facebook

# Managing MySQL at Scale

**Pradeep Nayak & Junyi (Luke) Lu**
Production Engineers - MySQL Infra

# Agenda

# Terminology

# Terminology

What's a instance ?

What's a shard ?

What's a replicaset ?

# Instance

foobar.prn:3307

foobar.prn:3309

foobar.prn:3306
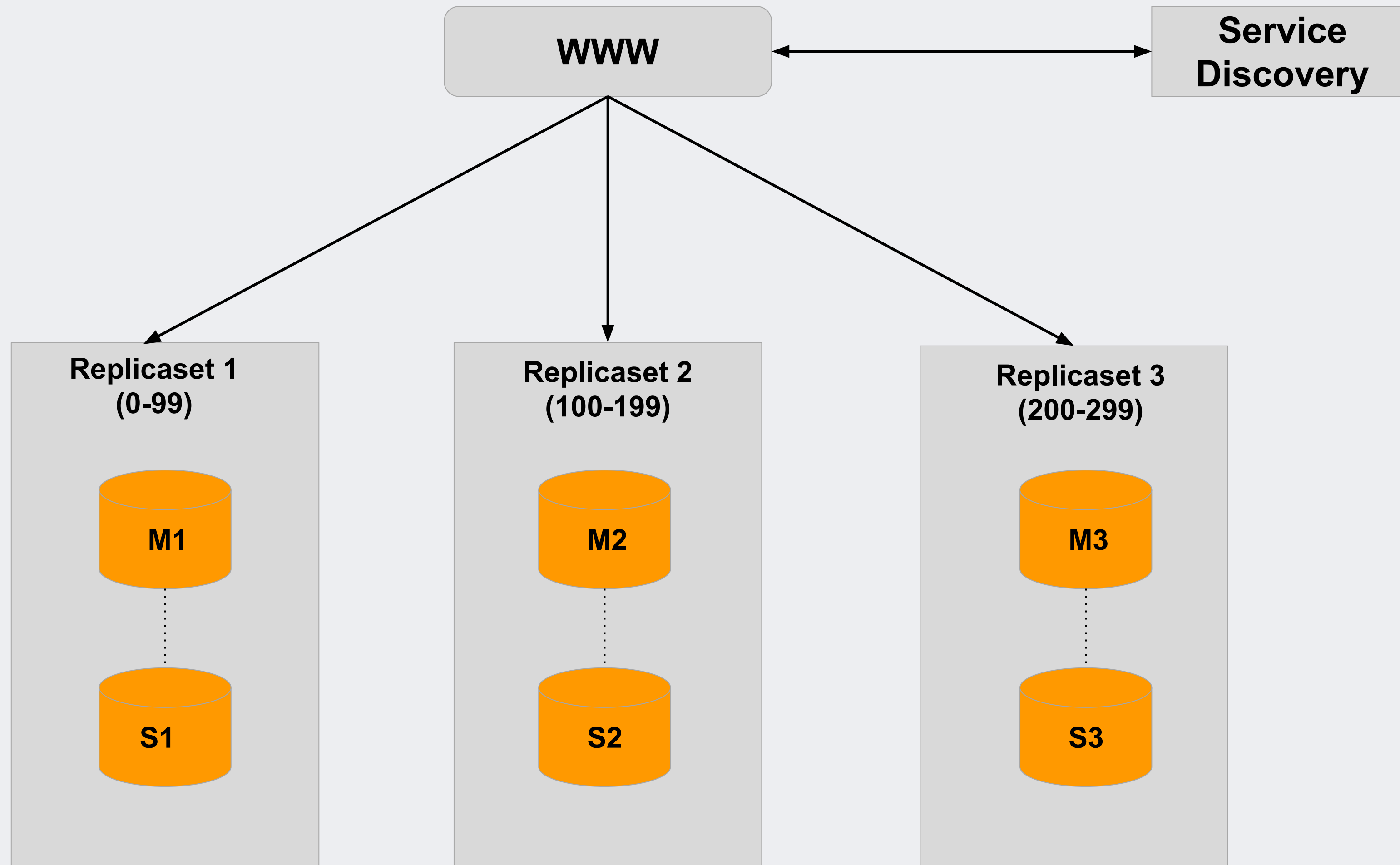
foobar.prn

3306

3307

3309

# Shard

db.helloworld
db.12345
db.44365

# Service Discovery

| Shard ID | Replicaset | Master | Slave |
|---|---|---|---|
| 0-99 | Replicaset 1 | db1234.prn1:3306 | db1234.frc1:3306 |
| 100-199 | Replicaset 2 | db4567.ftw1:3306 | db4567.prn1:3307 |
| 200-299 | Replicaset 3 | db1234.atn1:3306 | db1234.frc2:3308 |

# Service Discovery

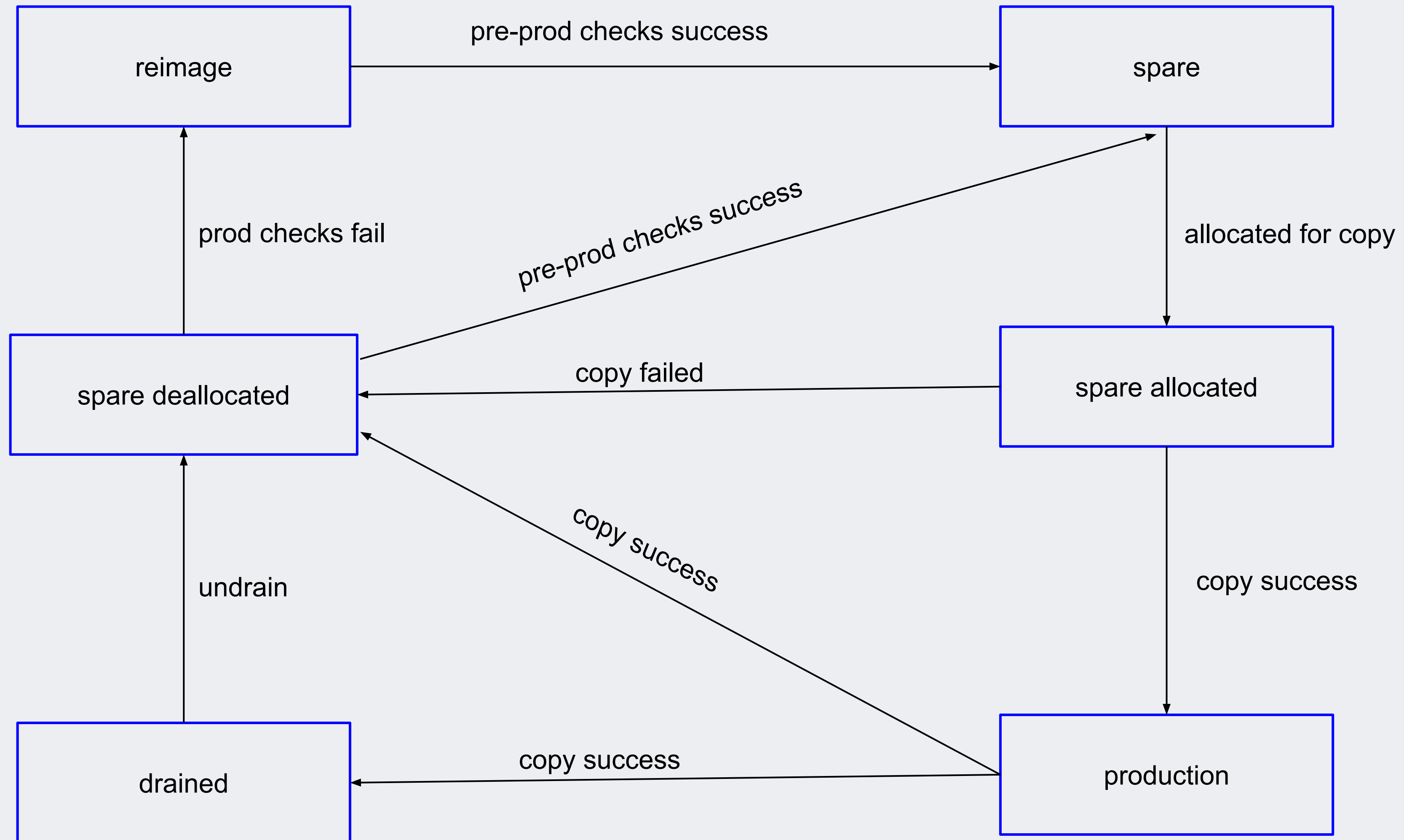| Shard ID | Replicaset | Master | Slave |
| --- | --- | --- | --- |
| 0-99 | Replicaset 1 | db1234.prn1:3306 | db1234.frc1:3306 |
| 100-199 | Replicaset 2 | db4567.ftw1:3306 | db1234.prn1:3309 |
| 200-299 | Replicaset 3 | db1234.frc2:3308 | db1234.atn1:3309 |

# Lifecycle of an instance

# Lifecycle of an instance

. States are production, spare, spare allocated, spare deallocated, reimage, drained

. Metadata includes instance properties like name, port, mysql rpm version, state etc

. A mysql shard hosts metadata of all instances in the fleet

# Lifecycle of an instance

# Lifecycle of an instance

. Each state has its own processor to do the work

. Each state has a queue where work is queued

. Runs constantly scanning the fleet

# Instance Migration

# Clone an instance

Use case of cloning a production MySQL instance

- Replace a broken instance/host

- Move data around for maintenance

- Balancing host utilization

# MPS Copy

A workflow system that manages the requests for cloning
MySQL instances

- spare allocation
- set up MySQL config
- copy data
- replication
- validation
- bring it online & remove the old instance if necessary

# MPS Copy - Allocation

Choose the best slot for the instance based on its footprint

- Disk usage
- CPU utilization
- Failure domain

| Allocation | Setup | Migration | Replication | Validation | Registeration |

# MPS Copy - Setup

Turn up an empty instance using the right configuration

- Install the right RPM version

- Bootstrap the correct directory

- Generate the right my.cnf based on its use case

- Make sure the empty instance is connectable

| Allocation | Setup | Migration | Replication | Validation | Registeration |

# MPS Copy - Data Migration

We support three different ways of cloning a production instance

- Physical copy: xtrabackup and myrocks_hotbackup
- Logical copy:
  - mysqldump
  - Restore from backup

| Allocation | Setup | Migration | Replication | Validation | Registeration |
|---|---|---|---|---|---|

# MPS Copy - Replication

- Setup replication
  - From current production master
  - From Binlog Server
- Catchup

| Allocation | Setup | Migration | Replication | Validation | Registeration |
|------------|-------|-----------|-------------|------------|---------------|

# MPS Copy - Validation

If the data migration is a logical one, we will use snapshot based checksum to verify the correctness of data by comparing to its current master

| Allocation | Setup | Migration | Replication | Validation | Registeration |
|:---:|:---:|:---:|:---:|:---:|:---:|

# MPS Copy - Service Registration

Register the new instance in our service discovery system so that the MySQL users will be able to notice this new instance that has been recently turned up

| Allocation | Setup | Migration | Replication | Validation | Registeration |
|------------|-------|-----------|-------------|------------|---------------|

# Online Shard Migration

# Online Shard Migration

Another fundamental piece of our infra to control the growth of each MySQL instance
- Instance can grow beyond the host level limit
  - Too big
  - Too hot

# Online Shard Migration (OLM)

Key concept: Move the data of a shard into other smaller/cooler instance through logical migration and register the new address into the service discovery system

# OLM



| Shard | Replicaset |
|-------|-----------|
| db.1 | **mysql.replicaset.2** |
| db.2 | mysql.replicaset.1 |
| db.3 | mysql.replicaset.3 |

# OLM Processor

Workflow management for massive OLM operations

- Conflict solver

- Picking the best destination replicaset

- Kickoff the actual move

- Proper retry and cleanup

# Balancing

# Balancing

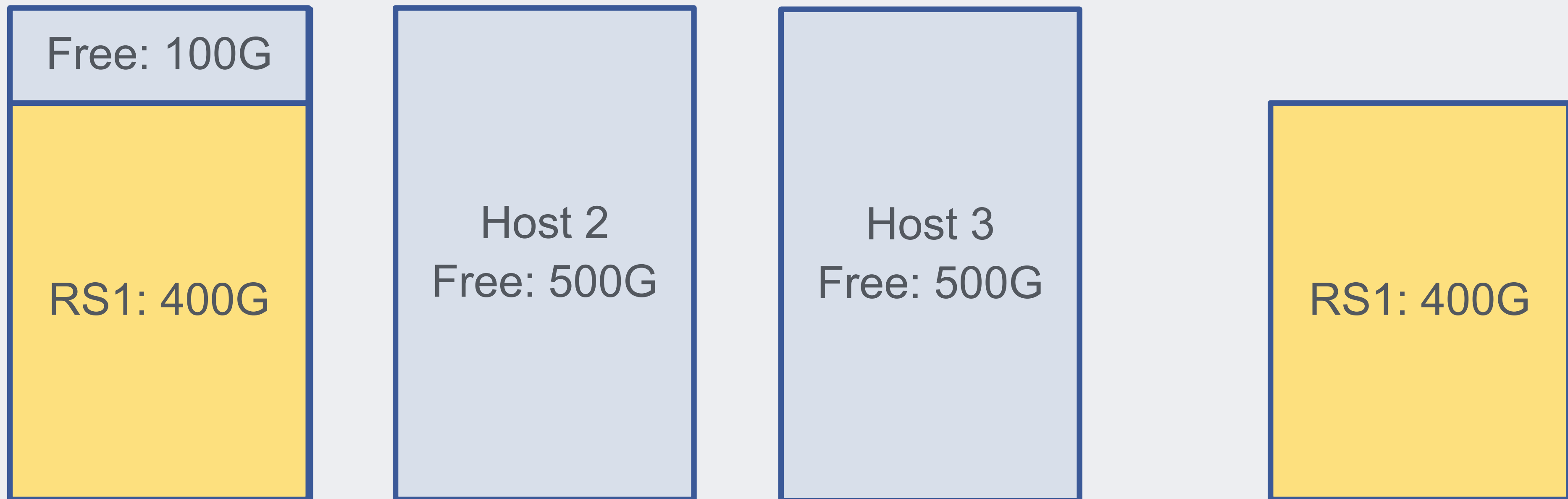Find the right place for the workload in order to achieve maximum sustainable resource utilization
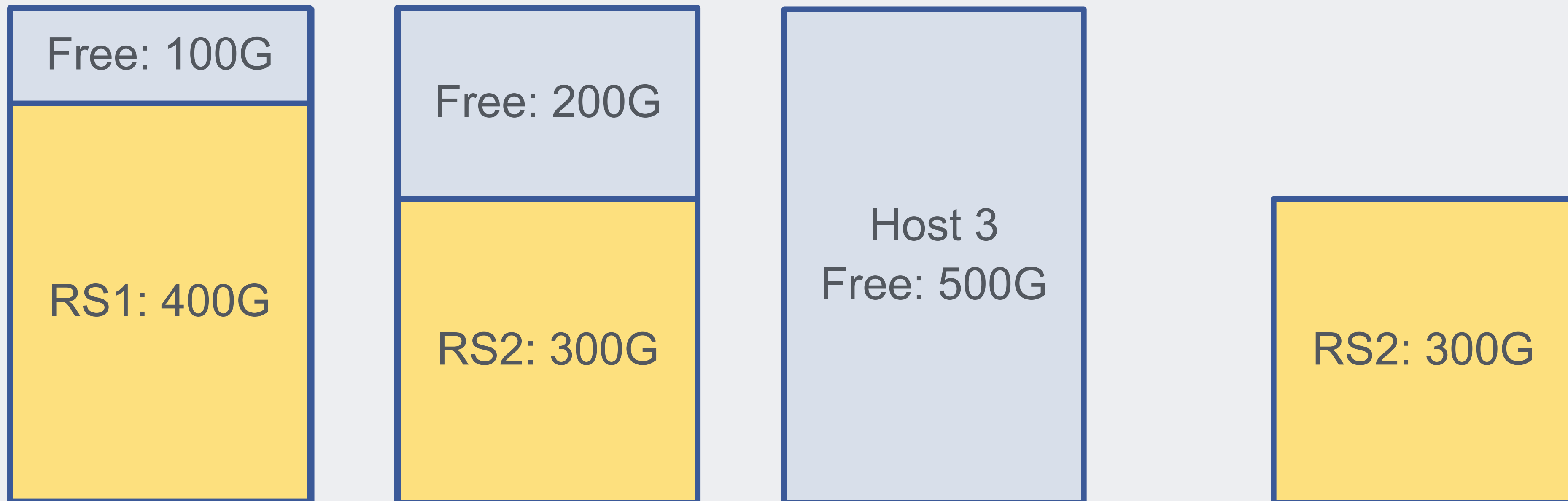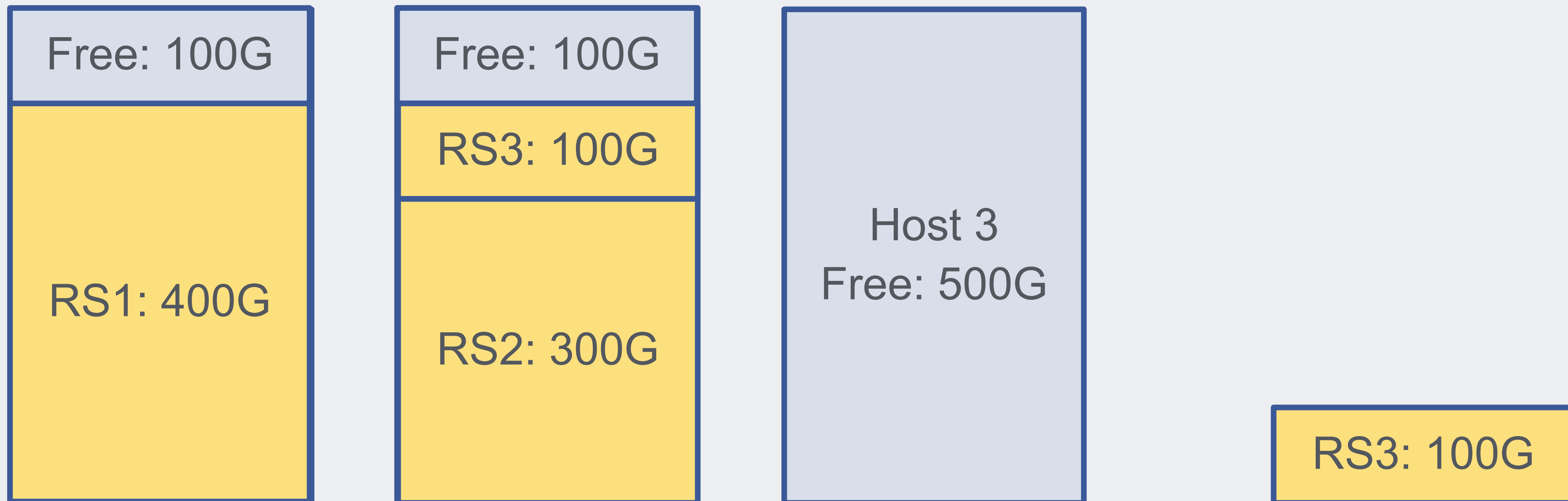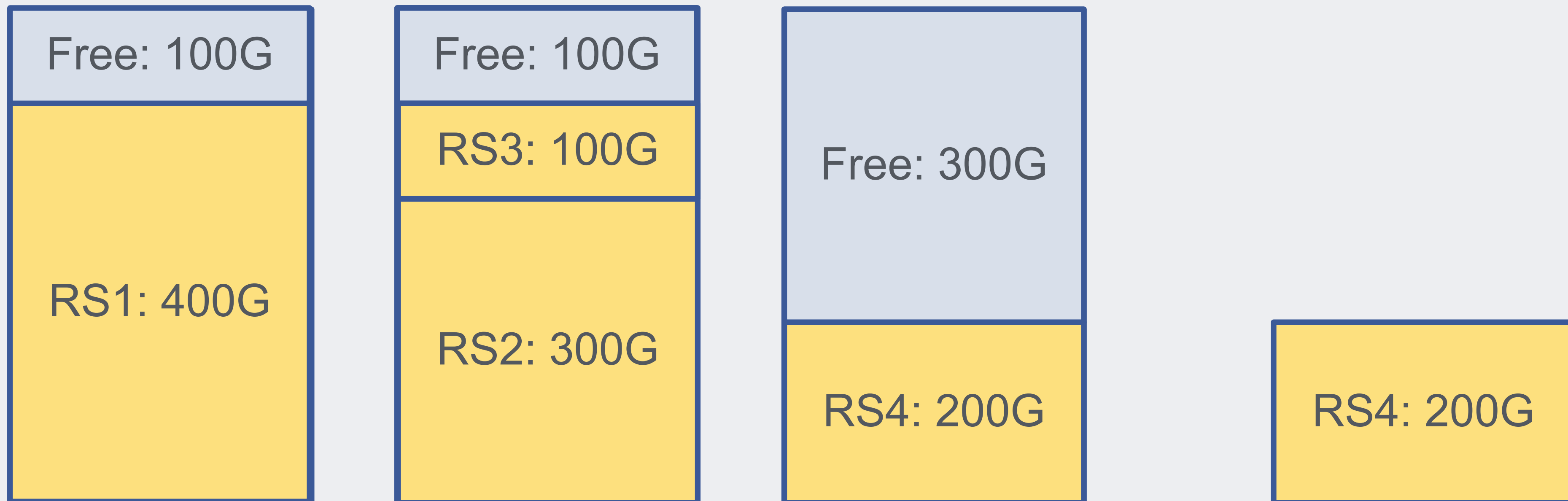
# Poor Stacking

# Poor Stacking

# Poor Stacking

# Proper Stacking

RS3: 100G

RS1: 400G

Free: 200G

RS2: 300G

Host 3
Free: 500G

RS3: 100G

# Proper Stacking

# Carve the Shape



RS1: 400G

RS2: 200G

RS3: 200G

RS4: 200G

Host: 500G

Total: 1000G

# Poor Shape



RS1: 400G — Free: 100G

RS3: 200G / RS2: 200G — Free: 100G

RS4: 200G — Free: 300G

X3

# Carve the Shape

RS1: 400G

RS2: 300G

RS3: 200G

100G

RS4: 100G

Host: 500G

Total: 1000G
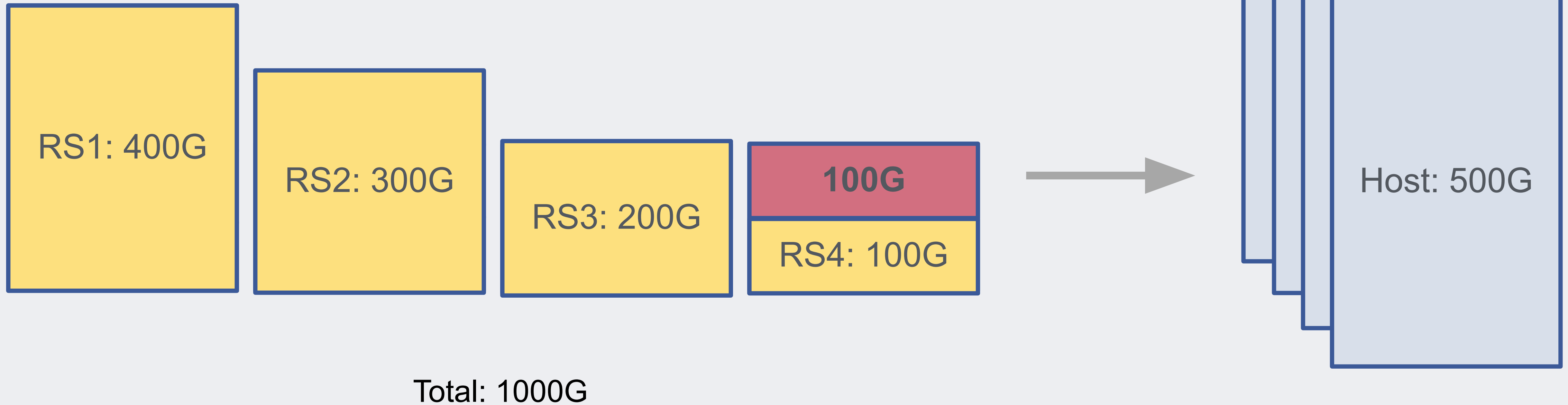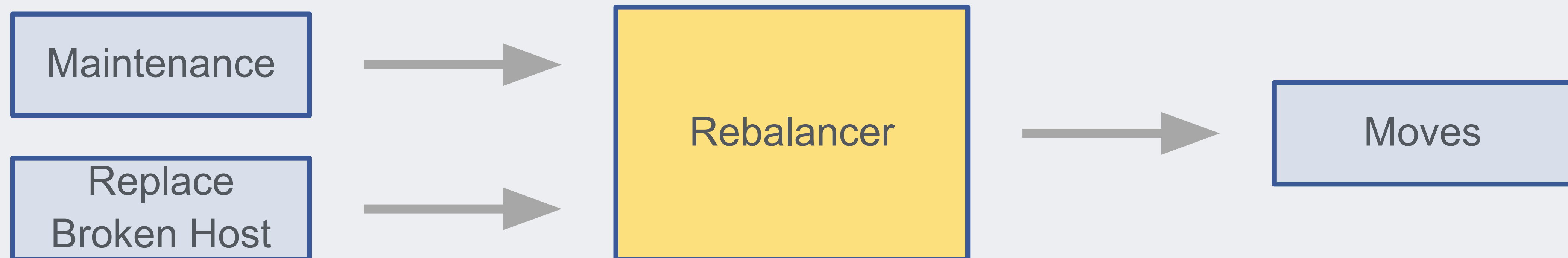
# Rebalancer

Goal: Find the best slot for hosting the given workload profile and reduce the imbalance score across the fleet to be minimum

# Rebalancer - Challenges

Multiple balancing factors

- CPU/Memory/Disk usage

- Fault domain spreading
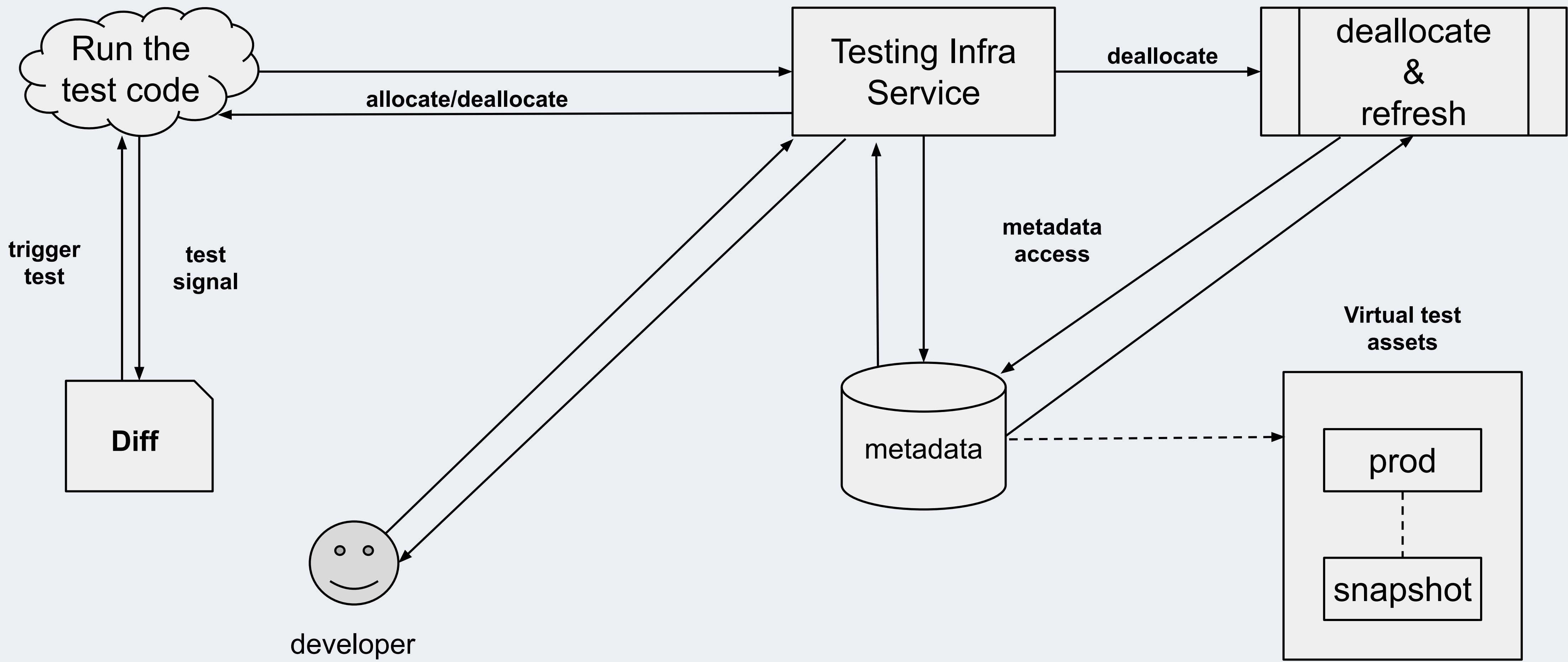
- MySQL vs LBU anti affinity

# Testing Infrastructure

# Testing Infra for Automations

. Lots of Automation code handling critical components of Infra

. UnitTests are good but mock backend connection

. Need to test end to end

# Testing Infra Goals

. build and canary packages based on the change

. provide signals at diff time for developer

. production like setup, but isolated environment

. iterate quickly with confidence

# Q&A