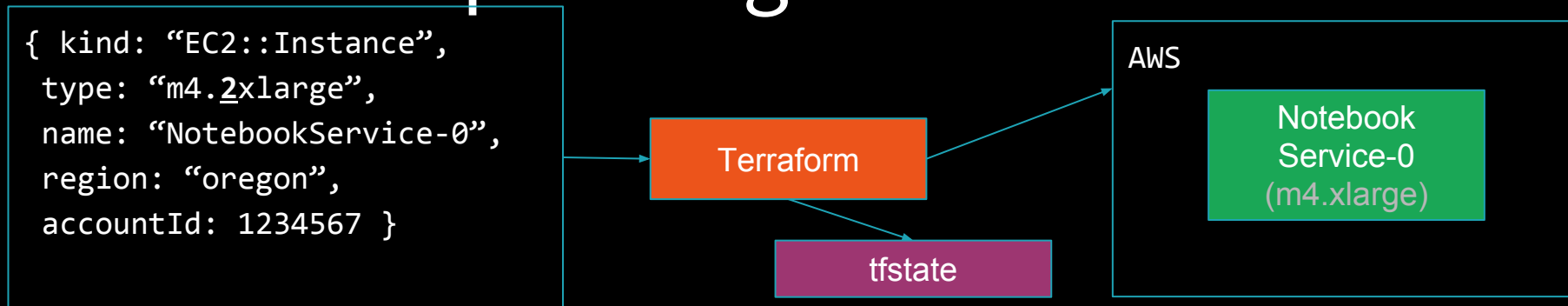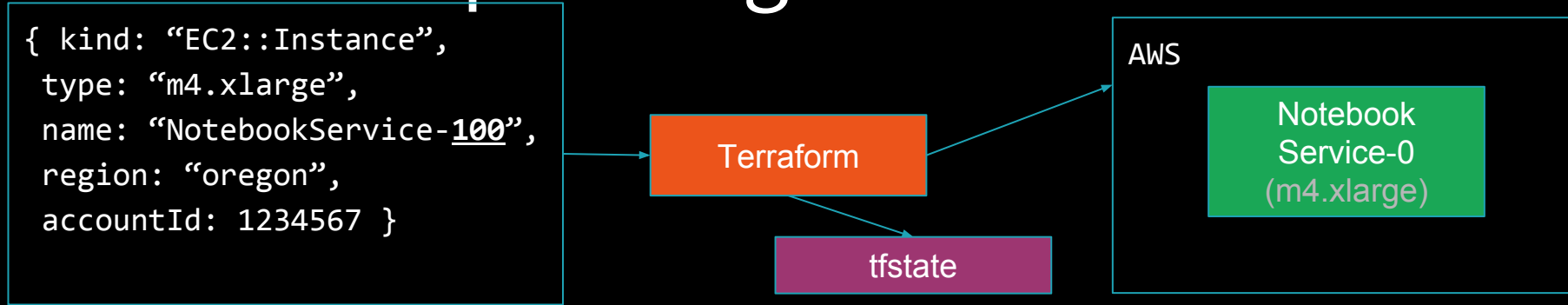Terraform

tfstate

AWS

Notebook
Service-0
(m4.xlarge)

- State file used so Terraform knows when it should delete objects.
- Otherwise, we would just create a second instance and keep the old one around.

databricks

# Declarative Deploy: Unsolved Problems

- Cloud agnostic terminology & semantics is elusive
- Declaring different classes of resources (e.g., cloud provider versus Kubernetes objects) requires different systems
- Enacting a certain change may require several intermediate templates
- No standard for templating.

databricks