# Caching in Cloud

# Why cache?



Improved Performance
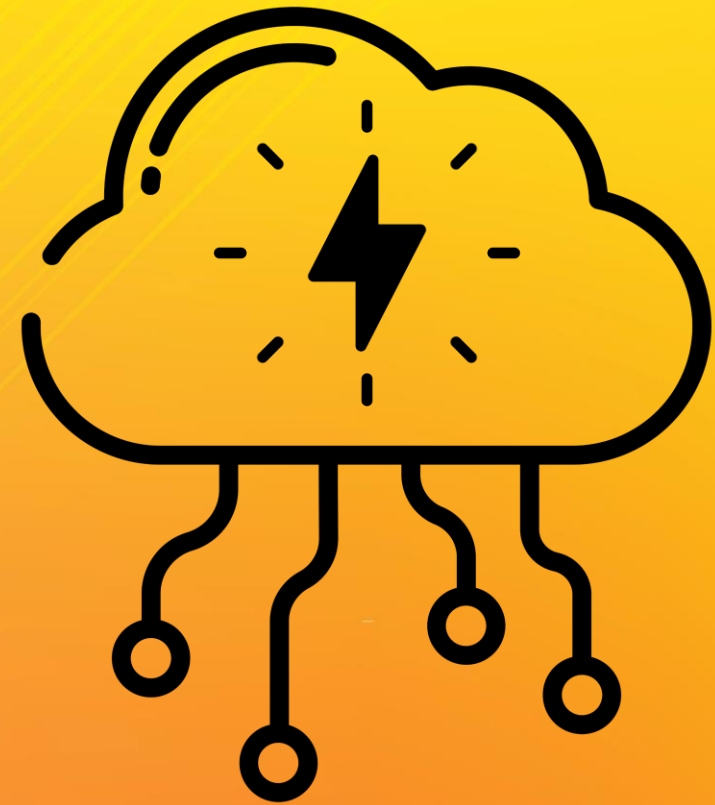
Additional Security
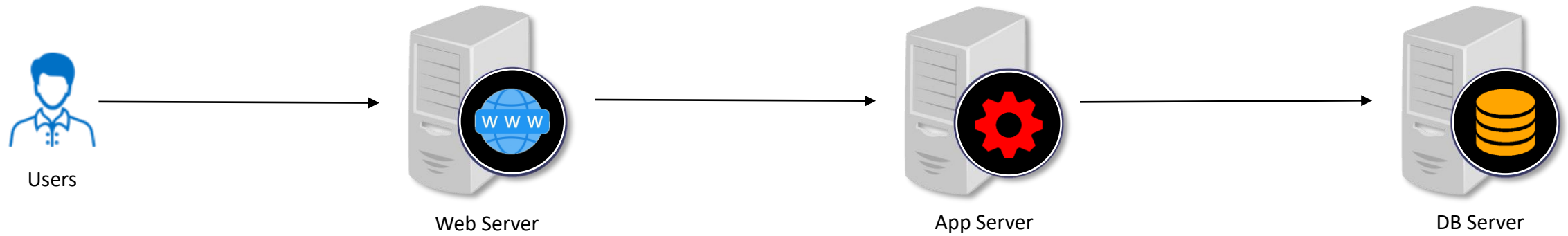
Cost Savings
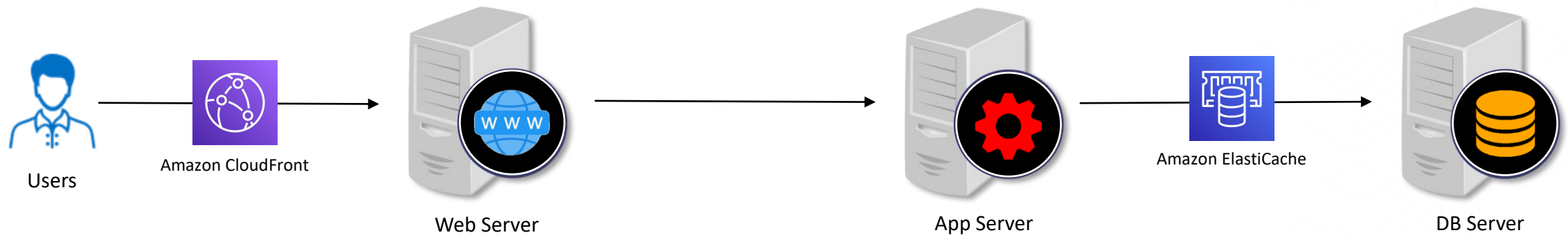
Caching in AWS

# A Typical Web Application



Users → Web Server → App Server → DB Server

# Caching on AWS



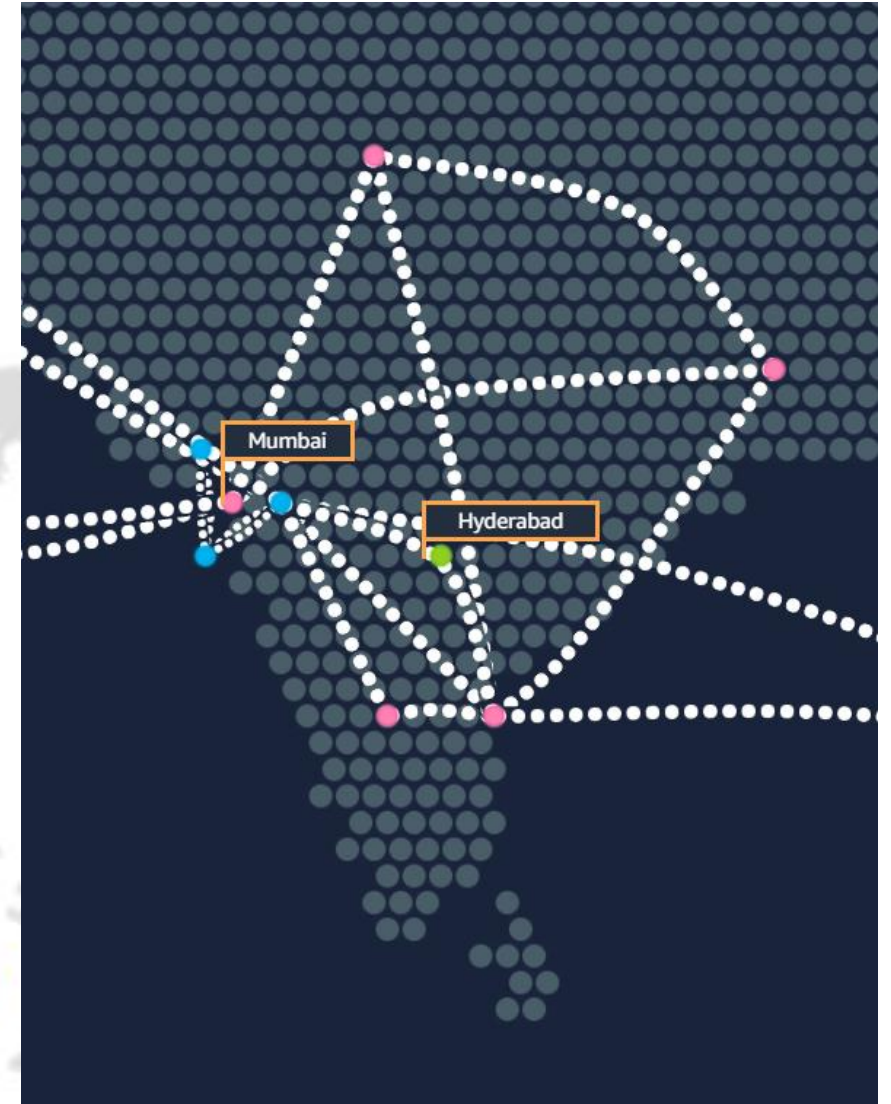Users → Amazon CloudFront → Web Server → App Server → Amazon ElastiCache → DB Server

Amazon CloudFront

# Amazon CloudFront – Global Content Delivery Network

- Uses Edge Locations for caching

# Region vs. Edge Locations (Analogy)

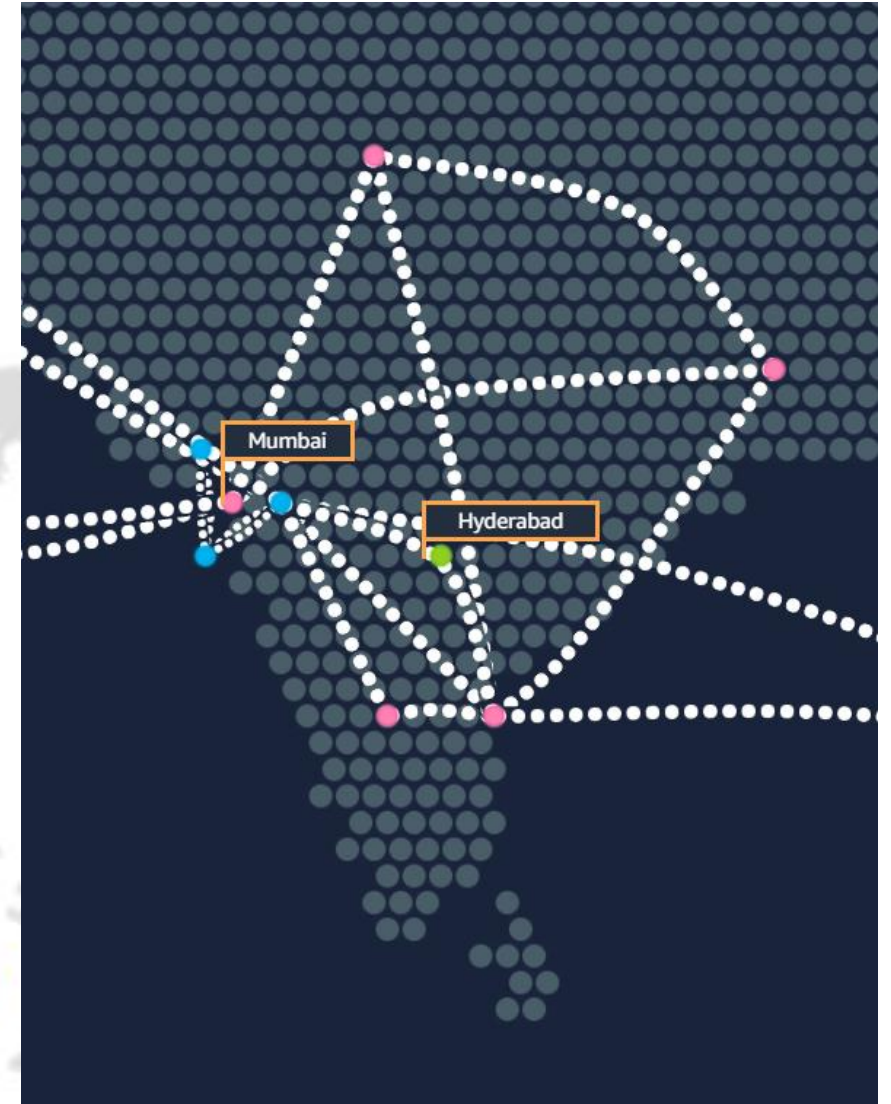**Large-format stores**
(Everything under one roof)

**Small-format stores**
(everyday essentials)
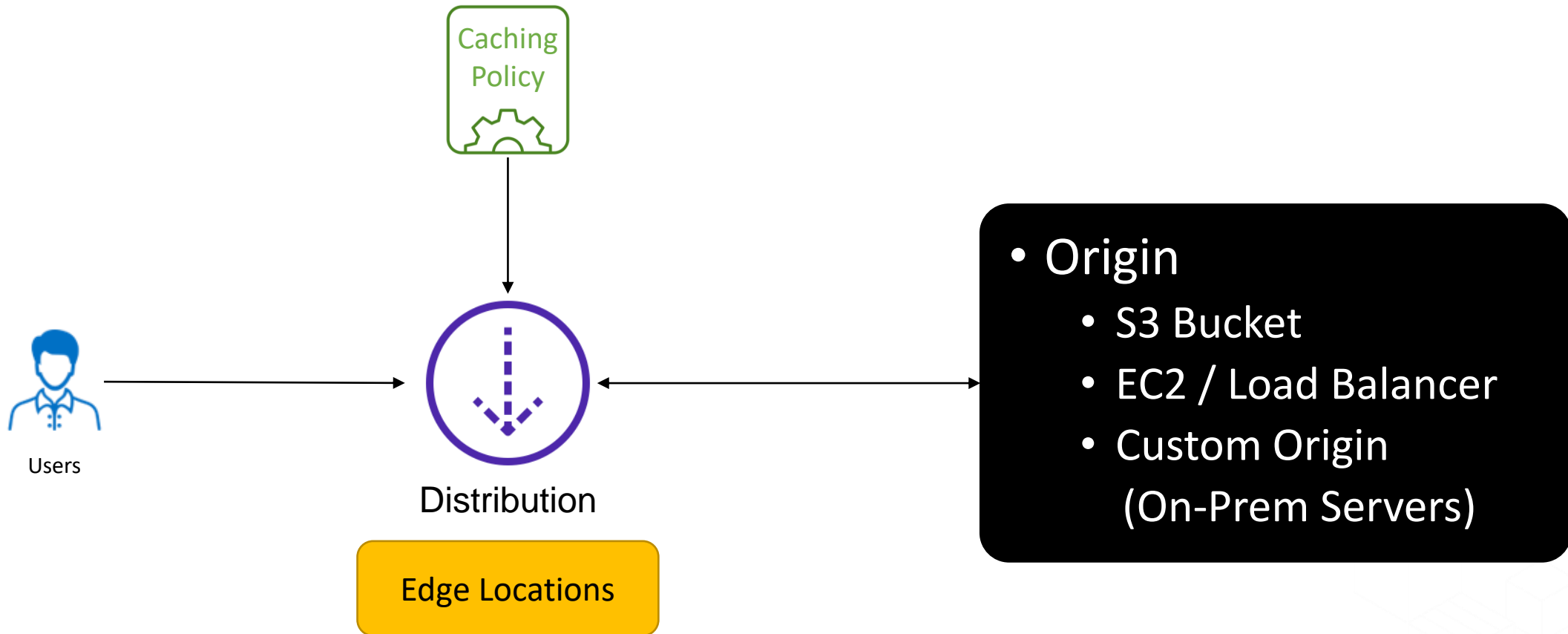
# Amazon CloudFront

- Uses Edge Locations for caching

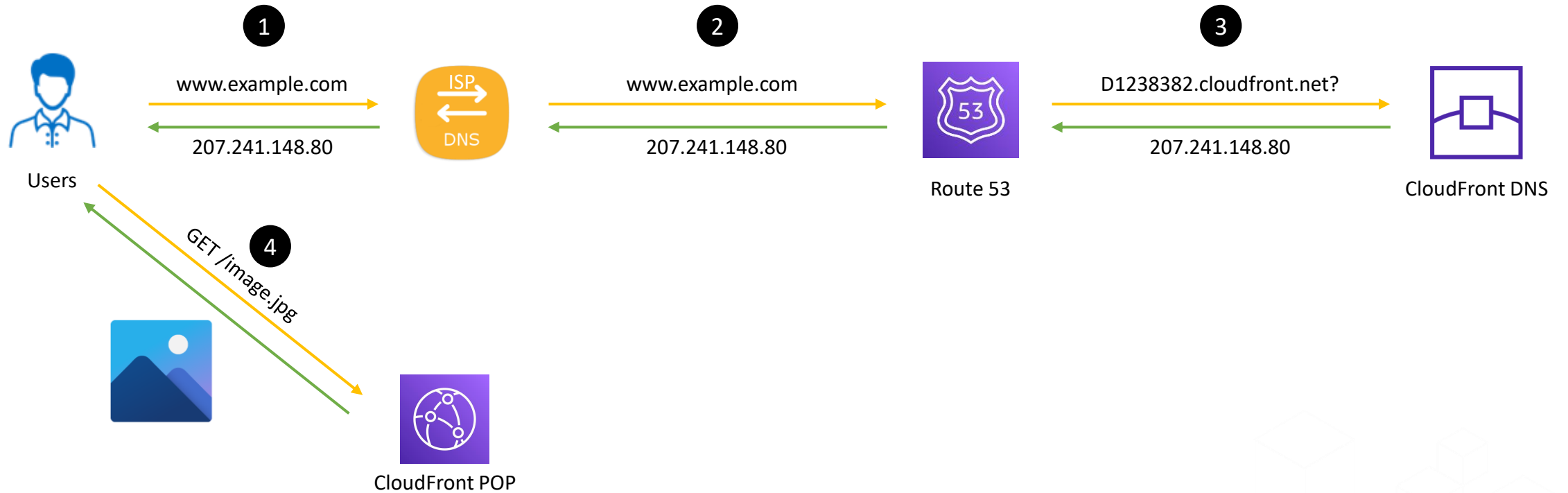- DDoS Protection

- Lambda@Edge

# Components



Caching Policy

Users

Distribution

Edge Locations

- Origin
  - S3 Bucket
  - EC2 / Load Balancer
  - Custom Origin
    (On-Prem Servers)

# HTTP Request – www.example.com/image.jpg

Amazon CloudFront

# How Amazon CloudFront works?

Origin (Source)

Amazon CloudFront

Get →

← Ok

1st Request

Cache Miss

Get →

← Ok

Amazon S3 Bucket

Index.html

[CSS]

[Images]

Ok ↓    ↑ Get

Subsequent Requests

Cache Hit

Database Caching

# Database Caching

# In-Line Cache vs Side Cache

## In-Line Cache

Application ←→ Cache ←→ Database

## Side Cache

Application ←→ Cache

Application ←→ Database

# Database Caching

- Open source caching engine
  - Memcached
  - Redis

Amazon Elasticache

# Amazon Elasticache


Fully Managed


Extreme Performance


Scalable


Secure & Compliant

# Amazon Elasticache – Engine choice

- Amazon Elasticache for Memcached

- Amazon Elasticache for Redis



ElastiCache for Memcached

Cluster



ElastiCache for Redis

Cluster

# Memcached vs. Redis

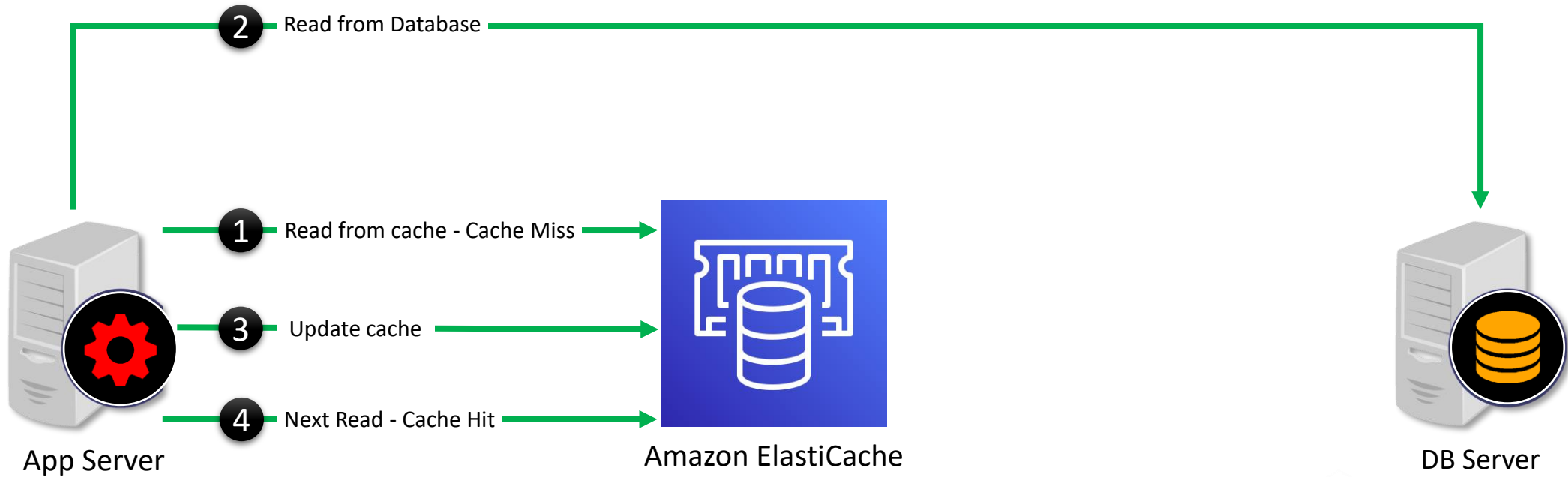| | Memcached | Redis |
|---|---|---|
| **Primary benefit** | Simplicity | Rich set of features |
| **Advanced data structures** | Not supported | Supported – strings, lists, sets, sorted sets, hashes, bit arrays, and hyperloglogs. |
| **High availability (Replication)** | Not supported | Supported |
| **Backup and Restore (Data Persistence)** | Not supported | Supported |
| **Pub/Sub capabilities** | Not supported | Supported |
| **Transactions** | Not supported | Supported |

Caching Strategies

# Lazy Loading Strategy



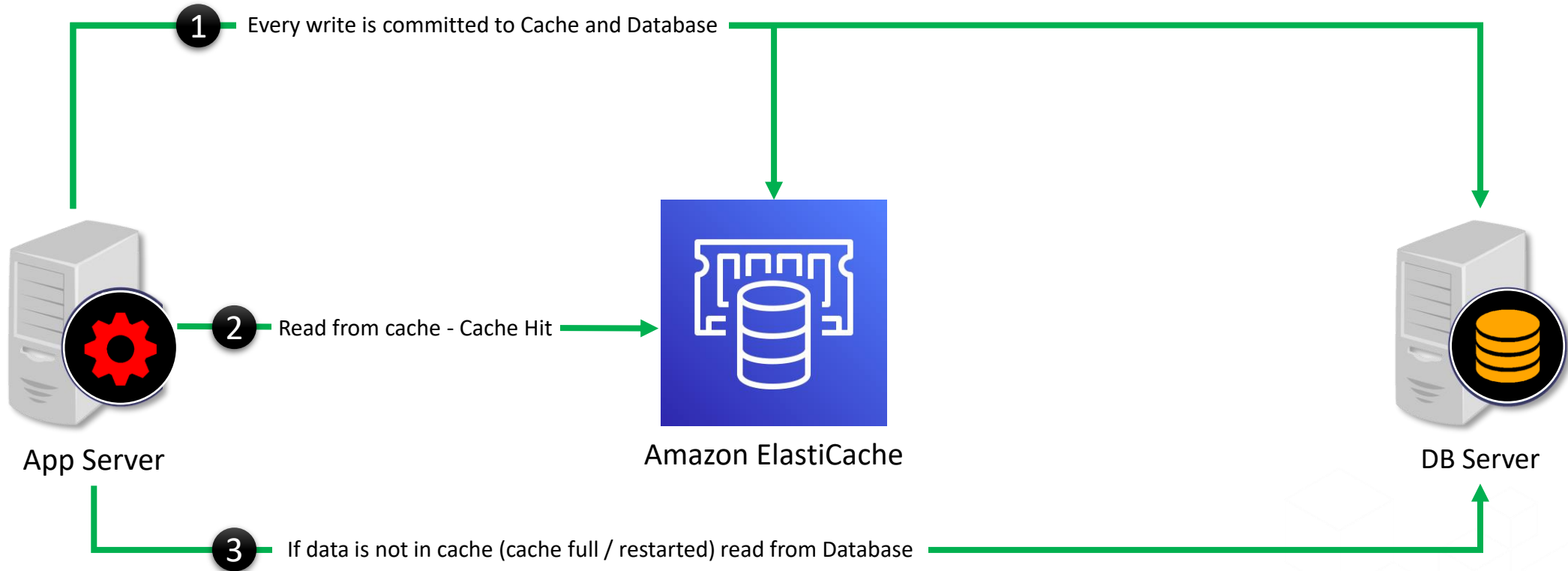| Advantage |
| --- |
| Only requested data is cached |
| Node failures are not fatal |

| Disadvantage |
| --- |
| Cache miss penalty. Each cache miss results in 3 trips |
| Application may receive stale data |

# Write Through Strategy



**① Every write is committed to Cache and Database**

App Server

Amazon ElastiCache

DB Server

**② Read from cache - Cache Hit**

**③ If data is not in cache (cache full / restarted) read from Database**

| Advantage | Disadvantage |
| --- | --- |
| The data in the cache is never stale | Write penalty - Every write involves two trips |
| Write latency better tolerated by customers vs. read latency | Unused data in cache |

# Cache Expiry - Time-to-Live (TTL)

Value 1
TTL = 15 Min

Value 2
TTL = 30 Min

Value 3
TTL = 45 Min

App Server

Amazon ElastiCache

DB Server

Amazon DynamoDB Accelerator (DAX)

# In-Line Cache vs Side Cache

## In-Line Cache
## (DynamoDB Accelerator)

Application ↔ Cache ↔ Database
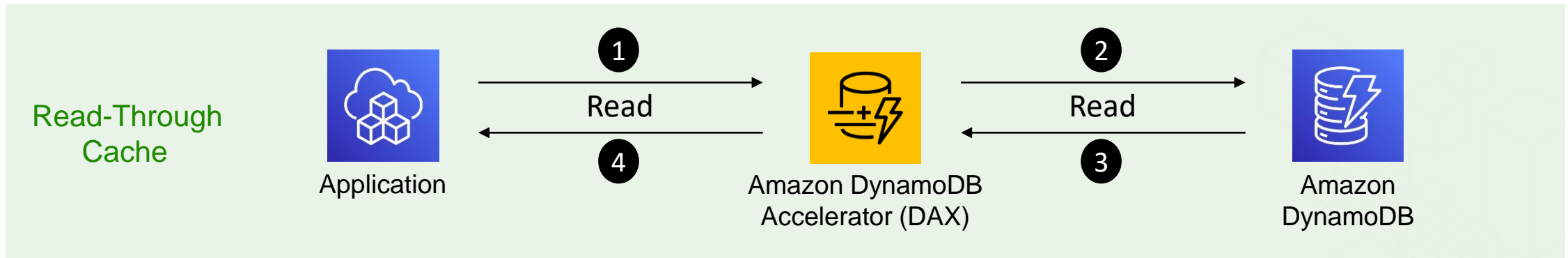
## Side Cache
## (Amazon Elasticache)

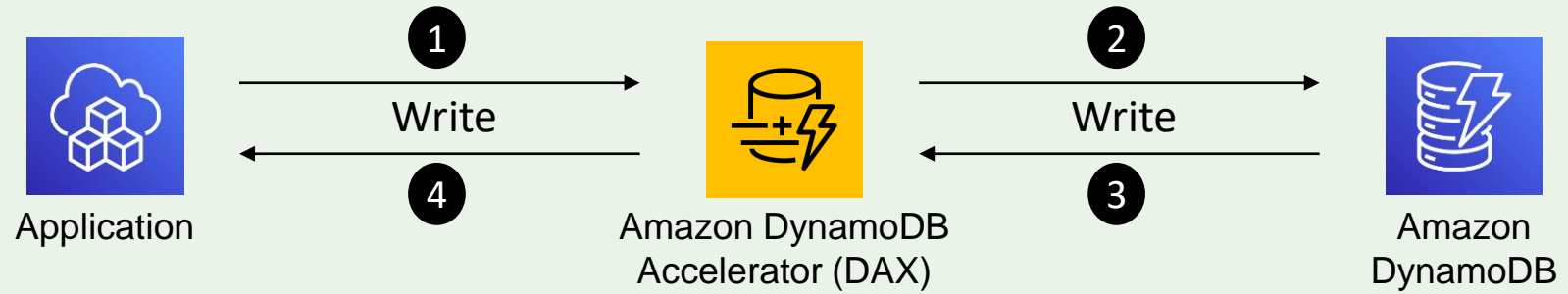Application ↔ Cache

Application ↔ Database

# Amazon DynamoDB Accelerator (DAX)

- Fully managed, highly available in-memory cache for Amazon DynamoDB

- Response times in microseconds (instead of milliseconds)

- API compatible with DynamoDB - you can simply point your existing DynamoDB application at the DAX endpoint, no need to rewrite the application.

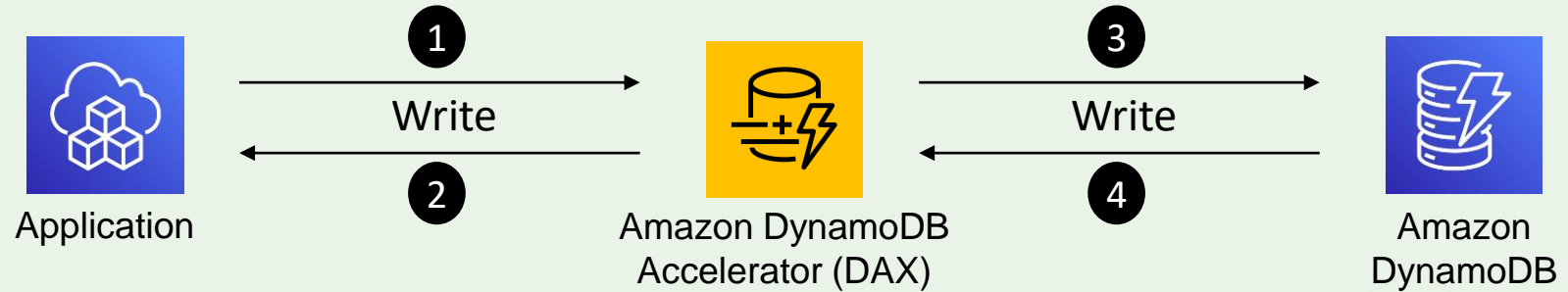- Security – Amazon VPC, AWS IAM, CloudTrail, AWS Organizations

Read-Through Cache

Application

1 Read →

4 Read ←

Amazon DynamoDB Accelerator (DAX)

2 Read →

3 Read ←

Amazon DynamoDB

# Write Operations