# Millions of Tiny Databases

NSDI'20
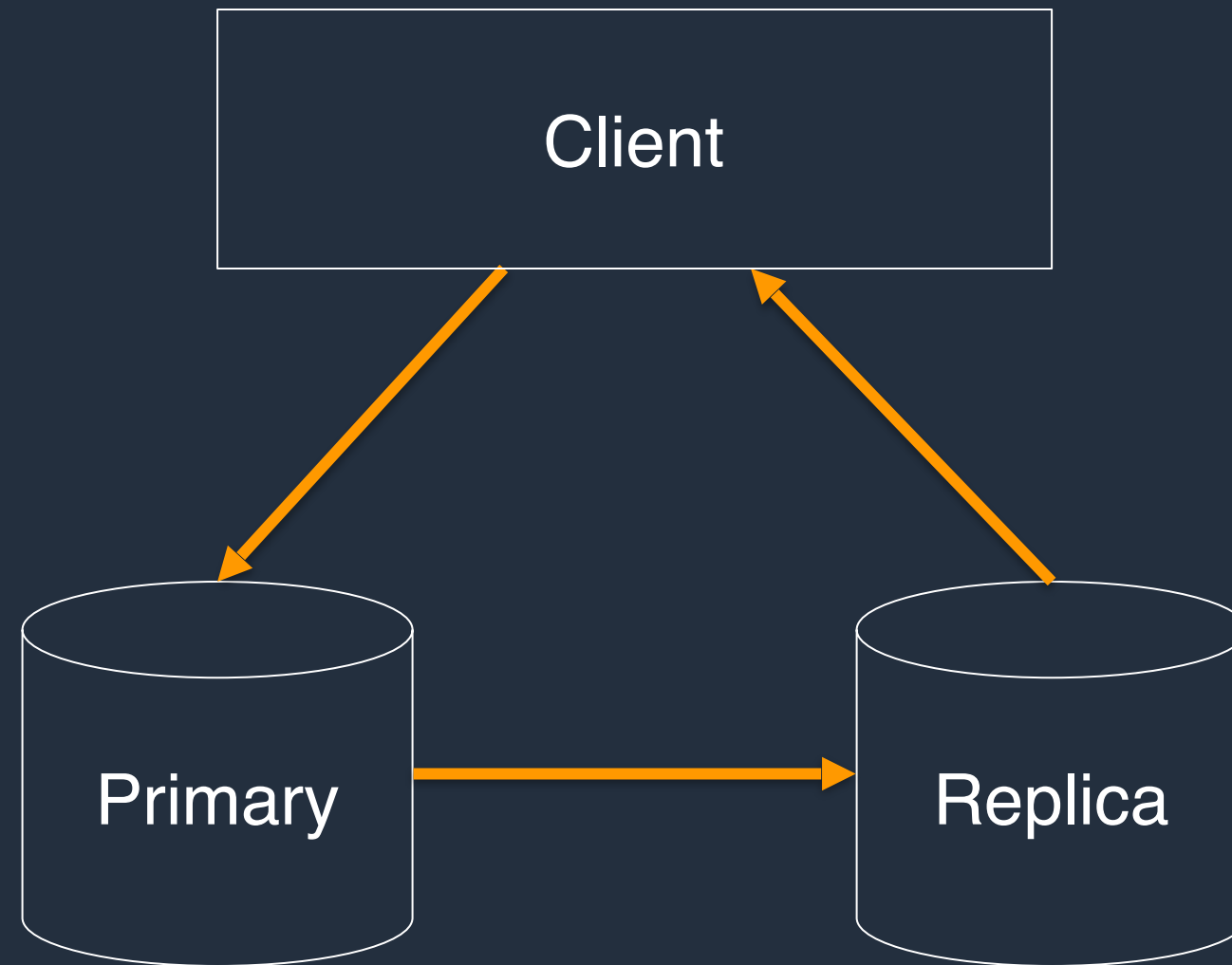
**Marc Brooker**, Tao Chen and Fan Ping

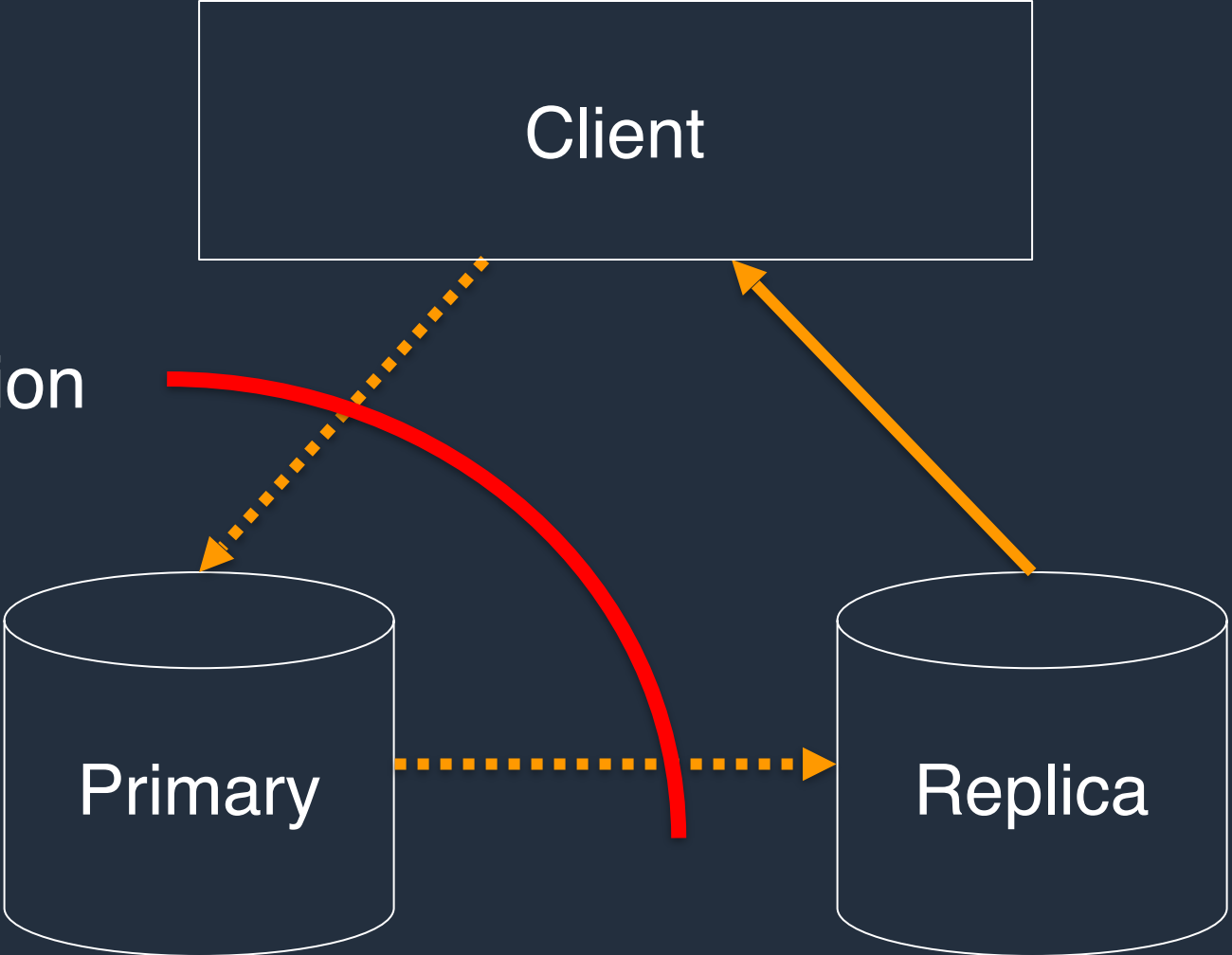February 2020

# Table of contents

- Tough CAP Tradeoffs
- Availability and Blast Radius
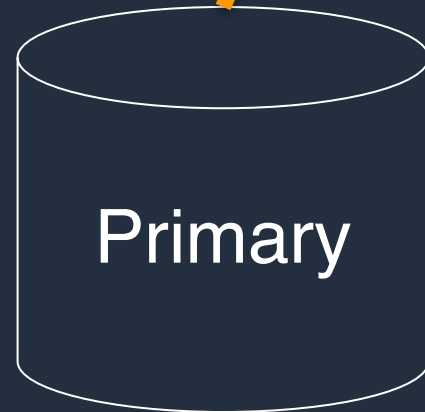- Physalia Architecture

# Simplified Storage System

100k+
SERVERS.

Here?

Or
Here?

aws

PARTITIONS ARE NOT CLEAN!

# Availability and Blast Radius

# Availability
## is typically improved with
# Redundancy

aws

# Availability
## is typically improved with
# Redundancy*

* Unless failures are correlated

aws

# Infrequent

# Short

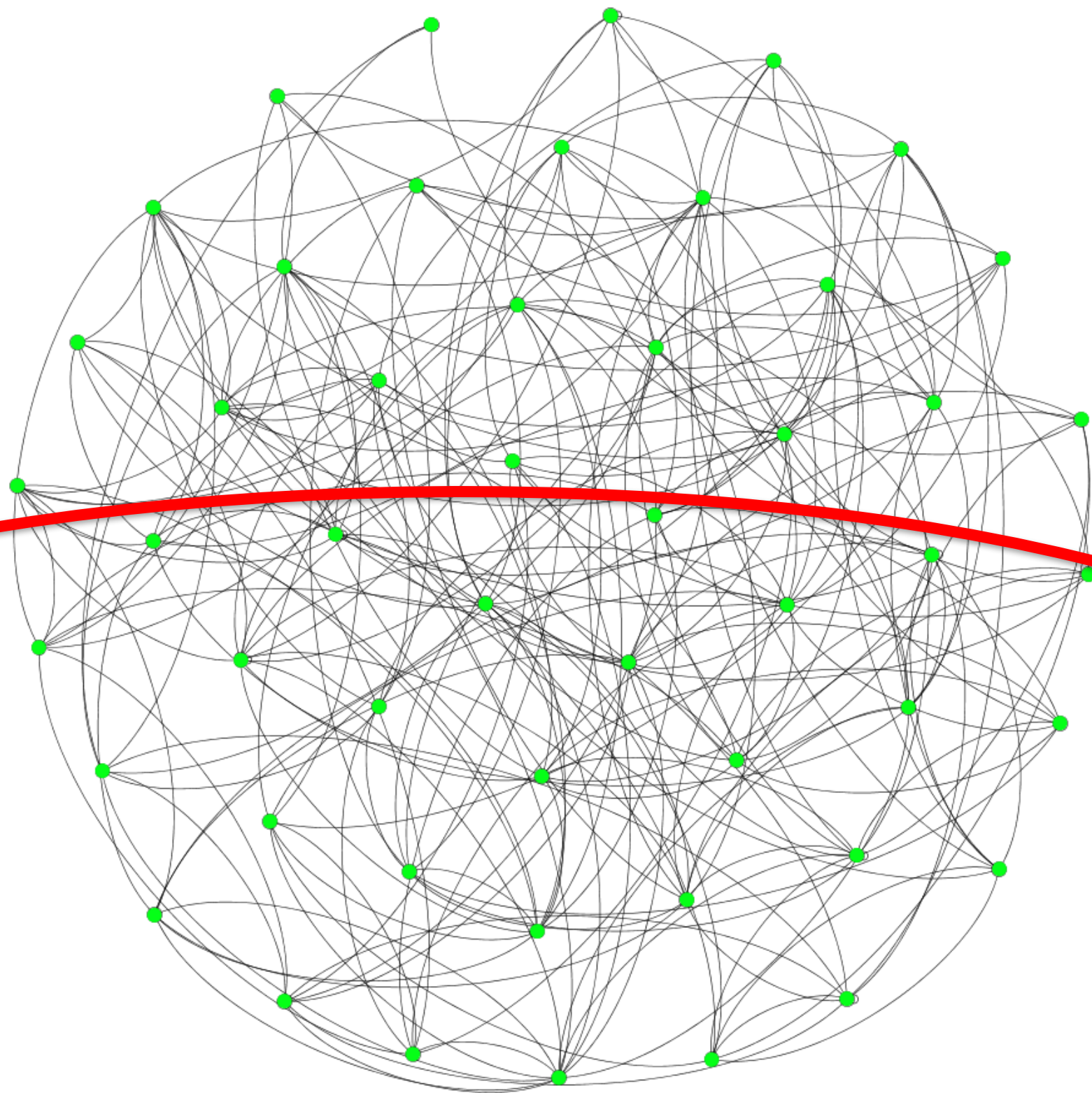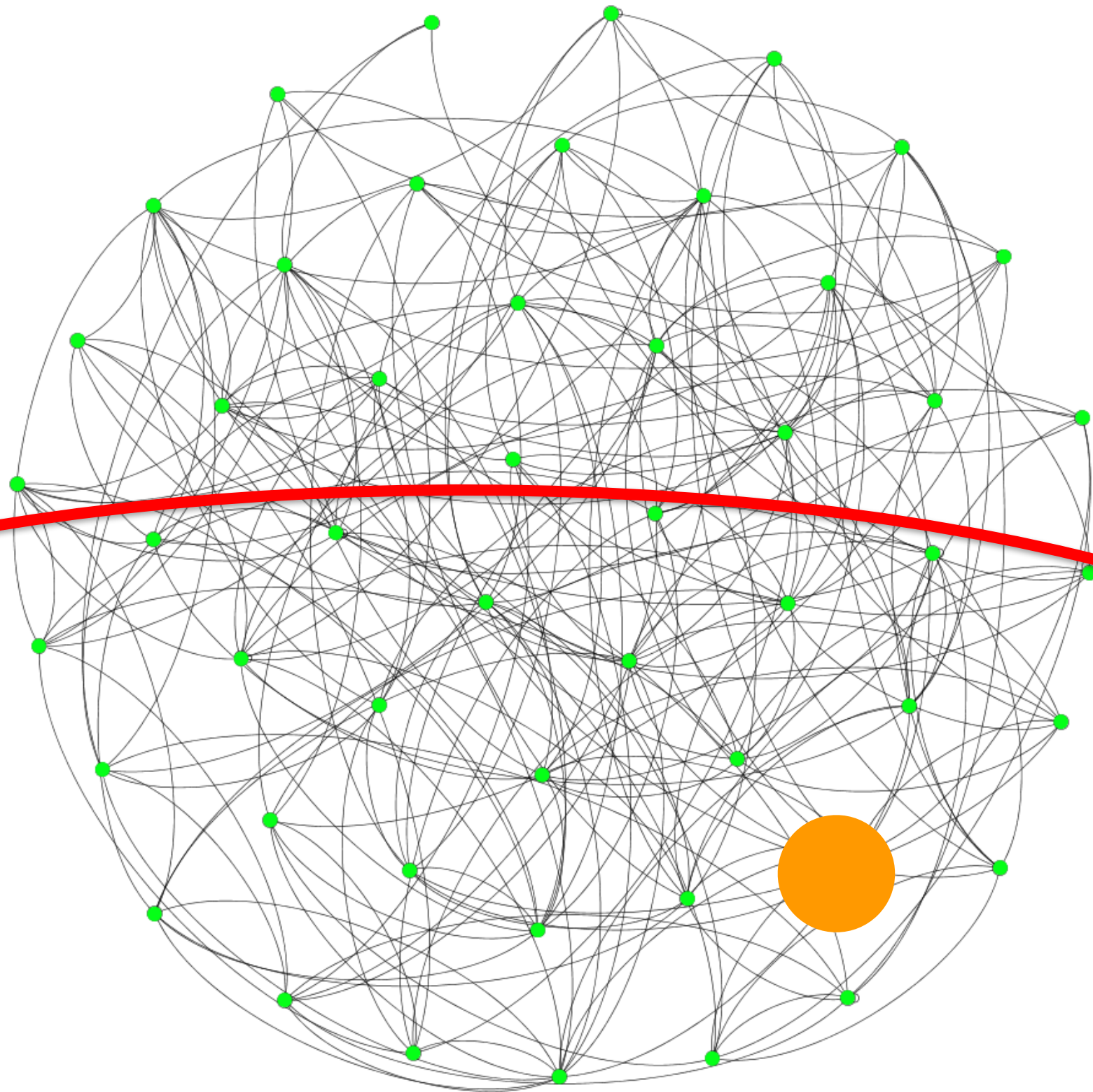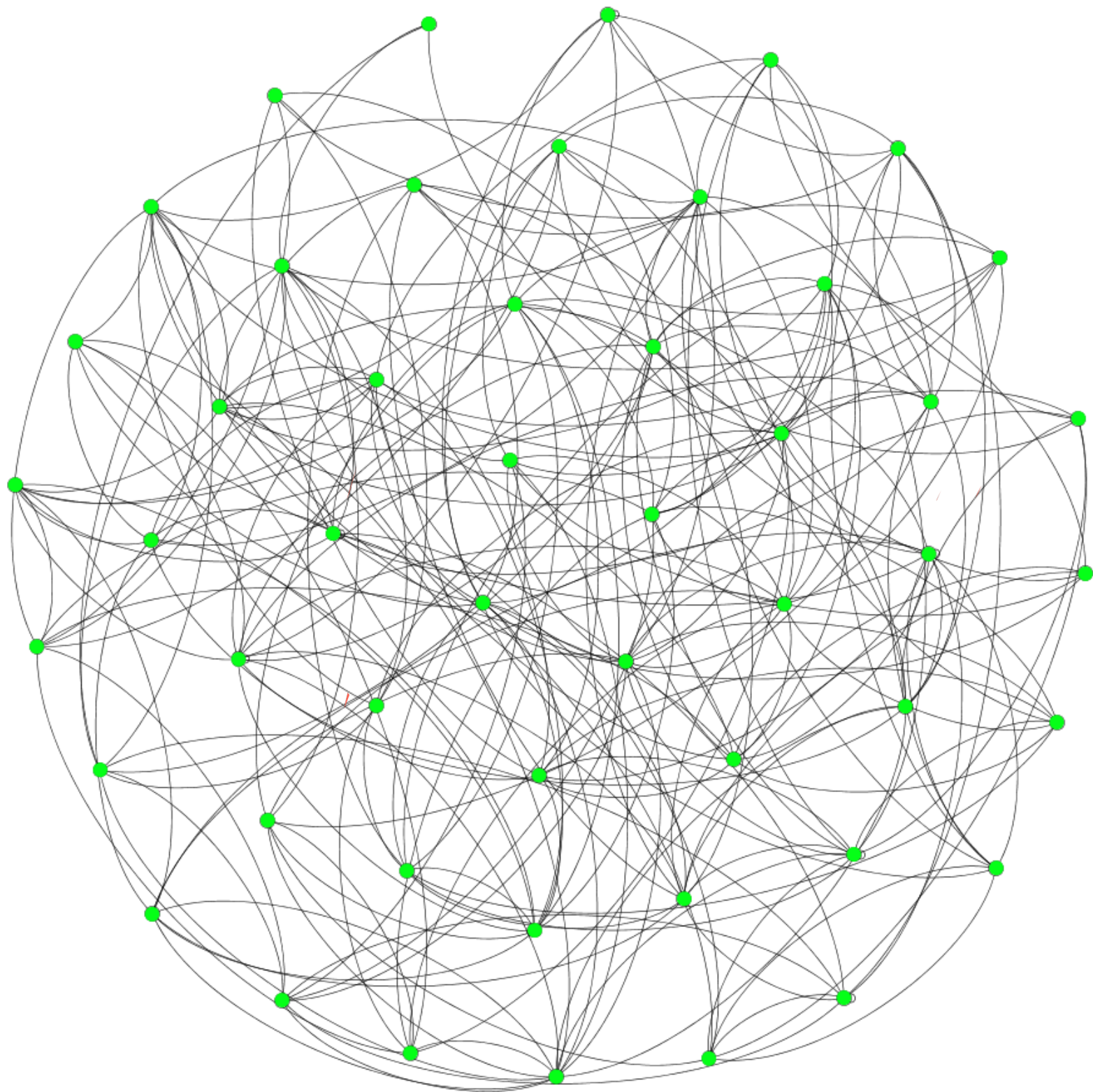# Small

aws

# Infrequent
# Short
# Small = UNCORRELATED

aws

# "Blast Radius"

# Physalia Architecture

Client

Primary

Replica
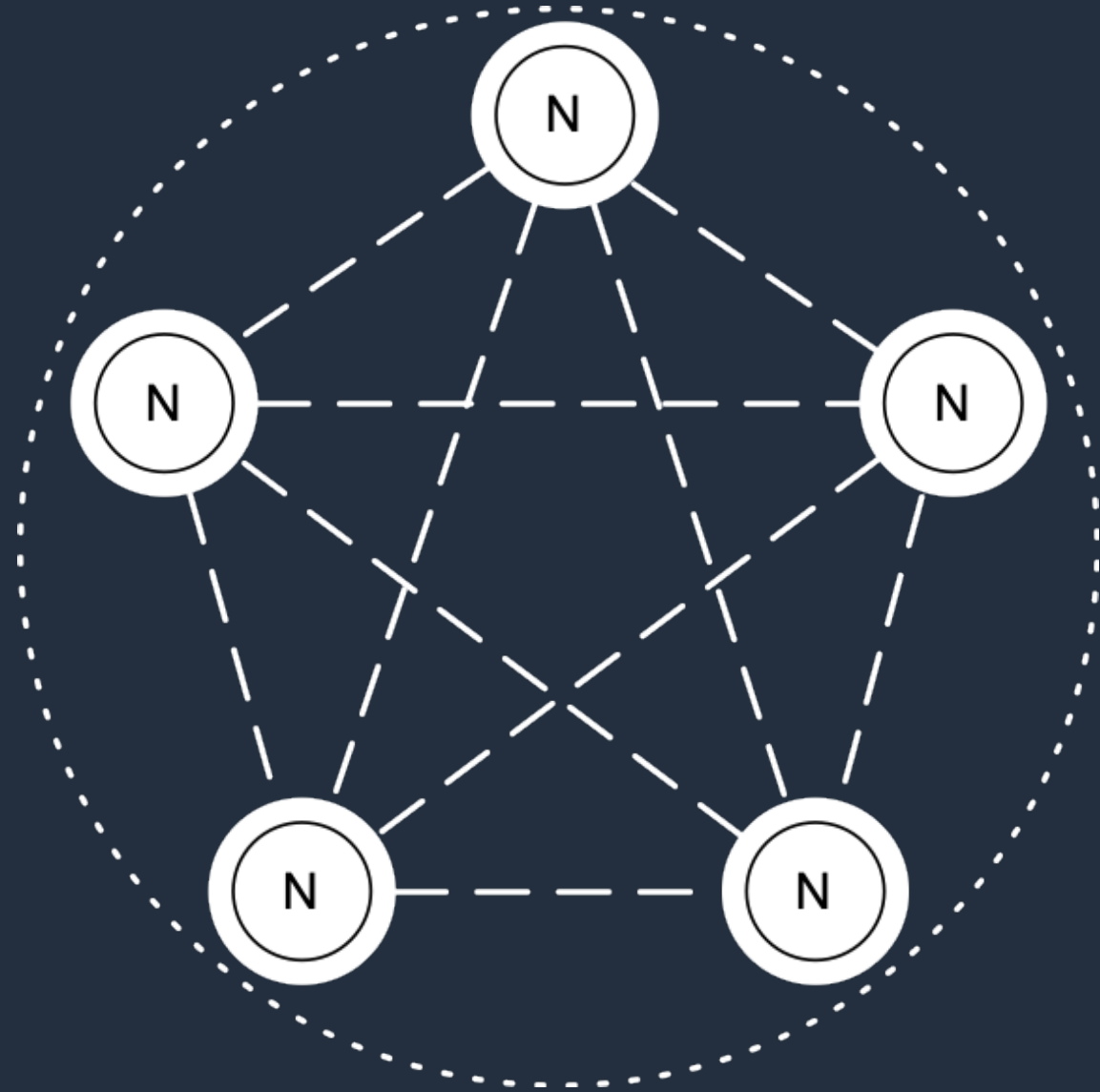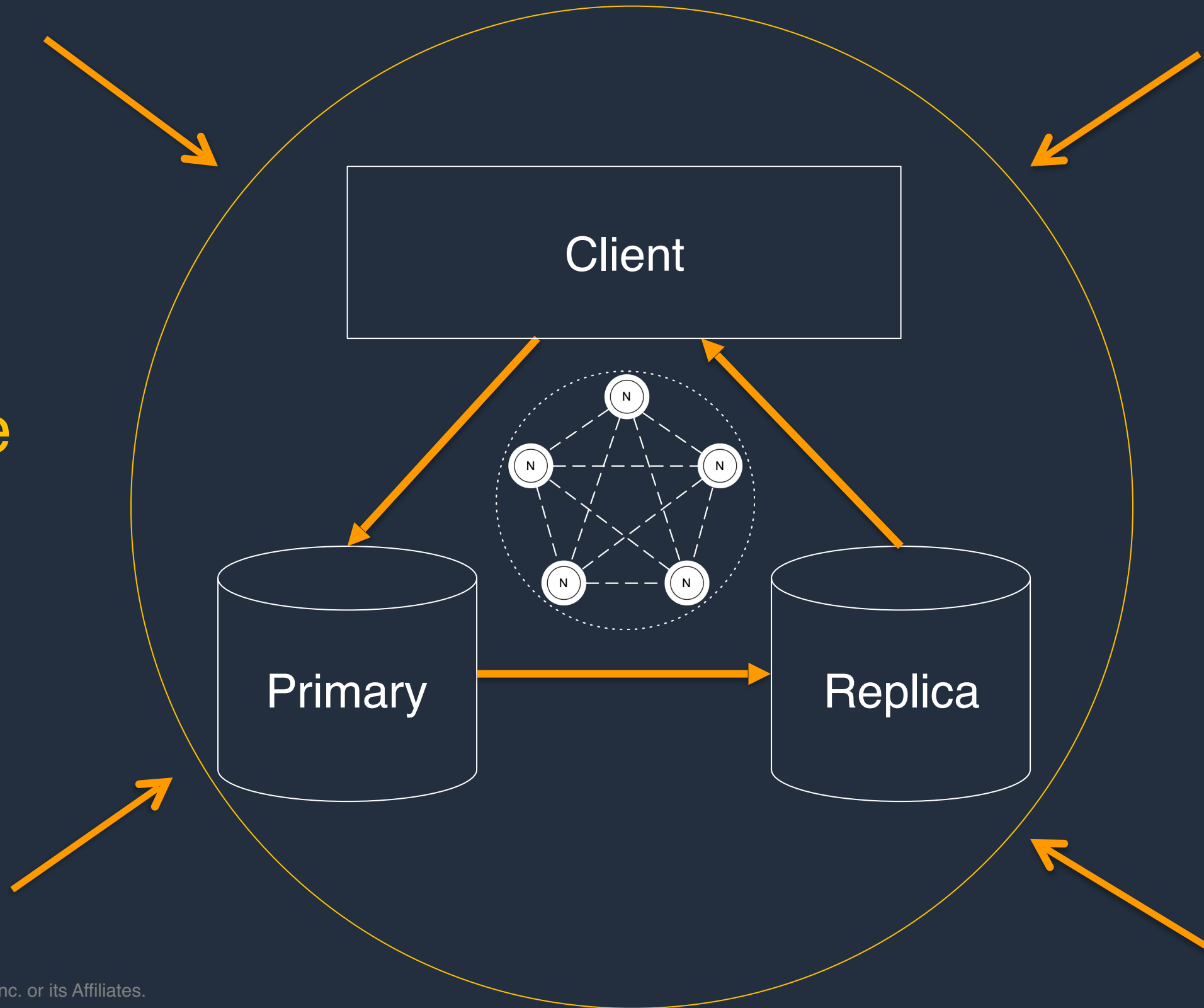
aws

# Physalia *Cell*

- Replicated state machine

- Configuration for one volume, or small set of volumes.

- K/V store API

- Strict serializable transactions

aws

Minimize The Radius

Client

Primary

Replica

# Topology Details Matter



Lower Partition Risk

Lower Availability Blast Radius

Radius

More Redundancy

More Bisection Bandwidth
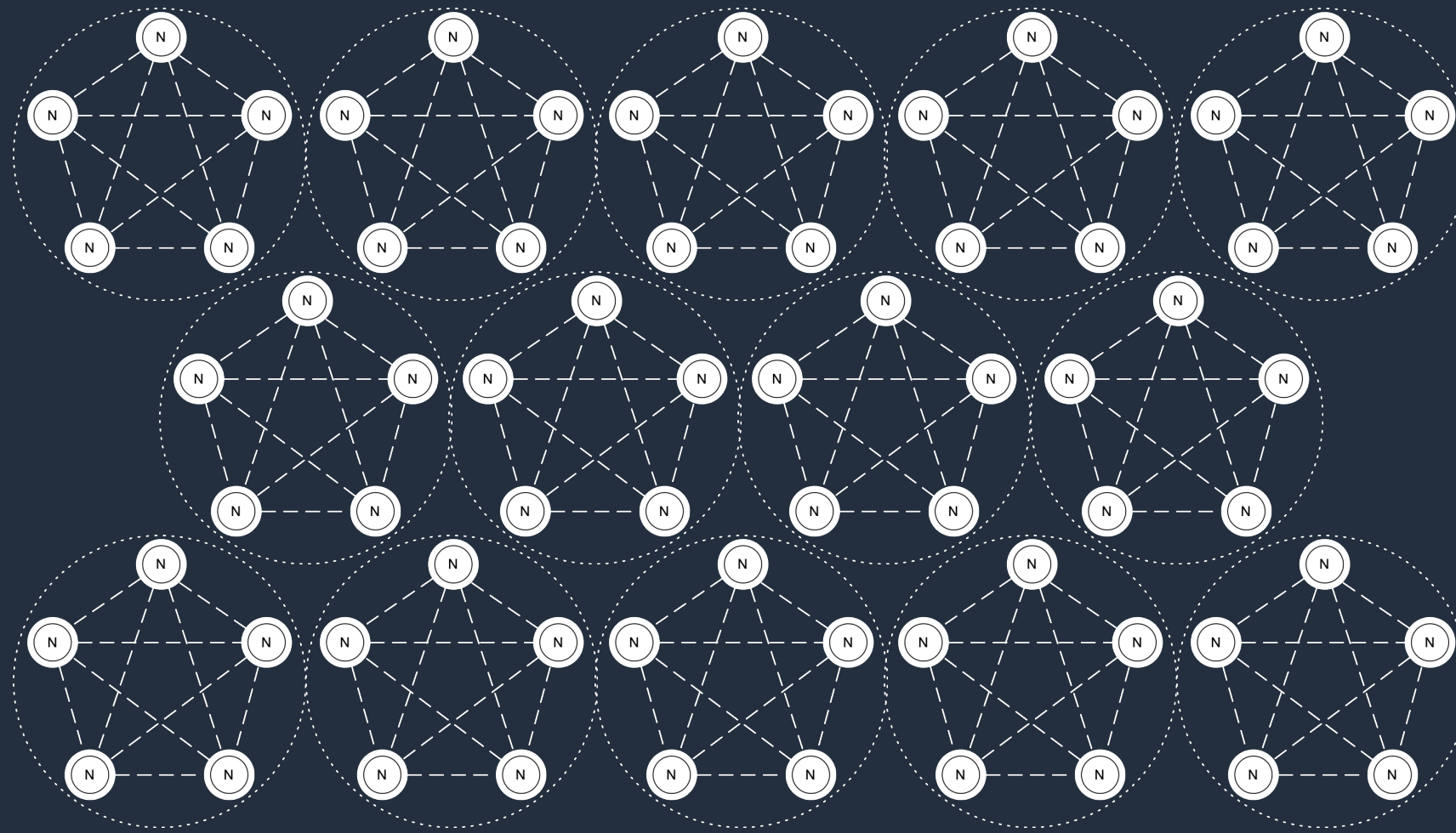
More Placement Options

aws

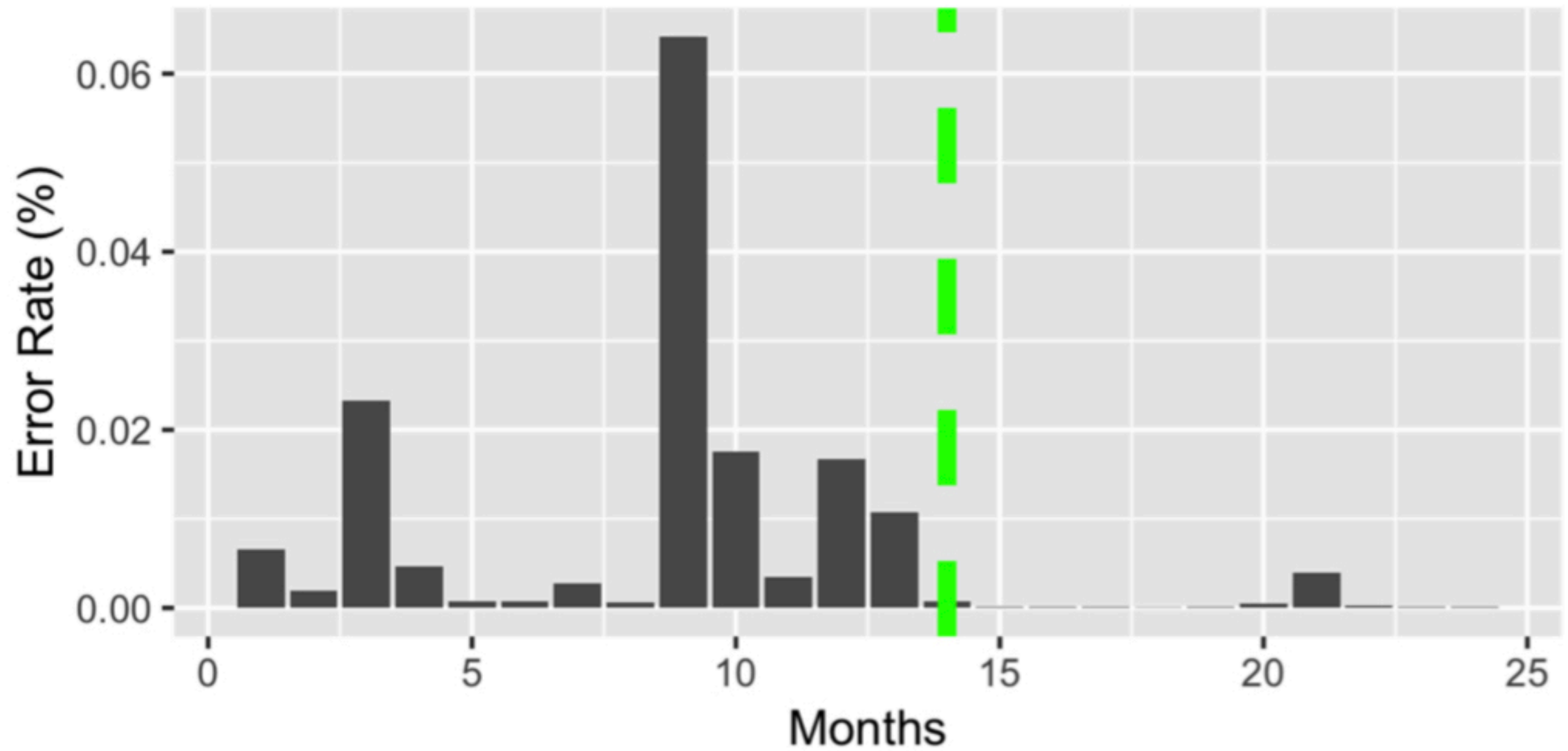# Take Advantage of Eventual Consistency When You Can!

- Discovery Cache (clients discover nodes)
- Monitoring
- "Meta" control plane

aws

# Optimize for Blast Radius

- Minimize impact of partitions (and CAP tradeoffs),
- overload,
- software bugs &
- operational issues.

Build *humility* into the system.

aws

# Q&A

Marc Brooker
mbrooker@amazon.com
@marcjbrooker

# The End

aws