

Forget everything you knew about Swift Rings

(HERE'S EVERYTHING YOU NEED TO KNOW ABOUT RINGS)

Your Ring Professors

- Christian Schwede
 - Principal Engineer @ Red Hat
 - STAND UP GUY
- CLAY GERRARD
 - Programmer @ SwiftStack
 - LOUD & ANNOYING

Rings 201

- **Why Rings Matter**
- **What are Rings**
- **How Rings Work**

- HOW TO USE RINGS
- NINJA SWIFT RING TRICKS
- MOAR AWESOME STUFF

Swift 101

Looking for more general intro to Swift?

- Swift 101: <https://youtu.be/vAEU0Ld-GIU>
- Building webapps with Swift:
<https://youtu.be/4bhdqtLLCiM>
- Stuff to read:
https://www.swiftstack.com/docs/introduction/openstack_swift.html

OR SIX

~~One Ring~~ To Rule Them All

Swift Dperceptors

CAN BE A WILD RIDE

RING MASTERS



Ring Features

- DEVICES & SERVERS
- ZONES
- REGIONS
 - MULTI-REGION
 - CROSS-REGION
 - LOCAL-REGION
- STORAGE POLICIES



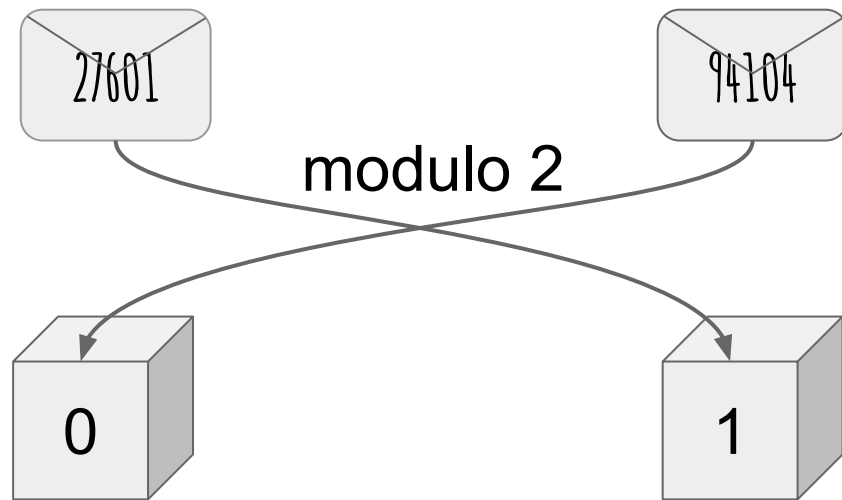
Swift's Rings use Simple Concepts

Consistent Hashing introduced
by Karger et al. at MIT in 1997

THE SAME YEAR HTTP/1.1 IS SPECIFIED IN RFC 2616

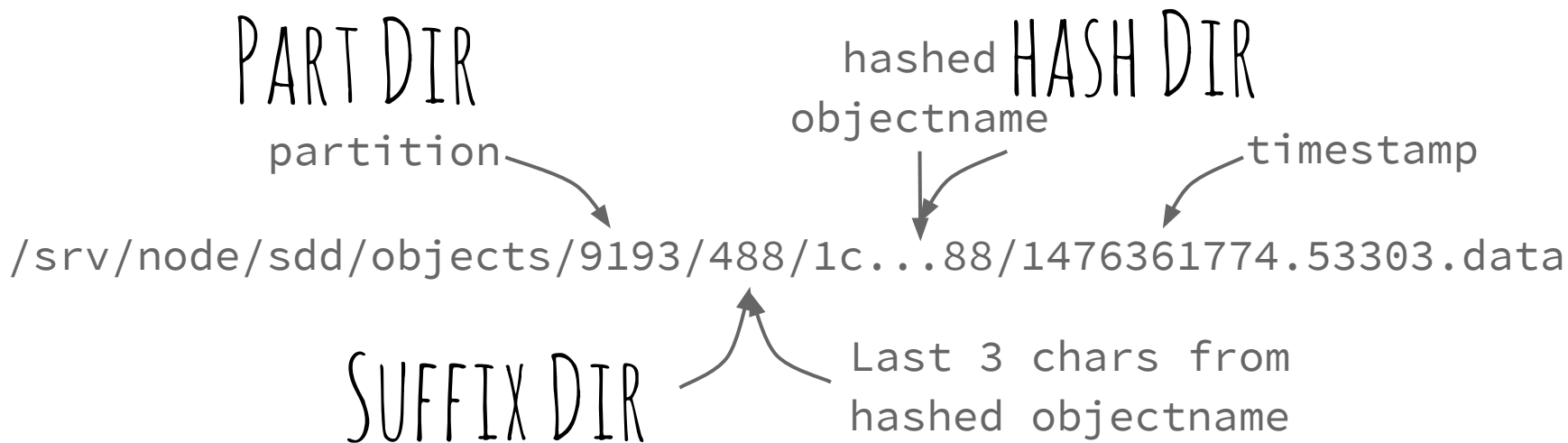
Consistent what?

- Just remember the distribution function
- No growing lookup tables!
- Easy to distribute!



Partitions in Swift

- Object namespace is mapped to a number of **partitions**
- Each partitions holds one or more objects



replica2part2dev_id

SWIFT'S ADDRESS BOOK



	Replica # 1	Replica # 2	Replica # 3
Part # 0	Device # 0	Device # 1	Device # 3
Part # 1	Device # 3	Device # 0	Device # 1
Part # 2	Device # 3	Device # 4	Device # 2
Part # 3	Device # 2	Device # 0	Device # 1
Part # 4	Device # 1	Device # 4	Device # 3
Part # 5	Device # 0	Device # 2	Device # 4
Part #

How to lookup partition

PRIMARY

`get_nodes(part)`

Part # 2	Device # 3	Device # 4	Device # 2
----------	------------	------------	------------

HANDOFF

`get_more_nodes(part)`



What makes a good ring

A good ring has good

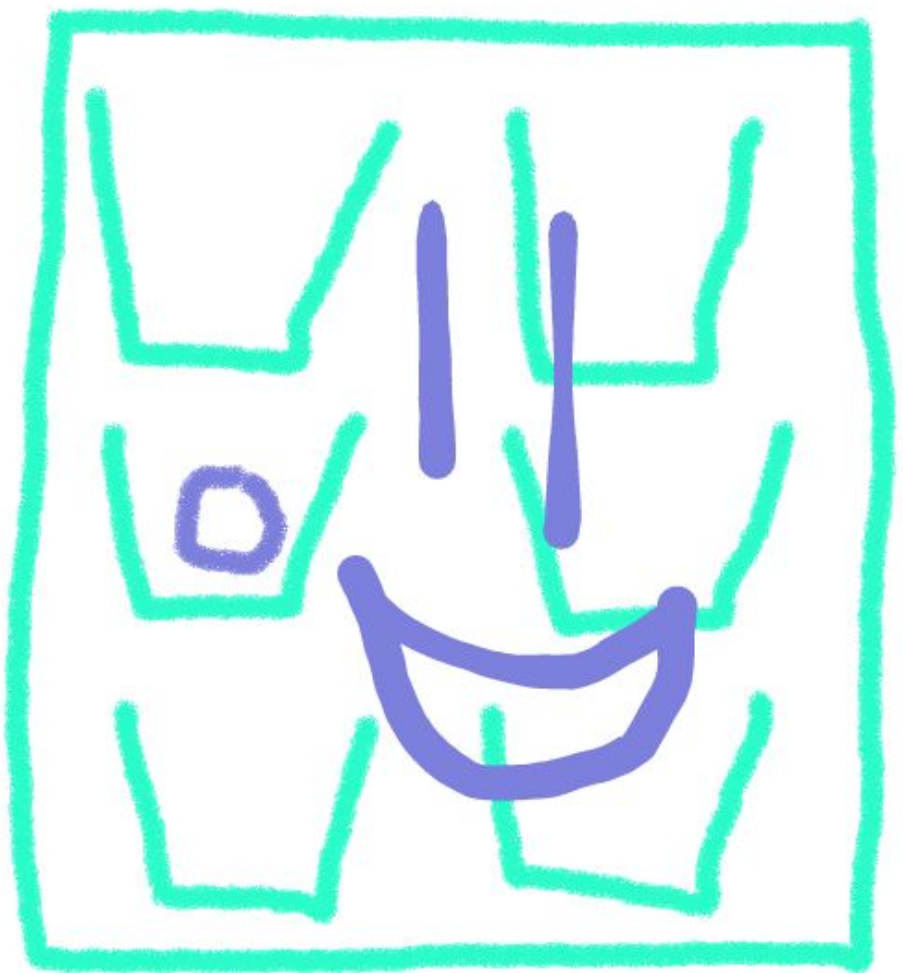
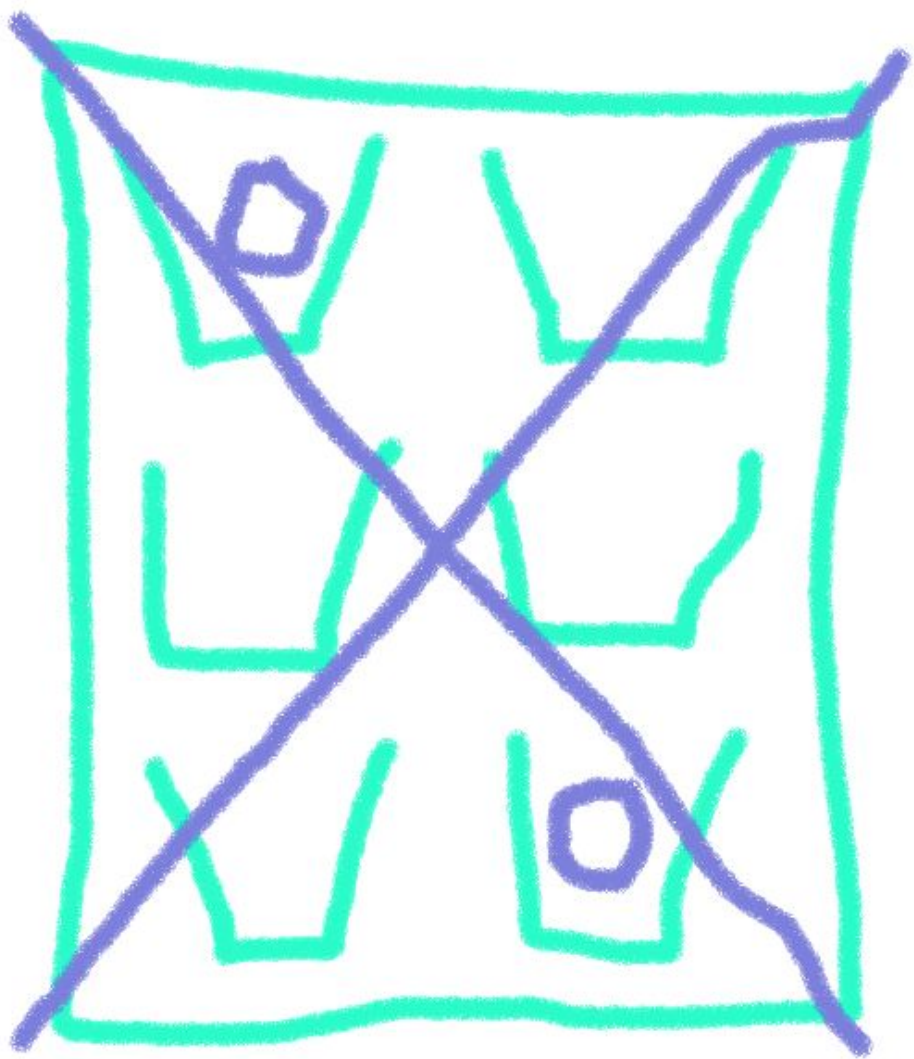
- Dispersion
- Balance
- Low overload (SOME, BUT NOT TOO MUCH!)

PC LOAD LETTER

Fundamental Constraints

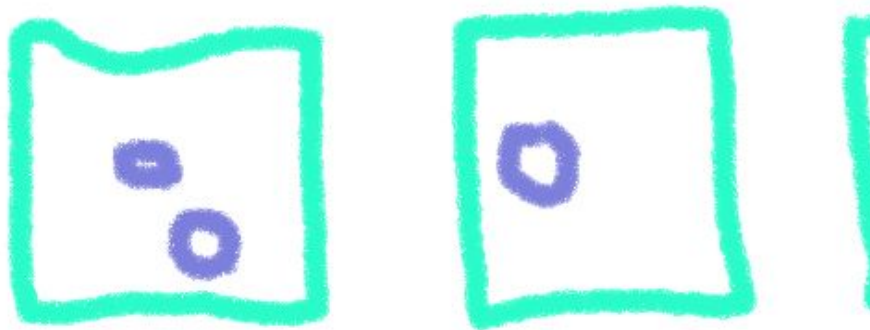
- **Devices (disks)** A FAILURE DOMAIN
- **Servers** FAILS TOGETHER
- **Zones (racks)**
- **Regions (datacenters)**

↑
THESE ARE TIERS



Dispersion

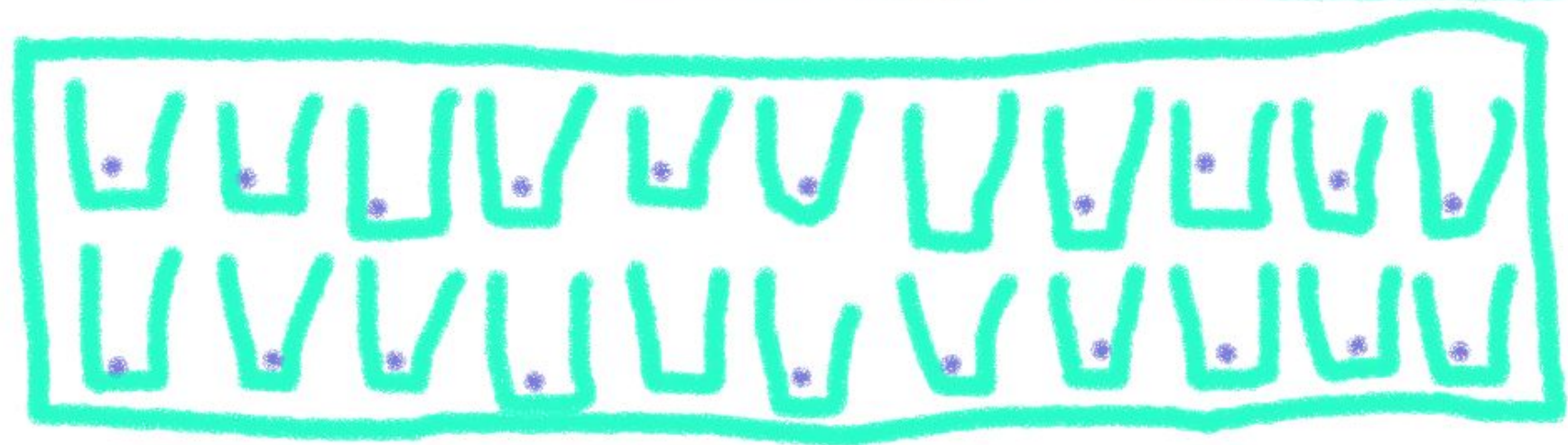
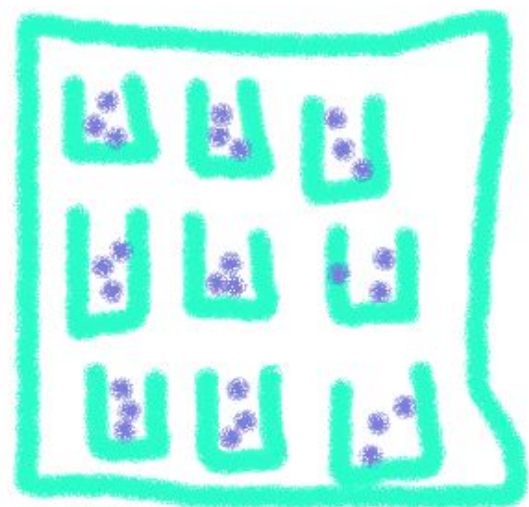
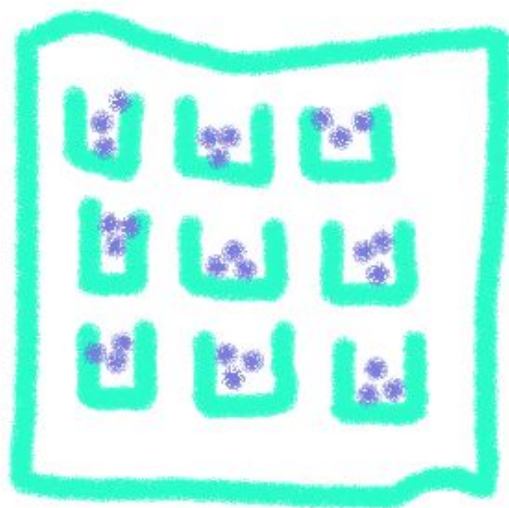
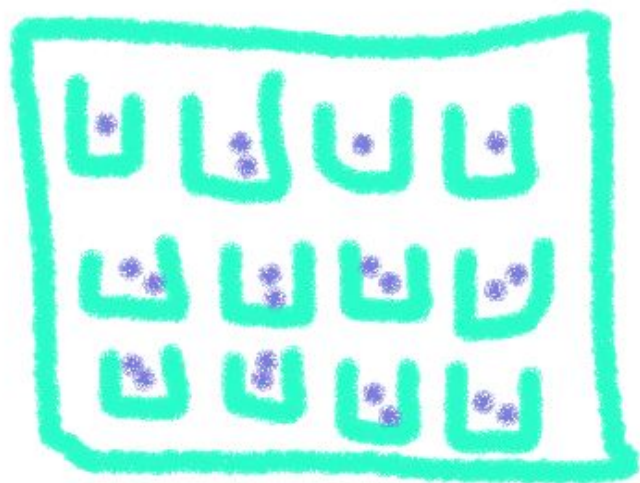
MEASUREMENT THAT THE FAILURE DOMAIN OF EACH REPLICA OF A PART
IS UNIQUE AS POSSIBLE



Fundamental Constraints



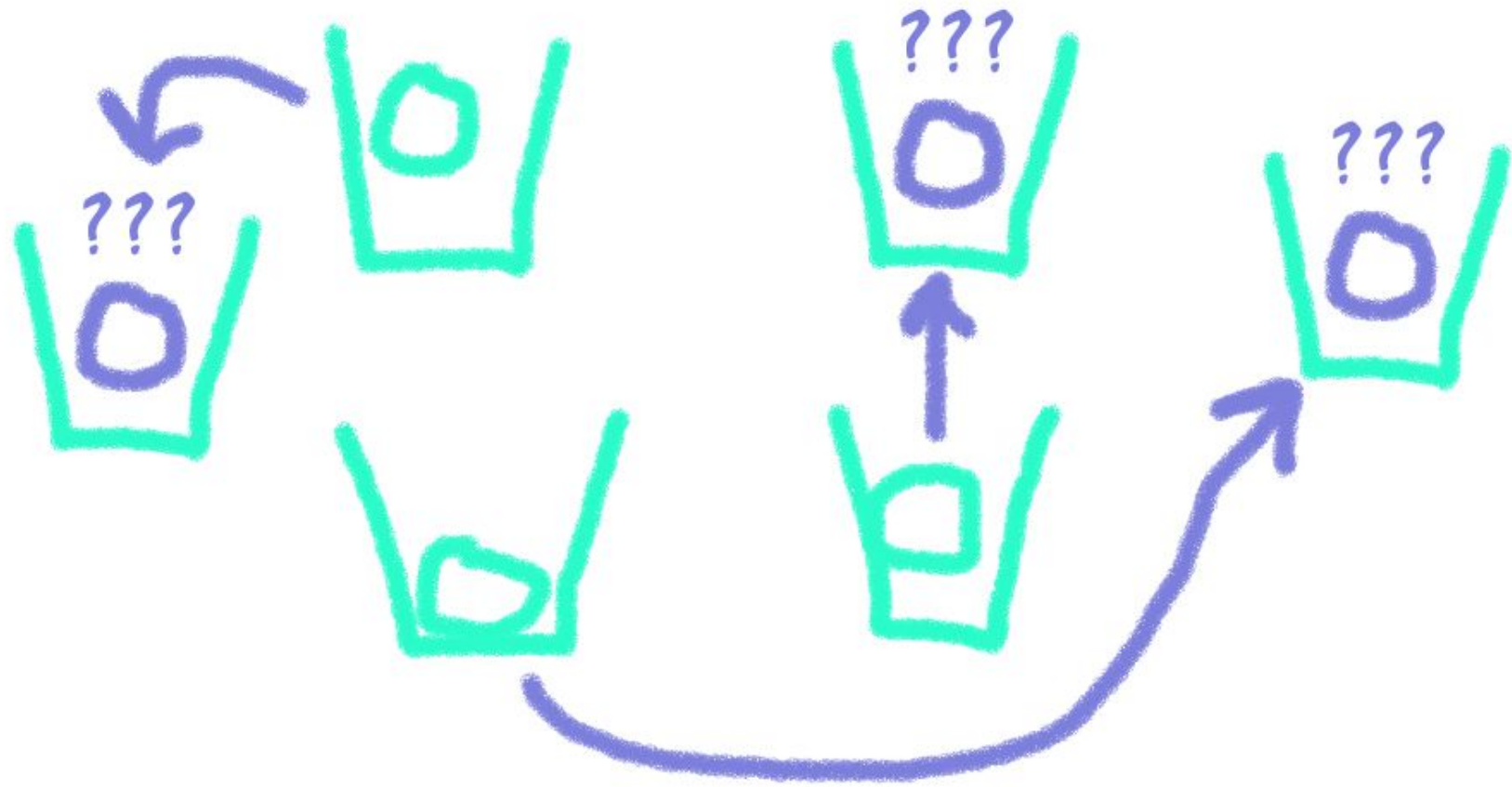
BALANCE



The Rebalance Process

"RINGS ARE NOT PIXIE DUST THAT MAGIC DATA OFF OF HARD DRIVES"

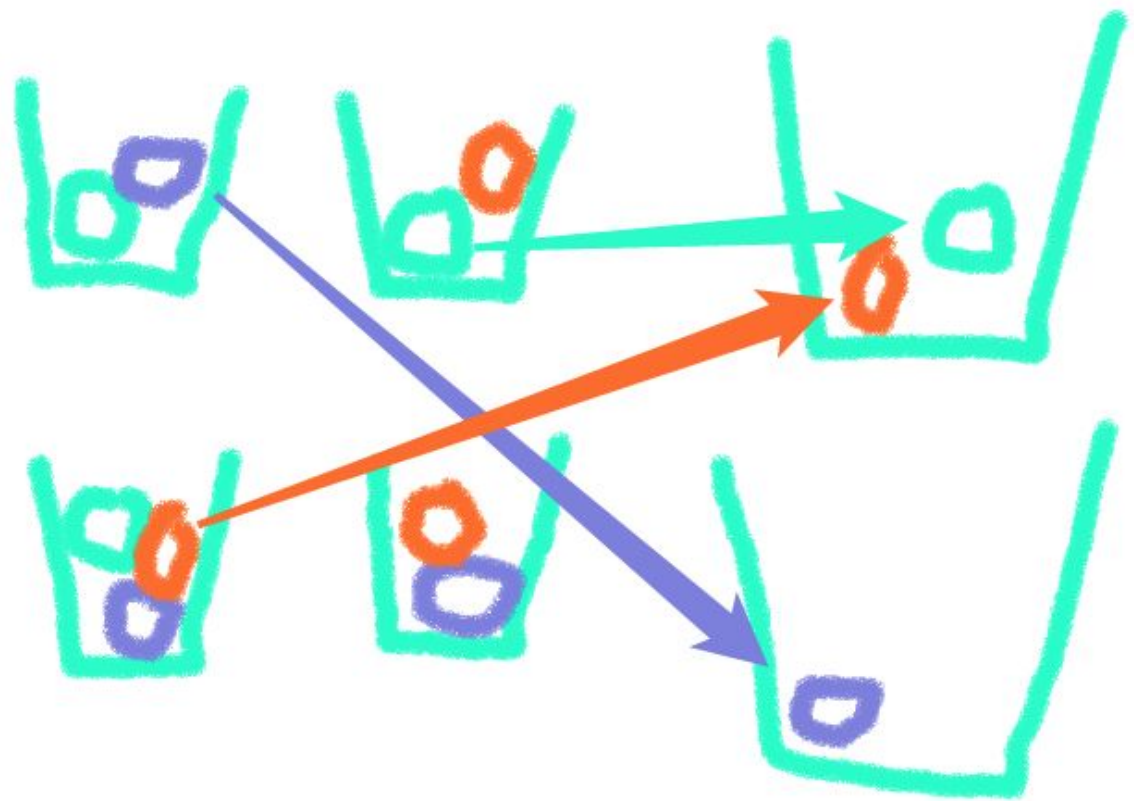
-- DARRELL



Fundamental Constraints

min_part_hours

ONLY MOVE ONE REPLICA OF A PARTITION
PER REBALANCE



Monitoring Replication Cycle

- ONLY REBALANCE AFTER A FULL REPLICATION CYCLE
- SWIFT-DISPERSION-REPORT IS YOUR FRIEND

```
Queried 8192 objects for dispersion reporting, ...  
There were 3190 partitions missing 0 copy.  
There were 5002 partitions missing 1 copy.  
79.65% of object copies found (19574 of 24576)
```

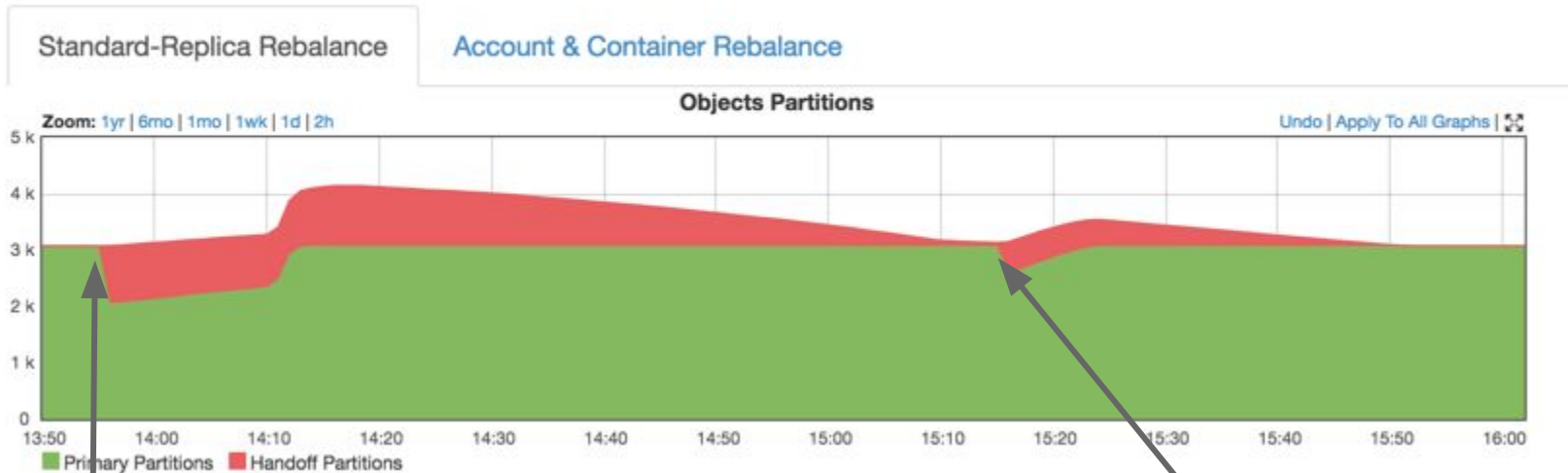
Hostname	Role	Zone	Storage
ss-rax-geo-syd-node1	Swift Node Account & Container, Standard-Replica (Default)	SYD Zone 1	
ss-rax-geo-syd-node2	Swift Node Account & Container, Standard-Replica (Default)	SYD Zone 1	
ss-rax-geo-syd-node3	Swift Node Account & Container, Standard-Replica (Default)	SYD Zone 1	
ss-rax-geo-dfw-node2	Swift Node Account & Container, Standard-Replica (Default)	DFW Zone 1	
ss-rax-geo-dfw-node1	Swift Node Account & Container, Standard-Replica (Default)	DFW Zone 1	
ss-rax-geo-dfw-node3	Swift Node Account & Container, Standard-Replica (Default)	DFW Zone 1	

PARTITIONS
ASSIGNED

GB USED

STARTING TO FILL!

Rebalance Status



RING PUSH



Primary Partitions



Handoff Partitions

FIRST CYCLE

FINISHED

OVERLOAD

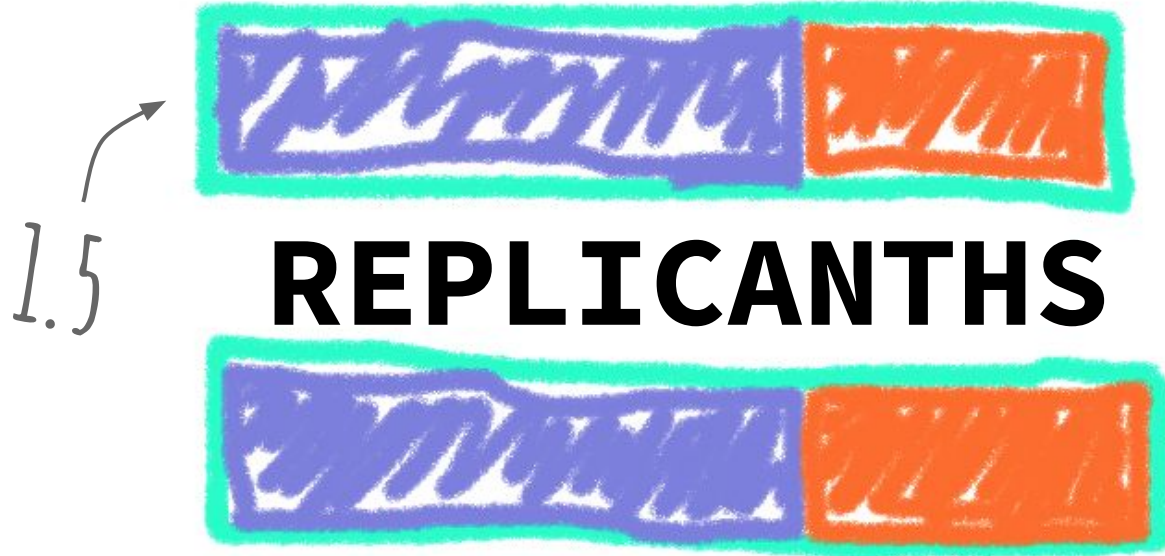
Balance vs. Dispersion

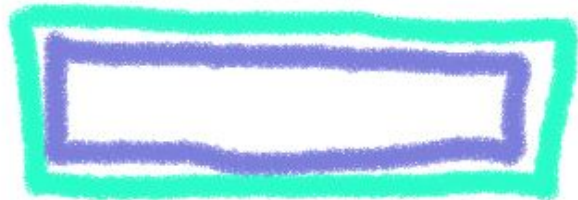
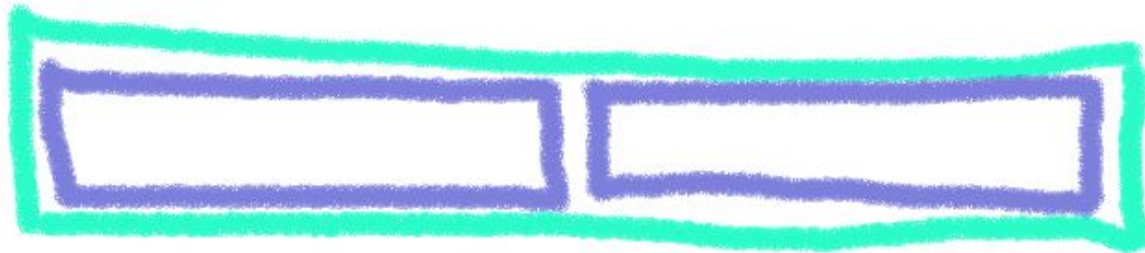


9.999999

FIGHT!

The decimal fraction of one
replicas worth of partitions

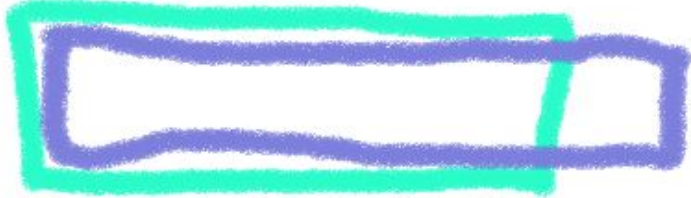
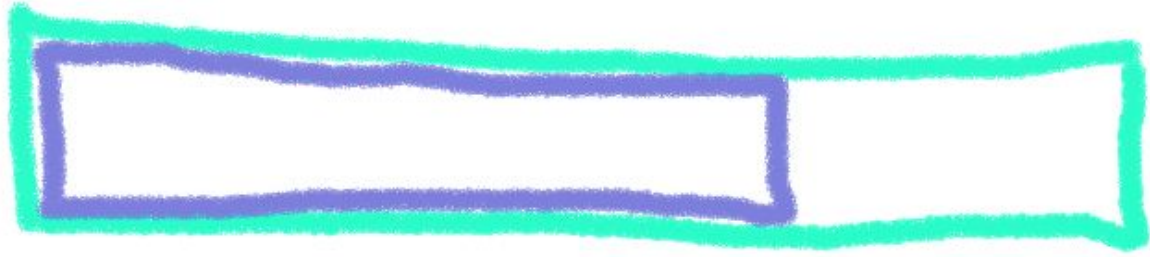




$$\frac{3}{5} = 0.6$$

REPLICAS
"UNITS"





2 REPLICAS

.6 => .66 ~ 11%

Overload

Too Much => DRIVES FILL UP



Not Enough => CORRELATED DISASTER

JUST USE 10% (HOPEFULLY IT WAS CAT PICS?)

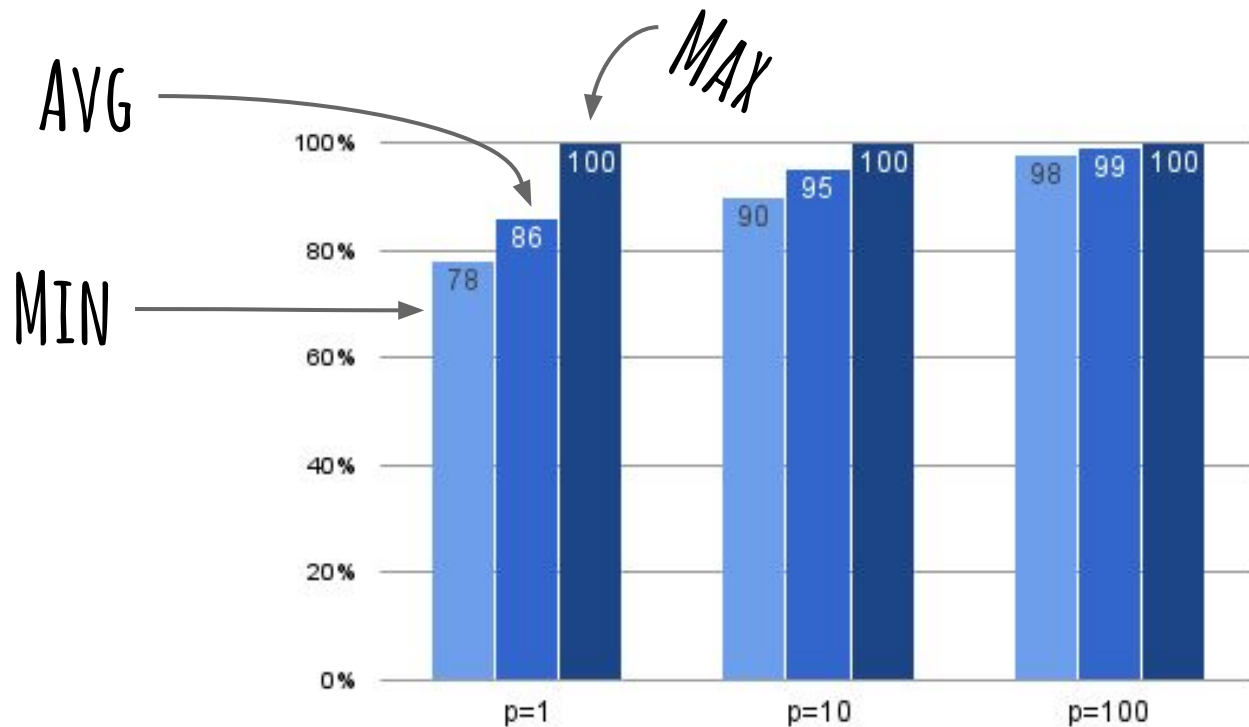
... IT'LL PROBABLY BE FINE

PARTITION POWER

BALANCING THE UNKNOWNNS

- How to distribute objects of unknown size well-balanced?
 - Objects vary between 0 bytes and 5 GiB in size
- => Store more than one partition per disk
- => Aggregation of random sizes balances out

DISK FILL LEVEL VS. PARTITION COUNT



CHOOSING PARTITION POWER

- Number of partition is fixed
- More disks => less partitions per disk
- Choose a part power with a ~ thousand partitions per disk
 - Based on today's need, not an imaginary future growth
- It is highly unlikely that your partition power is >> 20, and definitely not 32

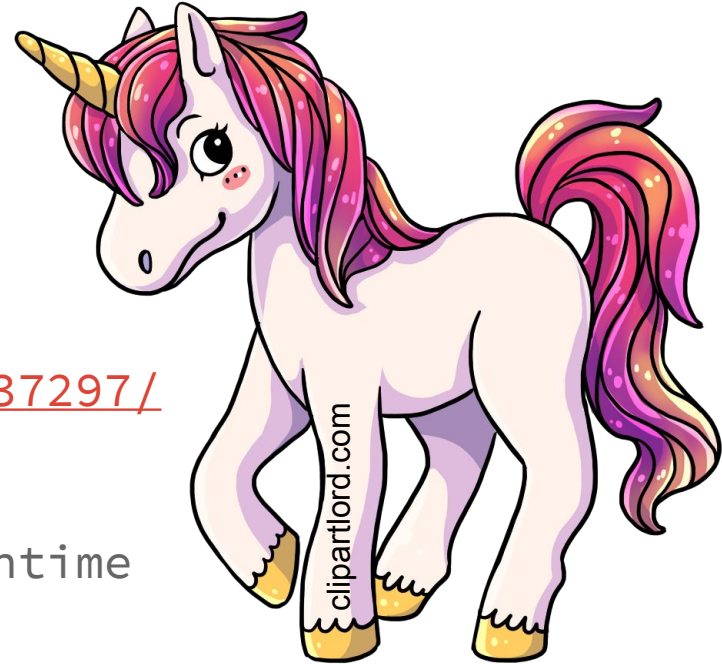
<https://gist.github.com/clayg/6879840>

YOU BECAME AN UNICORN

- Skyrocketing growth? Congrats!
- We're working on increasing partition power for you to keep your cluster balanced

<https://review.openstack.org/#/c/337297/>

- Decreasing won't be possible -
at least not without a serious downtime



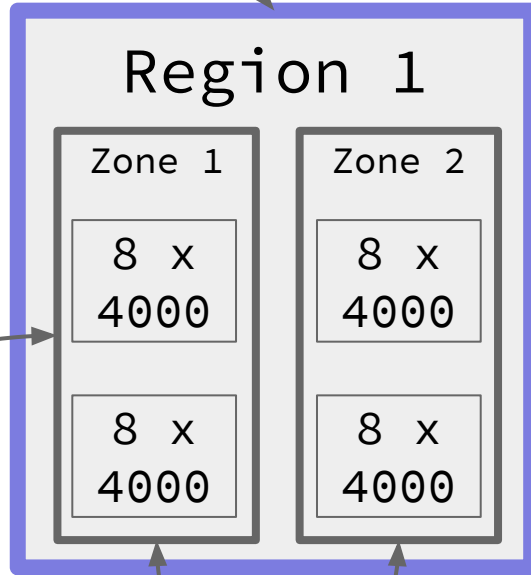
WRAPPING UP

What's a good cluster?

$$\text{PARTPOWER } 14 \rightarrow 2^{14} = 16384$$

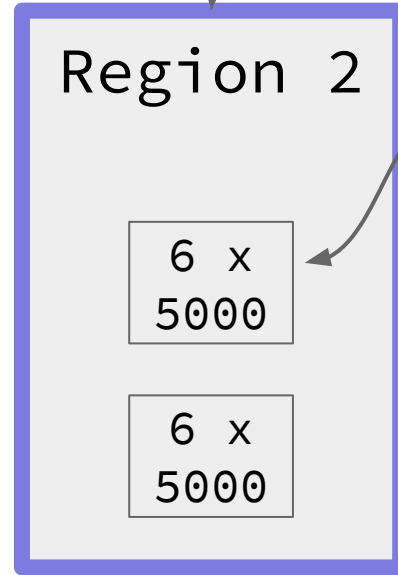
$$16384 \text{ PARTITIONS} * 3 \text{ REPLICAS} / 32 \text{ DISKS} = 1536 \text{ PARTS PER DISK}$$

MAIN DATACENTER



64 TB

2ND DATACENTER



60 TB

DISK WEIGHT

$$\frac{(64+64+60)}{3} = 62.66$$

OVERLOAD: 4.5%

DISPERSION: 0

BALANCE: 4.65

Questions?

THANKS!

CLAY@SWIFTSTACK.COM

CSCHWEDE@REDHAT.COM