# Migrating a Large Scale Search Dataset in Production in a Highly Available Manner

Leila Vayghan
Infrastructure Engineer

**shopify**

# **Agenda**

- Shopify
- Search Infrastructure
- Data Migration
- Challenges and Solutions
- Q/A

# Shopify

**Global commerce platform,** helping merchants do business online and offline.

# Shopify

**Millions**
of merchants worldwide
GymShark, FashionNova, FC Barcelona

**170+**
countries represented

**10%**
of total US commerce

**$444B**
global commerce activity

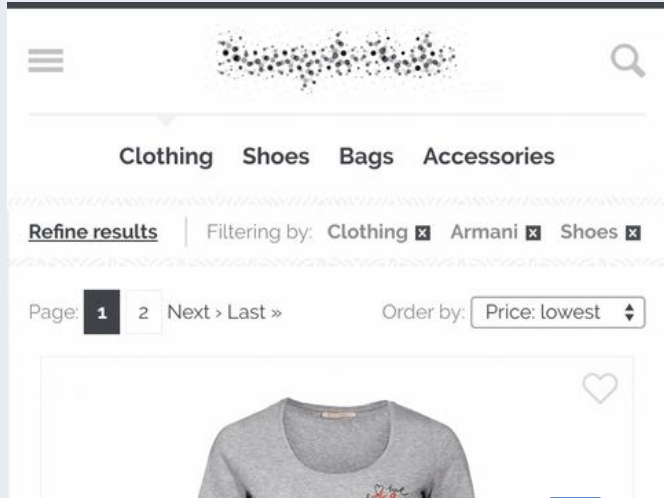# Shopify

**tobi lutke** ✔ 🛍
@tobi

Nerd BFCM stats:

Shopify's egress processed 145 billion requests on Friday. App servers handled peak of ~60 million requests per minute. Increase of 38%. Total GMV was $4.1b, up by 22% from last year.

# Search

Clothing   Shoes   Bags   Accessories

Refine results   Filtering by:  Clothing ☒  Armani ☒  Shoes ☒

Page: 1  2  Next › Last »          Order by: Price: lowest

**Search Engine**

# Elasticsearch

Distributed search and analytics engine

Built on top of Apache Lucene open-source search library

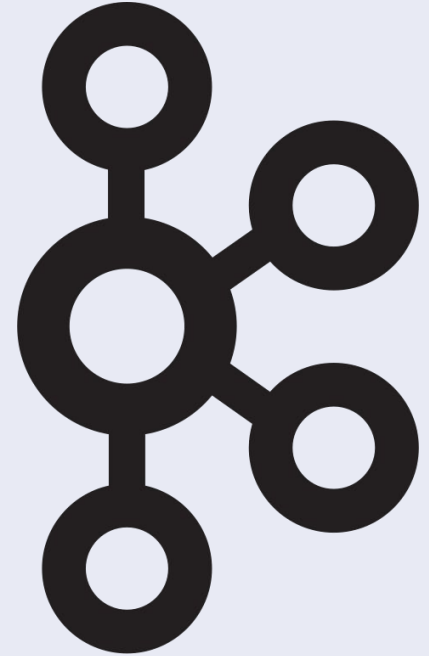Well suited for analyzing and searching large volumes of data
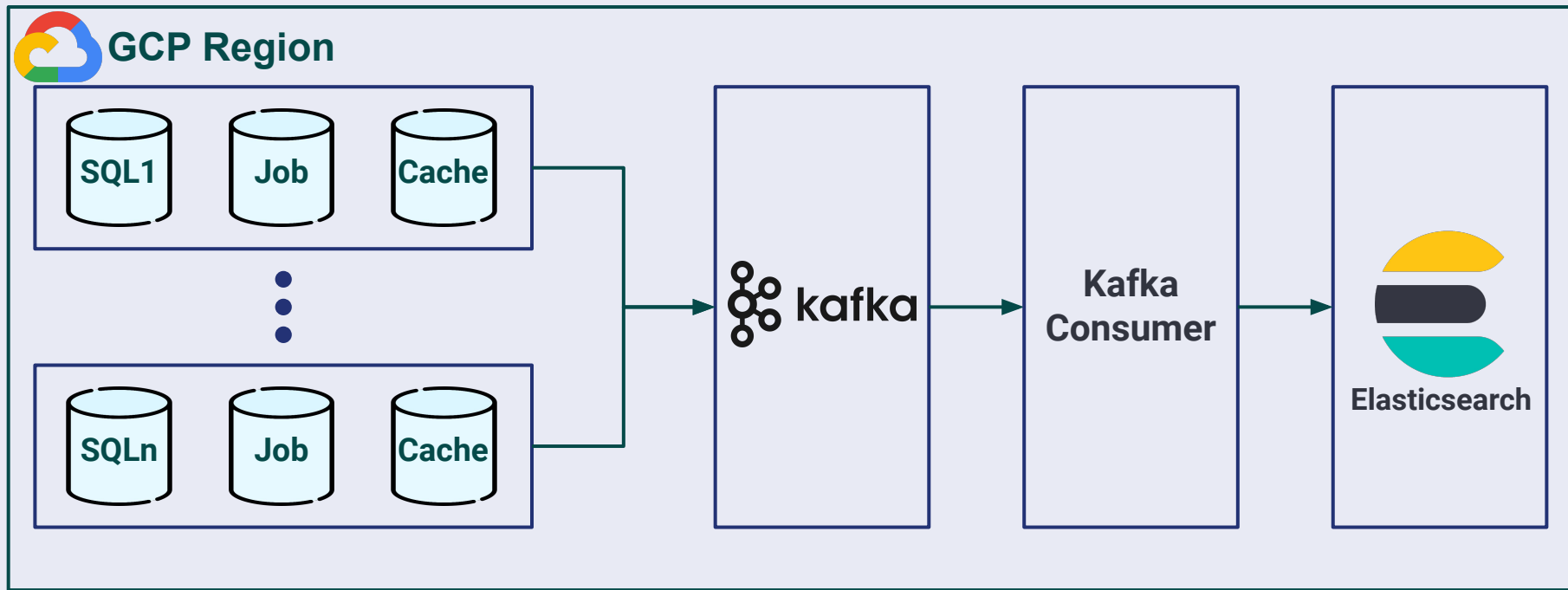
# Apache Kafka

Distributed event streaming platform
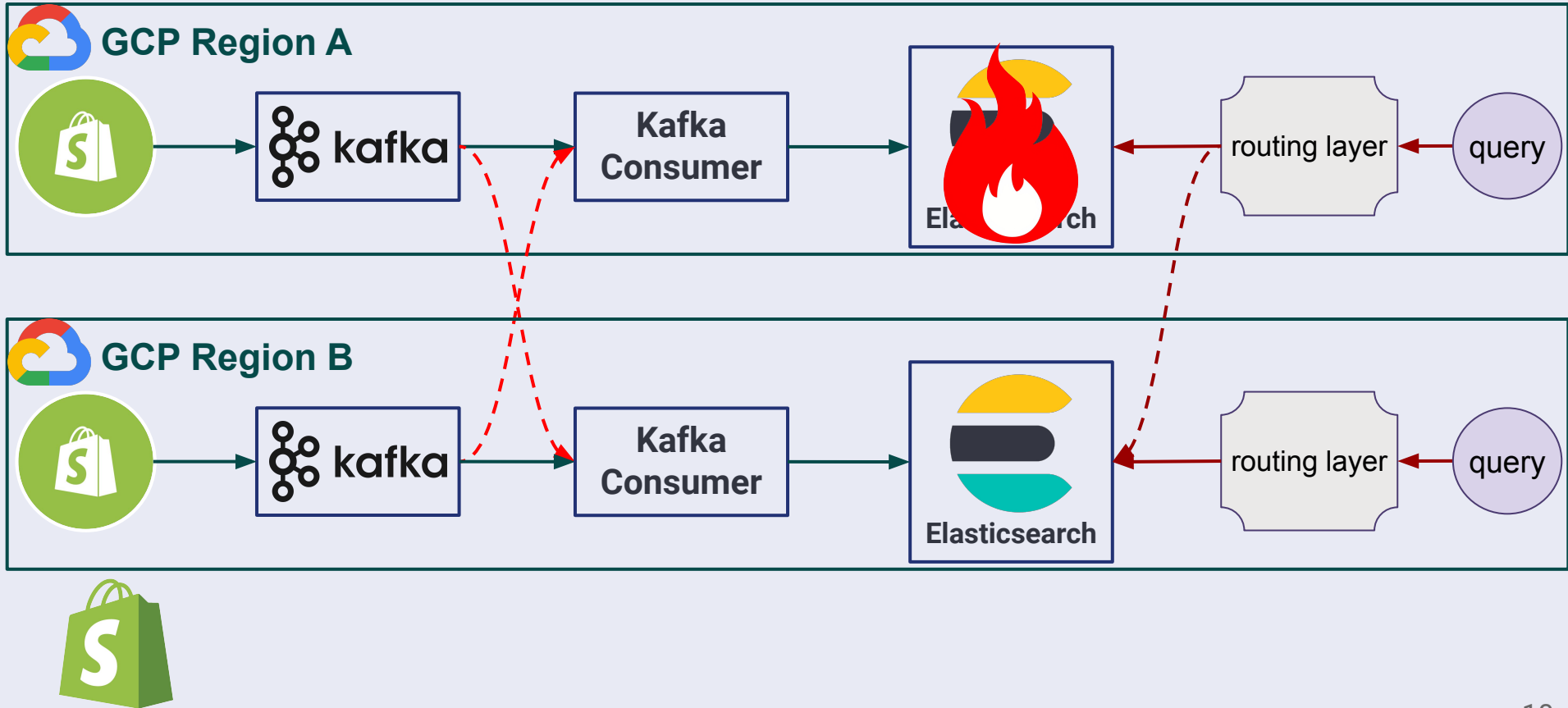
Used for building high performance realtime data pipelines

Shopify's primary messaging service, carrying messages between systems like MySQL & Elasticsearch
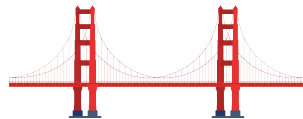
# Search at Shopify

GCP Region

SQL1 · Job · Cache

SQLn · Job · Cache

kafka → Kafka Consumer → Elasticsearch

# Search at Shopify

# Search at Shopify

- Elasticsearch is a secondary data store
- Indexing is the act of writing documents from the primary data store to Elasticsearch
- We have built two indexing pipelines with SQL as the primary data store
  - **Realtime pipeline:** captures changes to SQL records
  - **Reindex pipeline:** rebuilds an entire Elasticsearch index from SQL (maintenance)

# Search at Shopify

**90K**

**docs/sec
Realtime indexing rate**

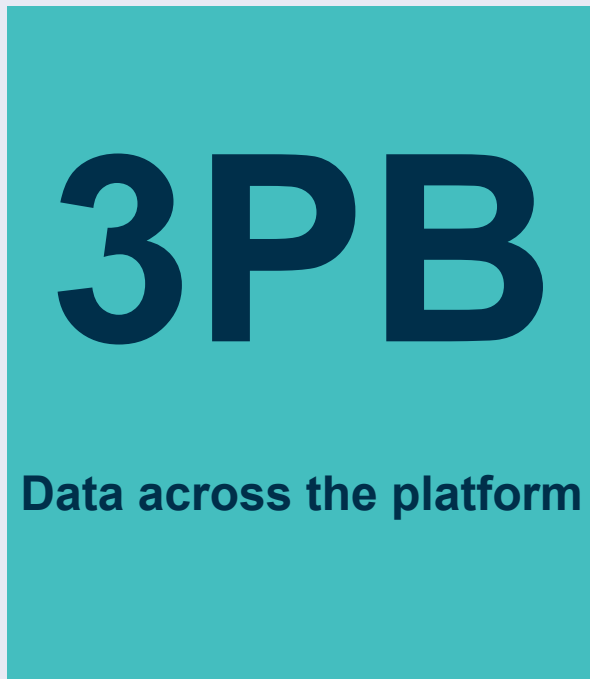**500K**

**docs/sec
Reindex indexing rate**
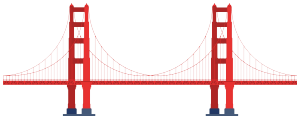
Data based on numbers as of March 2024

# Search at Shopify

~100 distinct Elasticsearch clusters

The platform supports large clusters with as many
as 260 nodes to small ones with only 3 nodes.

**3PB**

**Data across the platform**

Data based on numbers as of March 2024
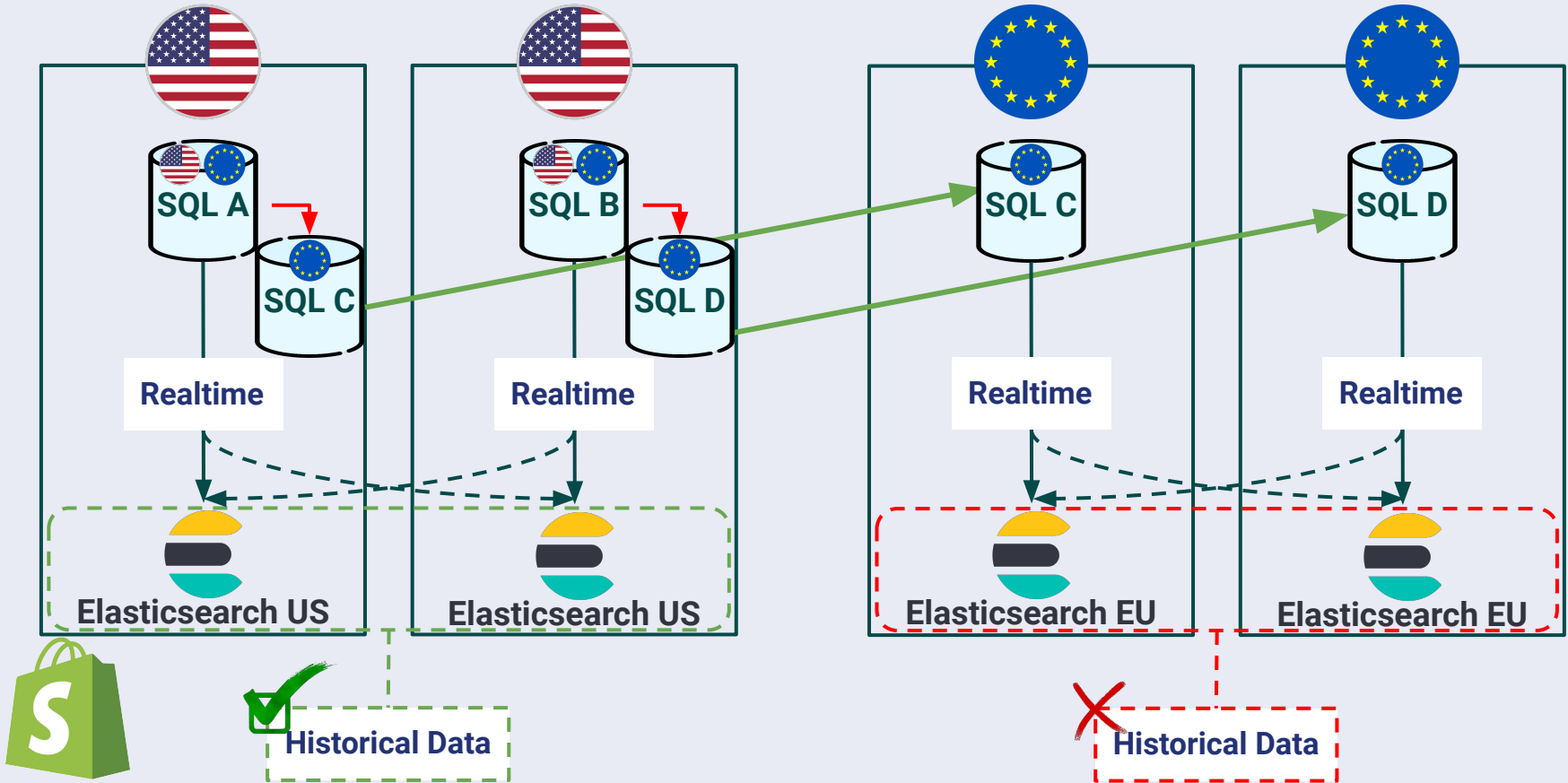
# Shopify, Before Covid

# Shopify, After Covid

- Number of merchants and buyers increased across the globe
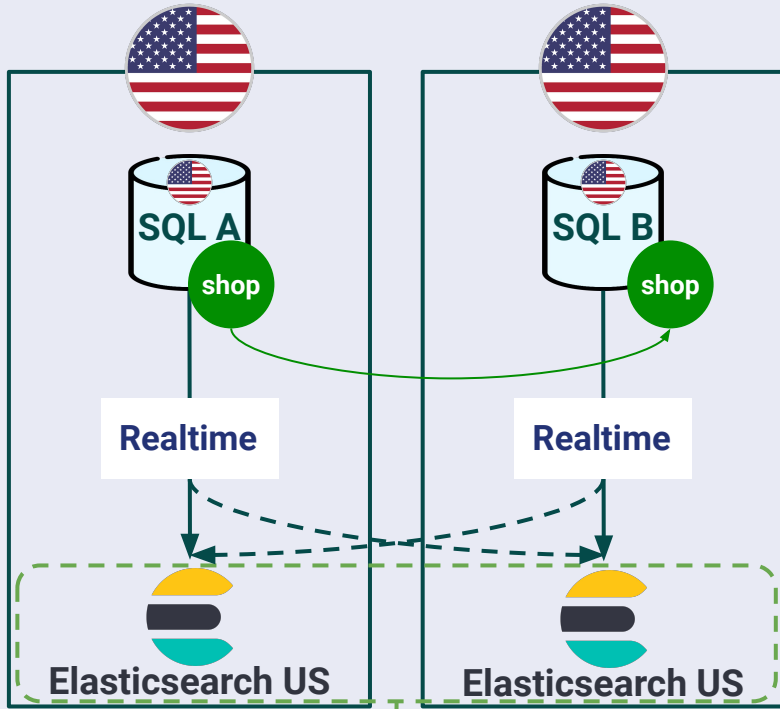  - Latency and data locality preferences became constraints
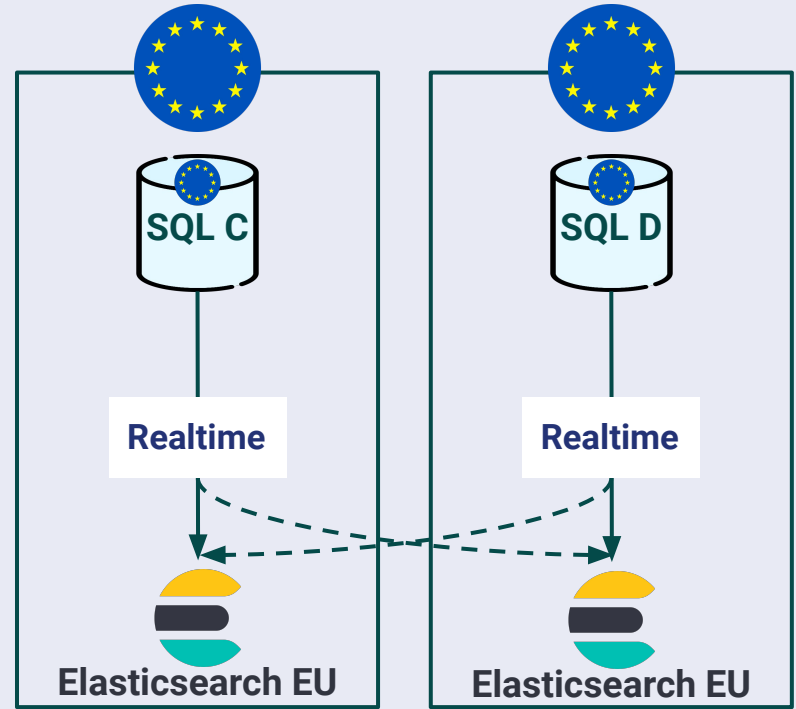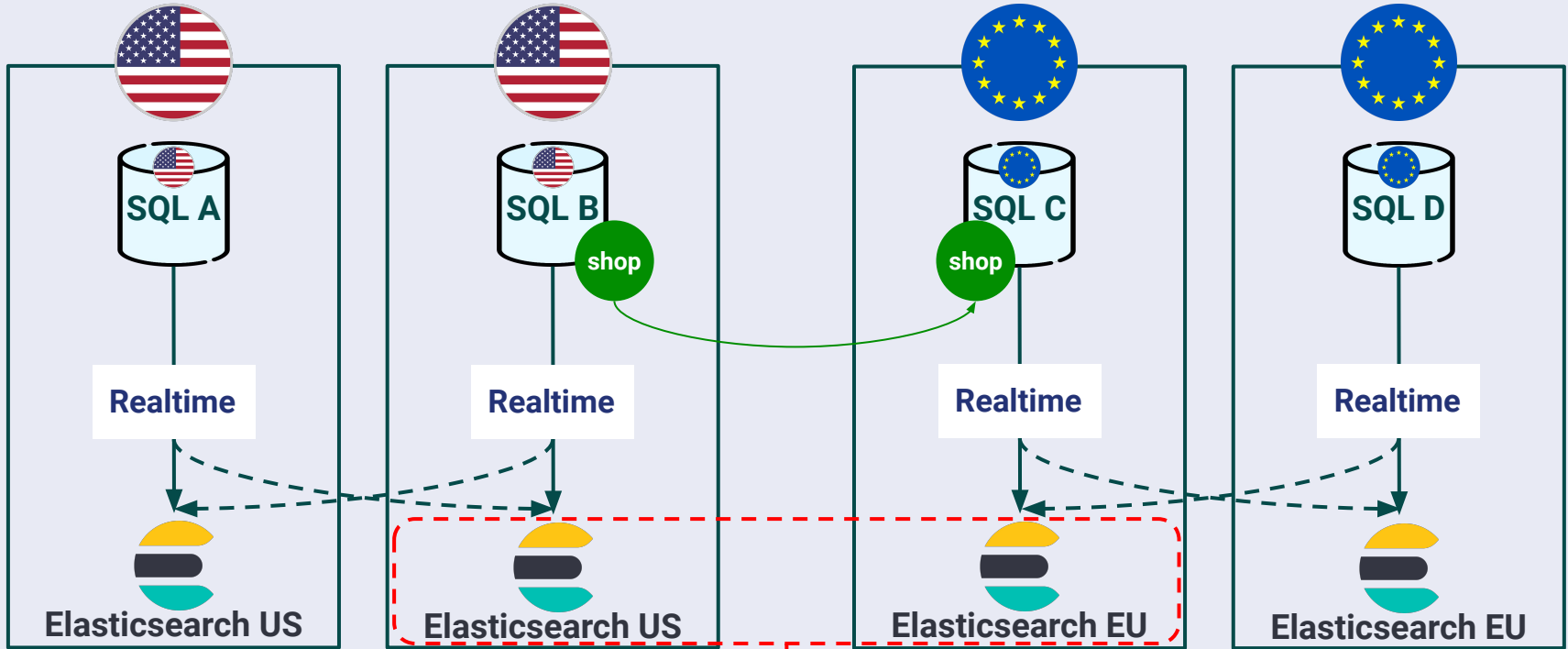
Shopify decided to expand to Europe

# Shop Moves

SQL A

shop

SQL B

shop

SQL C

SQL D

Realtime

Realtime

Realtime

Realtime

Elasticsearch US

Elasticsearch US

Elasticsearch EU

Elasticsearch EU

Same Dataset

# Shop Moves

SQL A   SQL B   SQL C   SQL D

shop   shop

Realtime   Realtime   Realtime   Realtime

Elasticsearch US   Elasticsearch US   Elasticsearch EU   Elasticsearch EU
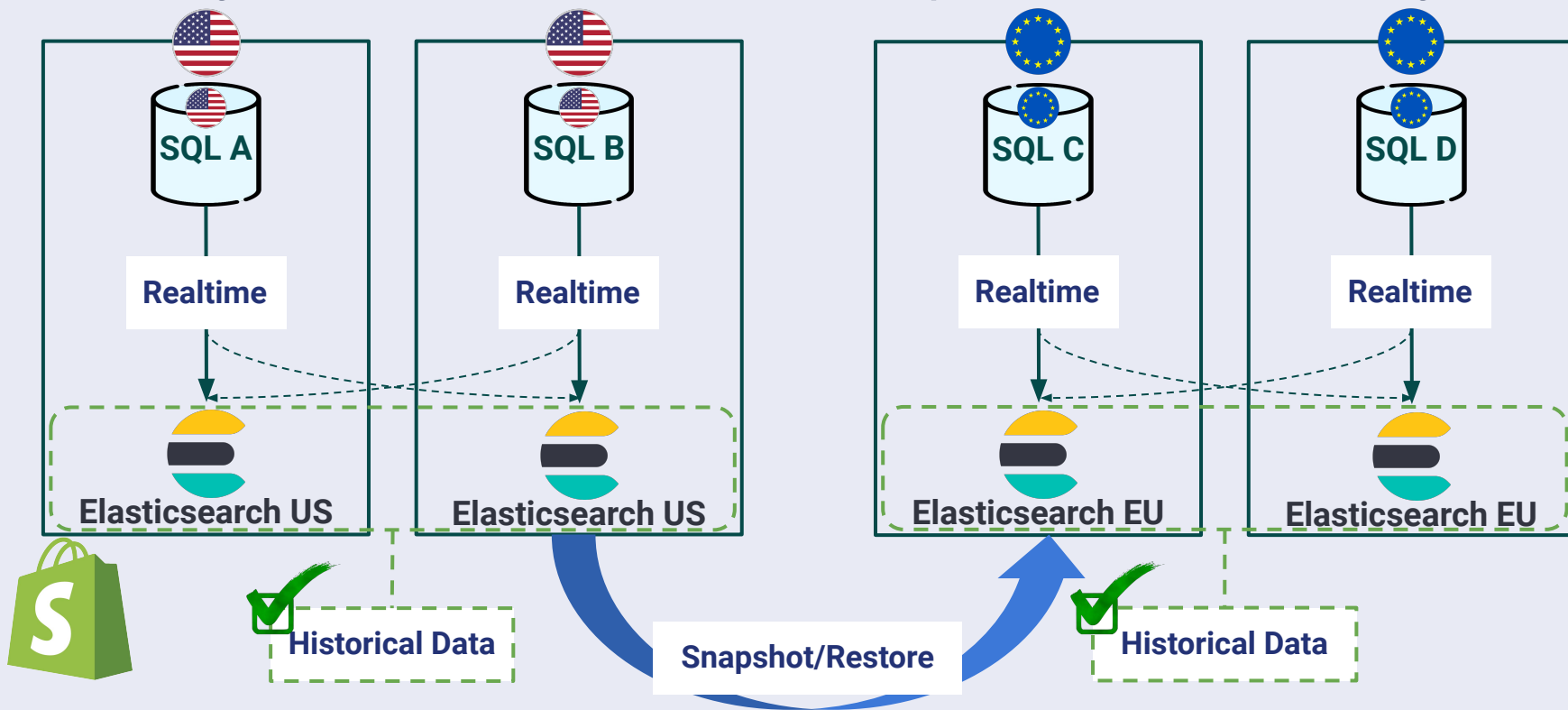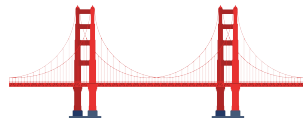
Different Datasets → shop missing data

# Solution A

Restoring historical data from a US snapshot in EU **after** migration

# Solution A

## Restoring historical data from a US snapshot in EU after migration

✓ Pros
➢  Simple operation

✗ Cons
➢  Search downtime during snapshot restore leading to revenue loss

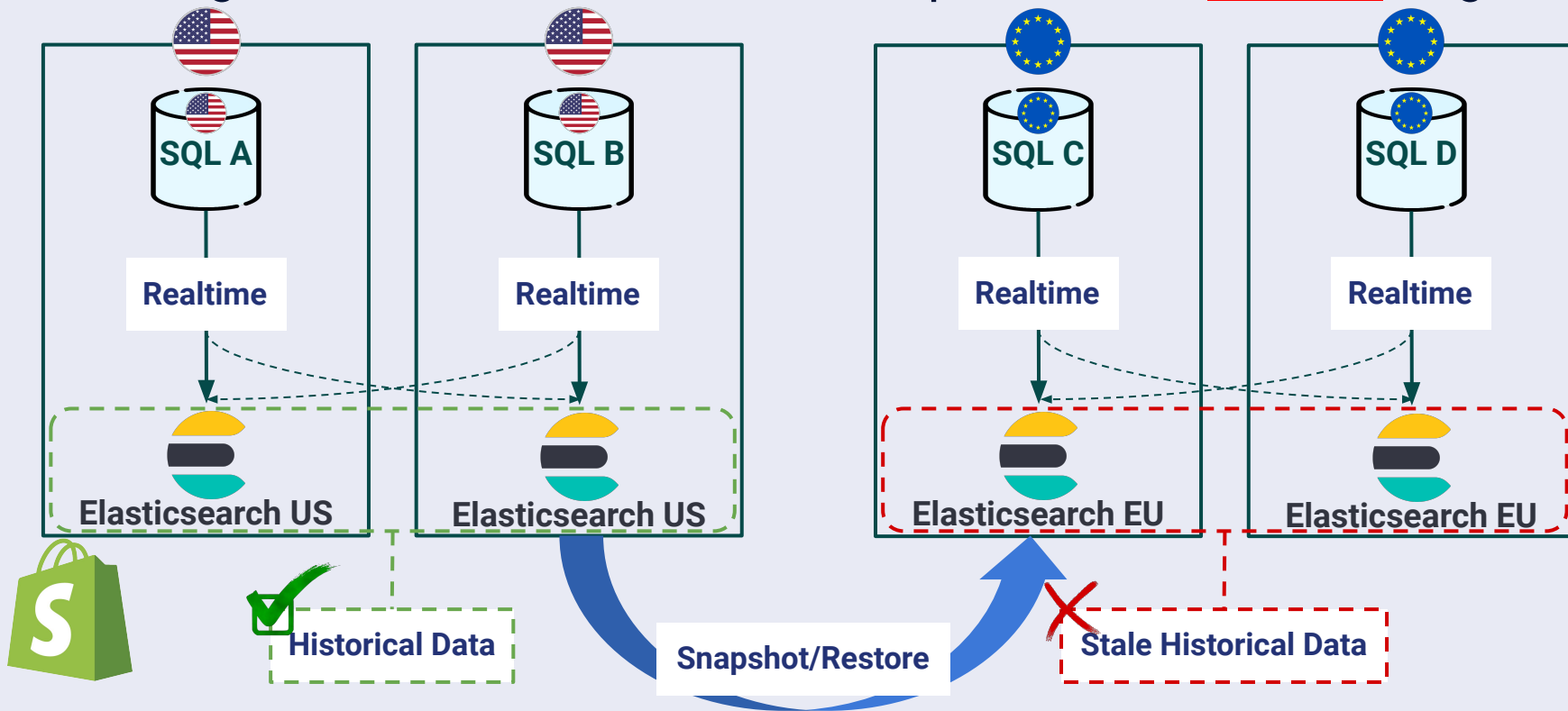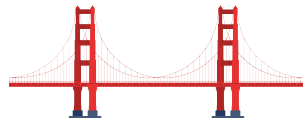➢  Search downtime for individual shops moving between jurisdictions

REJECTED

20

# Solution B

Restoring historical data from a US snapshot in EU **before** migration



SQL A
SQL B
SQL C
SQL D

Realtime
Realtime
Realtime
Realtime

Elasticsearch US
Elasticsearch US
Elasticsearch EU
Elasticsearch EU

Historical Data

Snapshot/Restore

Stale Historical Data

# Solution B

## Restoring historical data from a US snapshot in EU before migration

**Pros**

➤ Simple operation

**Cons**

➤ Realtime updates missing from search leading to revenue loss

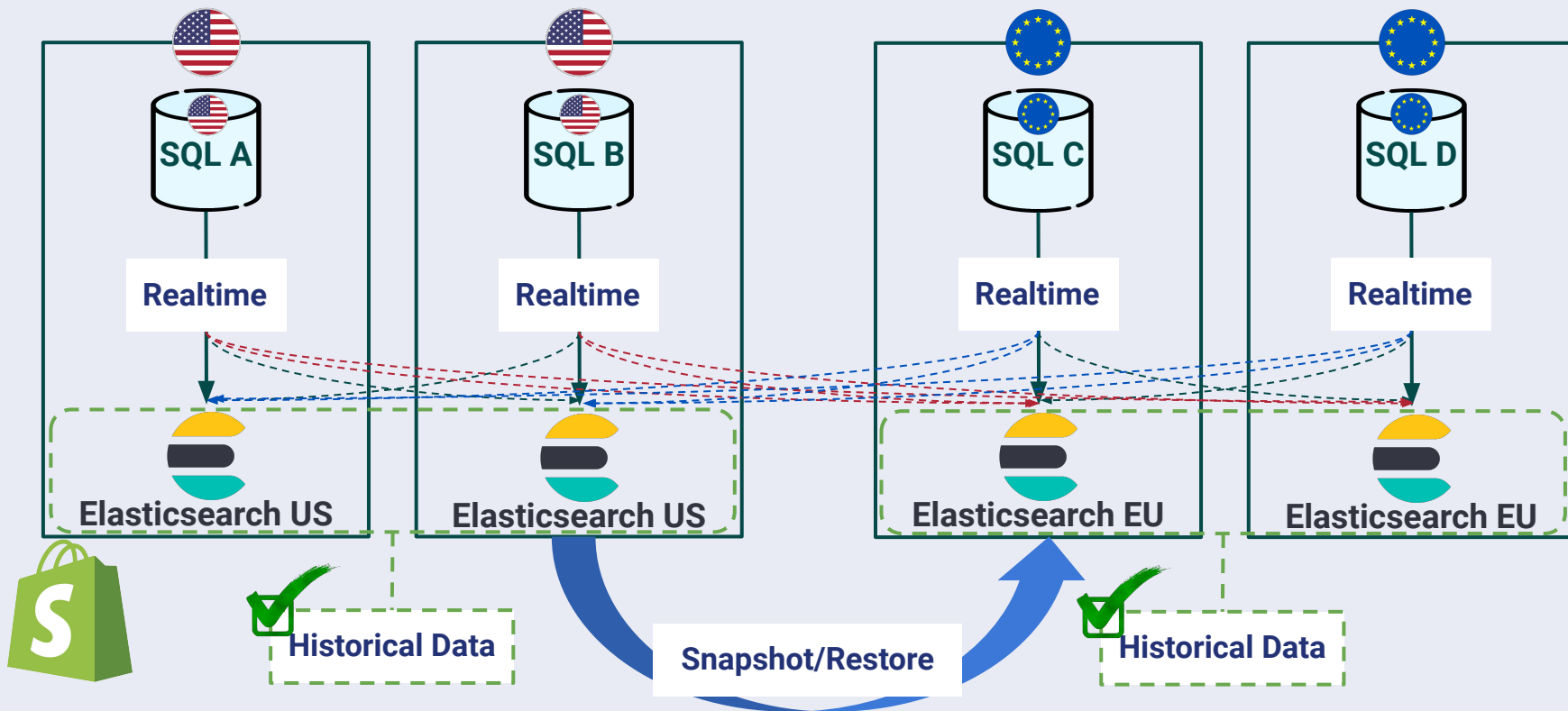➤ Search downtime for individual shops moving between jurisdictions

REJECTED

22

# Solution C

## Cross jurisdictional search data replication

# Solution C



## Cross jurisdictional search data replication

### Pros

➢ No revenue loss due to search downtime

➢ Decoupling from other stateful systems migration plans

### Cons

➢ Infrastructure cost

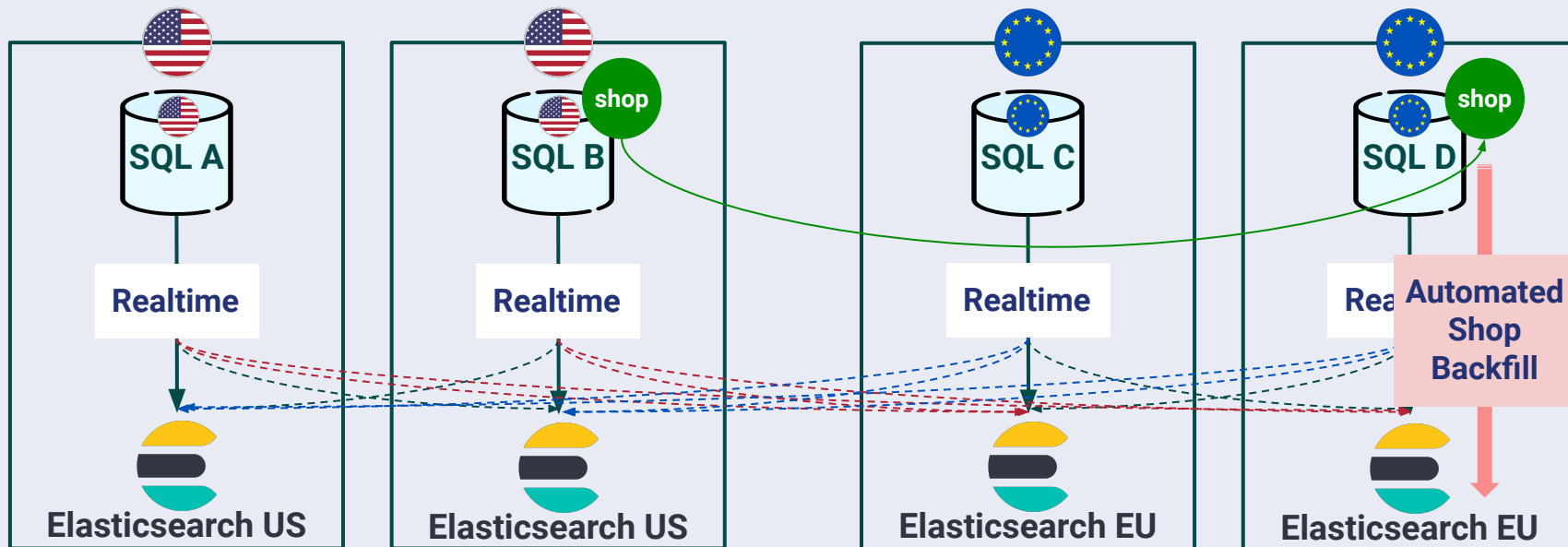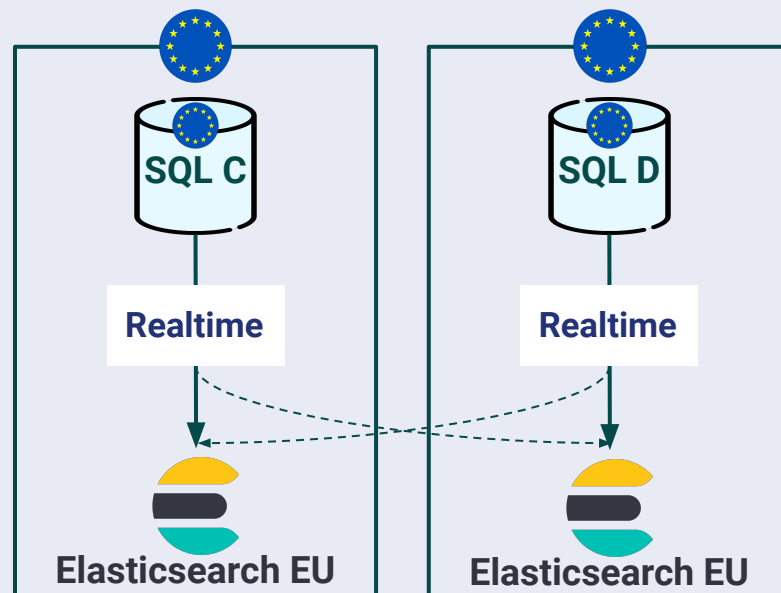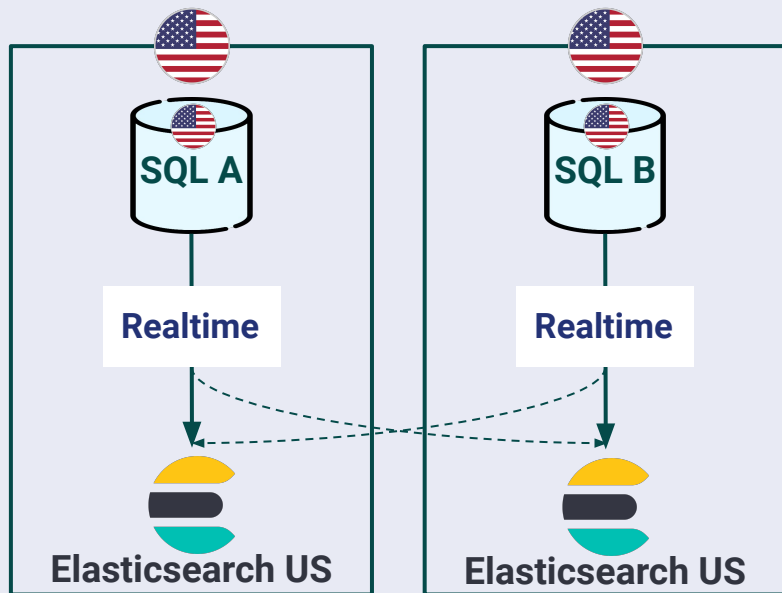➢ Complex infrastructure

# Temporary Topology

Cross jurisdictional search data replication

# Final Topology

# Shopify, Today