# Improving Kafka Resilience - Gray Failures Mitigation
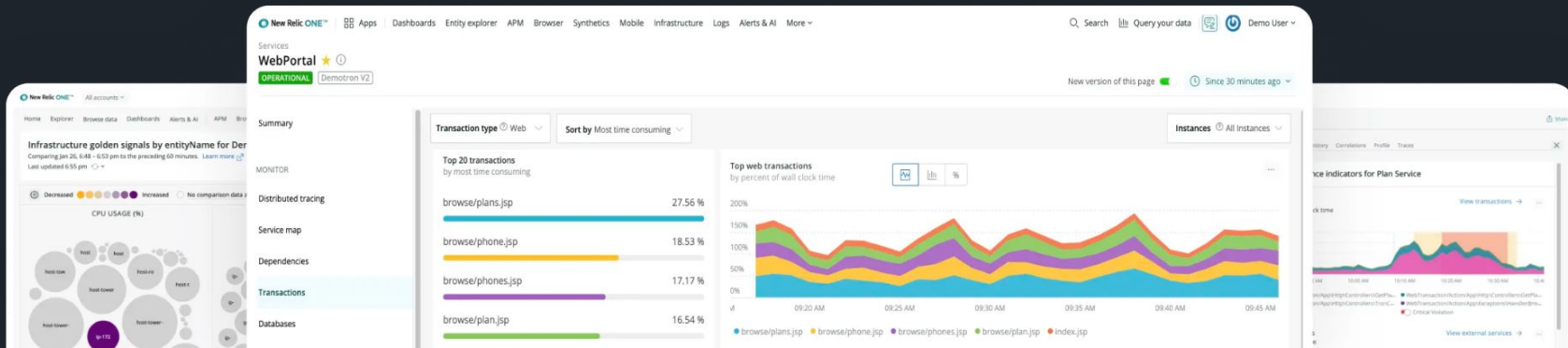
Michelle Valentinova, Senior Site Reliability Engineer

Kafka Platform Team
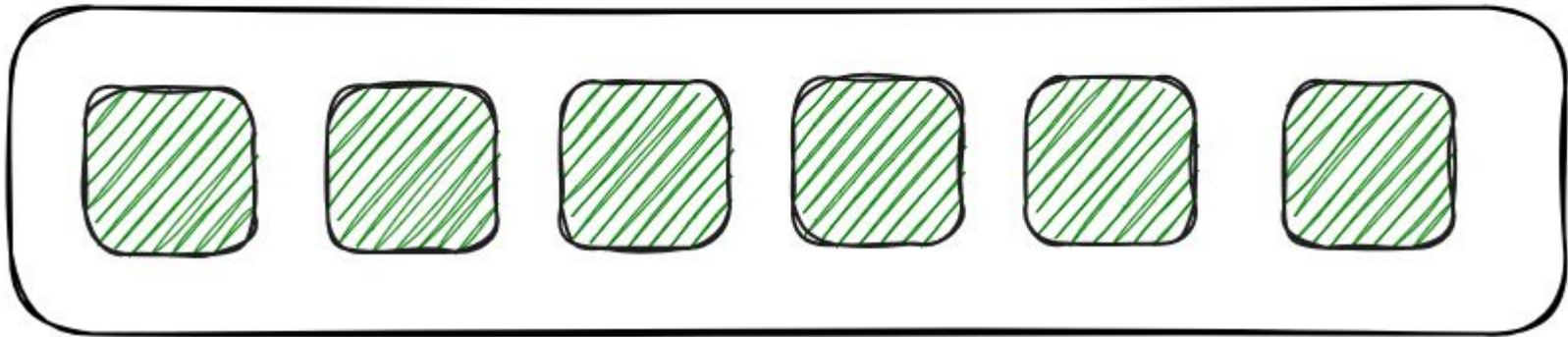
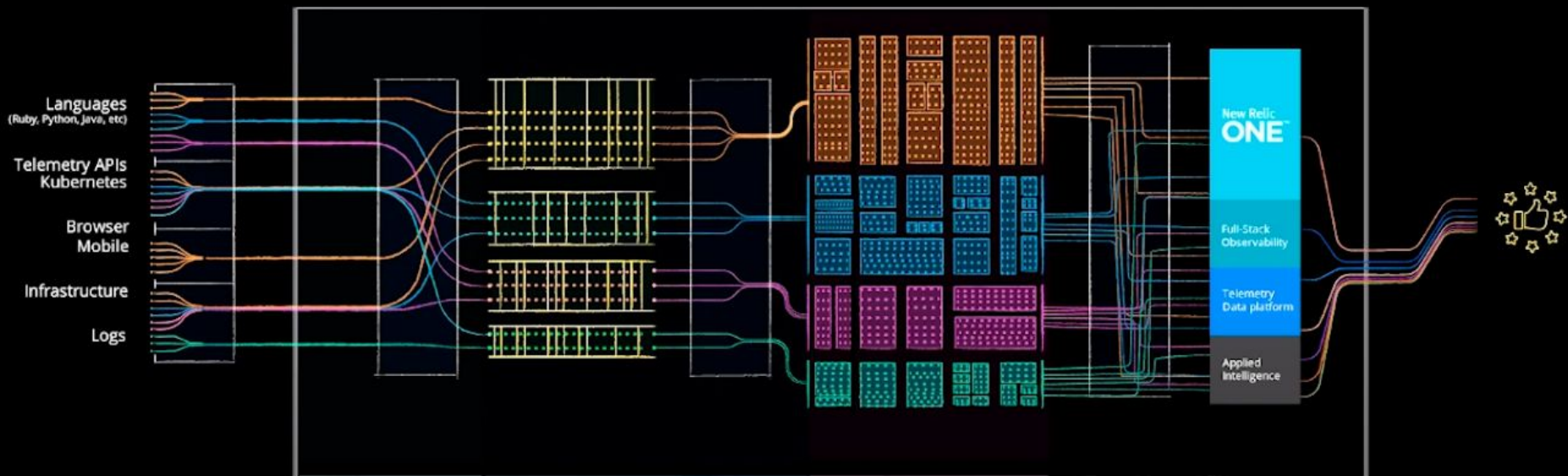# Apache Kafka - event streaming

Topic

new relic.

# Kafka at New Relic

# 97 clusters
# 2000 brokers
# 180M messages/s

new relic.

# Distributed systems

**It works**

**It doesn't work**

new relic.

# Distributed systems

## It kind of works

new relic.

# Kafka concepts

new relic.

# Topics are partitioned

- Parallel processing
- High availability
  - Failover between partitions

new relic.

# Topics, partitions, and brokers in Kafka

# Kafka in the cloud architecture

# Kafka broker's network storage very slow

# Kafka broker's network storage very slow
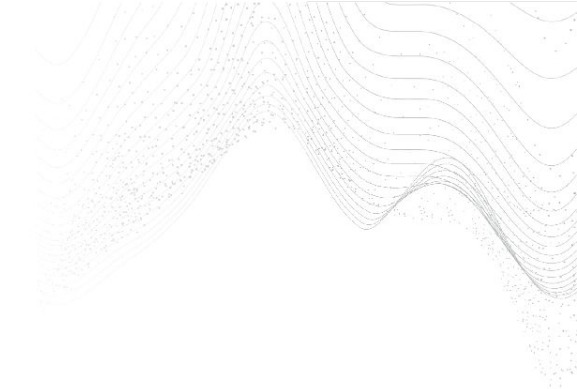
# Issues with a single broker affect all producers



**Bytes in per second (by broker)**
Since 1 hour ago

new relic

# Producers optimise for throughput (send buffer)

# One broker is slow

Producers retry infinitely (configurable)

Don't want to lose customer logs, metrics, miss alerts, etc.

new relic.

# Send buffer fills up

# The issue is resolved

# The issue is resolved continued

new relic.

# Issue with a single broker affects all producers

# Impact for all producers:
# 25 minutes

# Customer metrics, logs, alerts, etc. delayed for 25 minutes

new relic

# Why do we even use network storage if it causes problems?

new relic.

Our busiest Kafka clusters have 1.7PB of storage

We use a managed service and can't choose local disk

new relic

# Batching records using a partitioner

new relic

# What happens when there is a failure

# Adaptive random partitioning

# Strictly uniform sticky partitioner

- Available in Apache Kafka 3.3.1 - 1 year ago
  - Our busiest clusters have 280 topics
  - Rolling out changes can take some time
- Partition probability to get records is inverse of queue size
- The partitioner is configured on the Kafka producers

new relic.

# Strictly uniform sticky partitioner

- Works well out of the box
- But can be configured to improve even further

new relic.

# Send buffer is shared

# Partition availability timeout - disabled vs 5s



**Producer buffer availability**
Sep 6, 3:10pm – Sep 6, 3:40pm

450 kB/s

70 MB
60 MB
50 MB
40 MB
30 MB
20 MB
10 MB
0 B

3:10pm    3:15pm    3:20pm    3:25pm    3:30pm    3:35pm

● vortex-arm-0    ● vortex-arm-2    ● vortex-arm-1

**disabled**

## 5s timeout

**Producer buffer availability**
Sep 6, 4:50pm – Sep 6, 5:10pm

70 MB
60 MB
50 MB
40 MB
30 MB
20 MB
10 MB
0 B

4:50pm    4:55pm    5:00pm    5:05pm

● vortex-arm-2    ● vortex-arm-1    ● vortex-arm-0

new relic.

# Partition availability configuration continued

- Value too low can result in too much flapping with usable brokers
  - Less throughput (Kafka specification is 10MB/s per partition)
- Value too high can result in still exhausting the send buffer

new relic.

# Recovery requires even more throughput



**Bytes in per second (by broker)**
Since 1 hour ago

# More throughput

- Analysed load
  - CPU had headroom - no need for more nodes
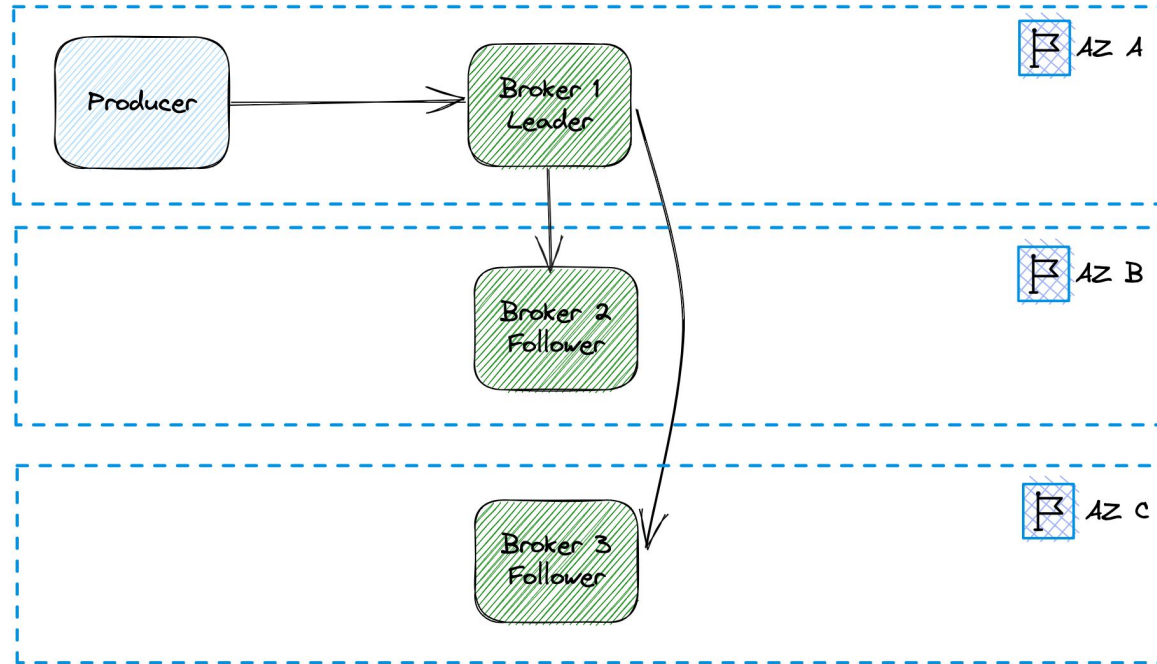  - Disk write throughput could be improved

new relic.
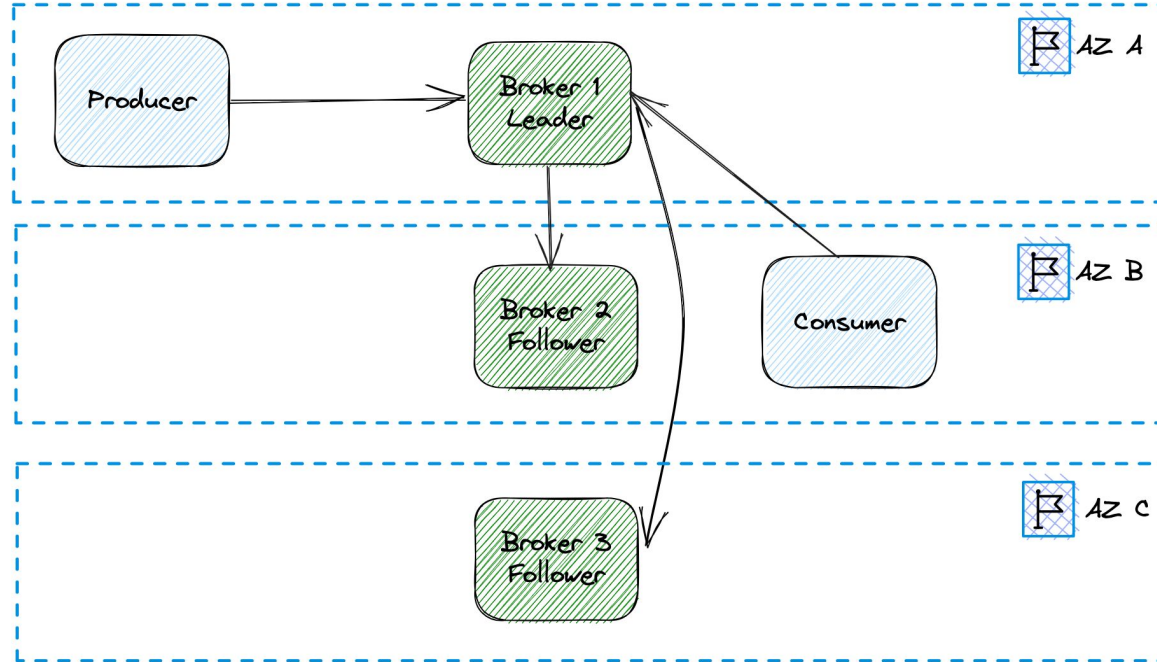
# More throughput continued

- Originally we could only use EBS gp2 volumes
  - Limited to 250MB/s write throughput
- New EBS gp3 provisioned throughput costs 30-40x less compared to extra EC2 instances with EBS gp2 volumes

new relic.

# Does this solve all gray failures in Kafka?
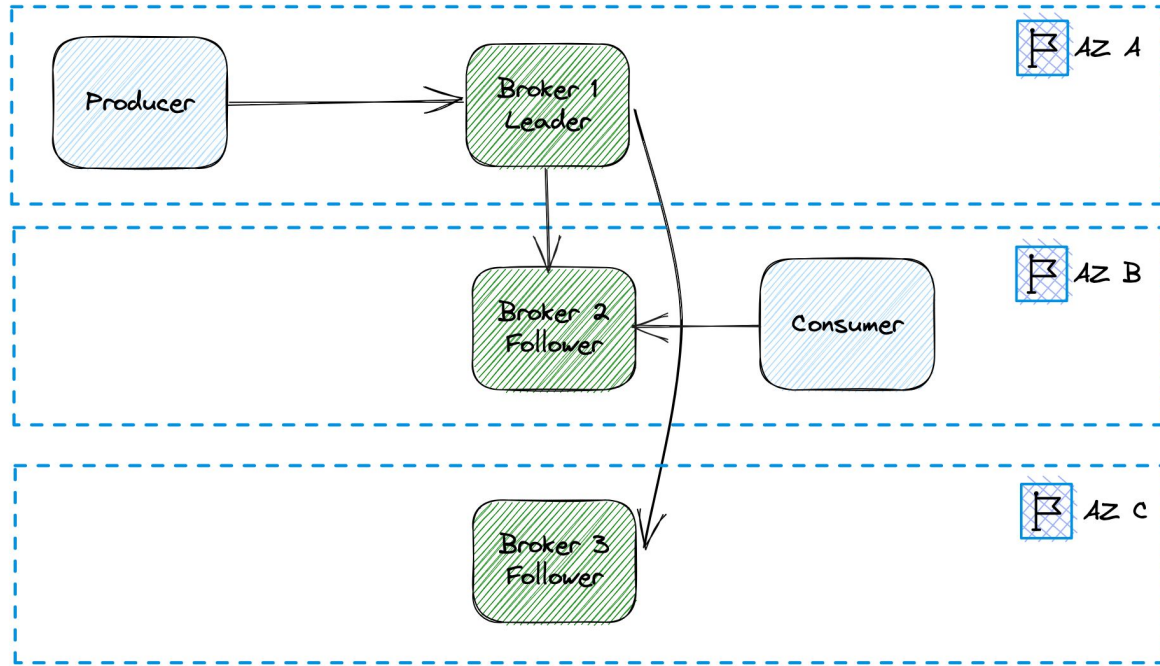
new relic.

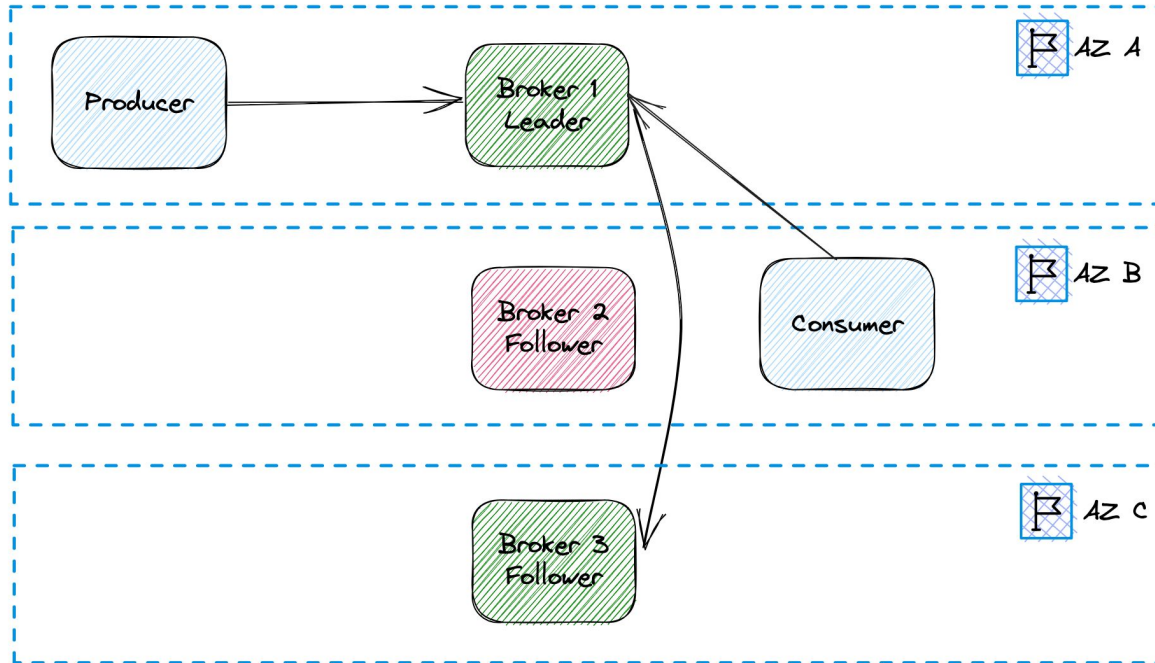# High availability + durability with replication factor 3

# Data path from producer to consumer

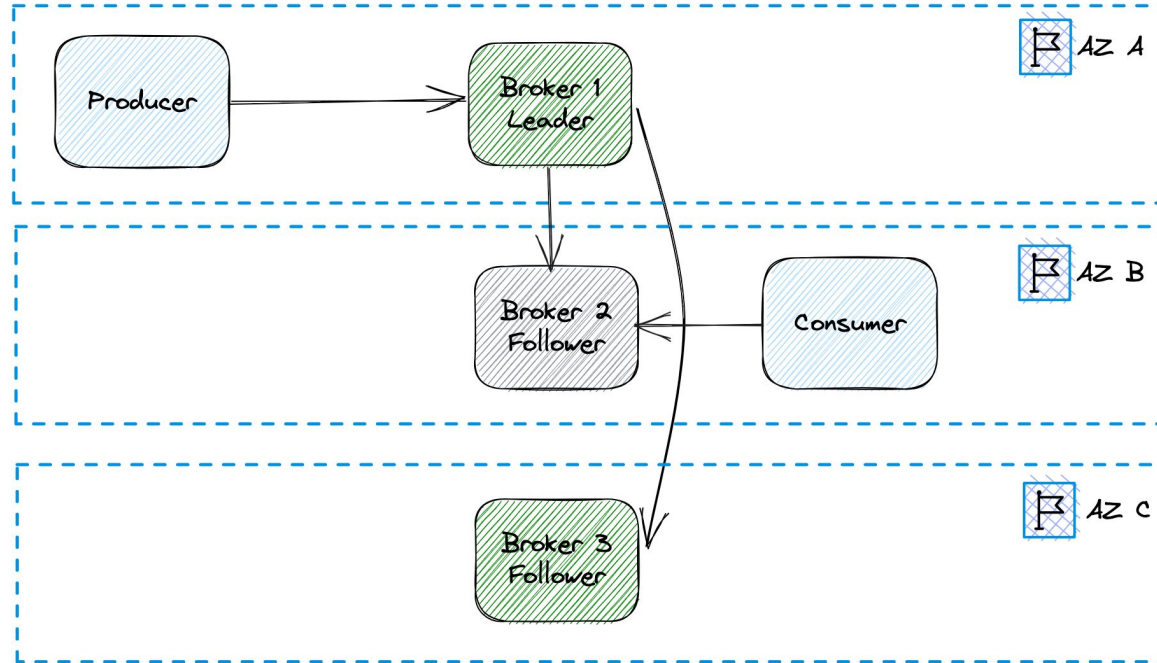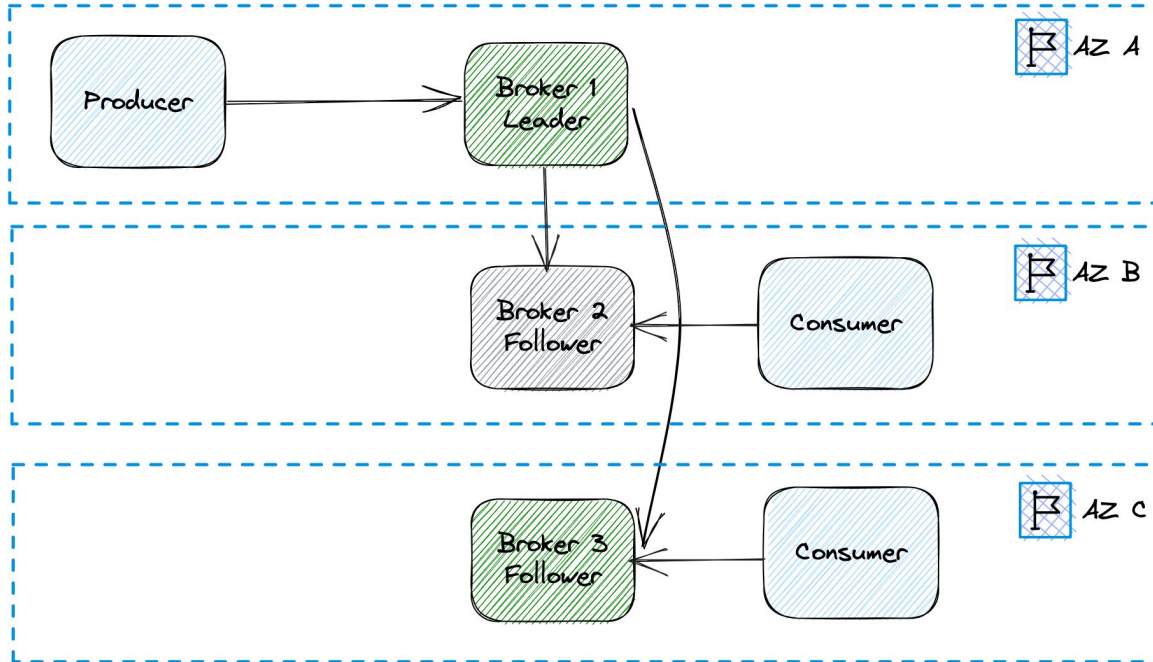# Fetch from closest replica / less cross-AZ data

# Broker in AZ B has a problem

new relic.

# Broker in AZ B recovers but it's still out of sync

# Other AZs are NOT affected

# Why not disable fetch from closest replica?

new relic

# Data path from producer to consumer

# Fetch from closest replica / less cross-AZ data

# Peak daily ingest 85GB/s

new relic

# Cross AZ traffic is only about ⅔ of all

# Why we didn't disable fetch from closest replica

- Finance would still be unhappy even with "just" 57GB/s
- The fix was already in the works
- Can be disabled either in client or server
  - But 95 clusters
  - And close to 300 topics in busiest clusters

# What's needed

- Kafka clients 3.3.2 / 3.4.0 - 8 months ago
- Kafka server 3.3.2 / 3.4.0
- Either helps, but full fix requires both

new relic.

# Alternative - if you can't upgrade

- If you are OK with more cross-AZ traffic
  - disable fetch from closest replica

new relic.

# Does this solve all problems?
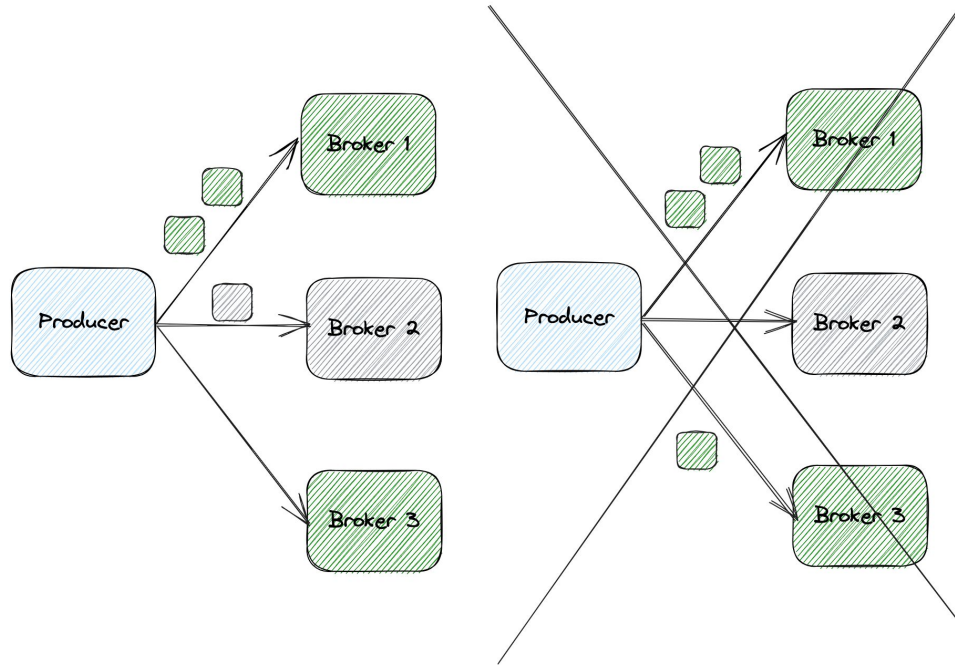
new relic

# Producing with a partition key

new relic.

# Takeaways - producers with random partitioning

- Use Kafka clients 3.3.1 - 1 year ago
- Use the Strictly Uniform Sticky Partitioner
  - It's the default - don't override the partition class
  - Improves send buffer exhaustion on its own
  - Even better with partitioner.timeout.availability.ms

# Takeaways continued

- Use Kafka clients 3.3.2 / 3.4.0 for consumers
  - 8 months ago
- Use Kafka server 3.3.2 / 3.4.0


Alternative:

- If you can't upgrade but you are OK with more cross-AZ traffic - disable fetch from closest replica

new relic.

# Takeaways SRE

- Gray failures are hard to deal with
- Gamedays to reproduce gray failures
- Analyse system architecture
- Plan capacity according to your needs

new relic.

# Special thanks

- Alex Thengumpalli

- Anton Rodriguez

- Luke Kirby

- Christopher Wildman

- Alex Lindeman

- Kafka Platform Team

- Streaming SRE team

new relic.

# new relic

# Thank you.

https://www.linkedin.com/in/michellevalentinova/