aws academy

**CSCI 5902 Adv. Cloud Architecting**
**Fall 2023**
**Instructor: Lu Yang**

**Module 13 Caching Content (Sections 4-5)**
**Dec 4, 2023**

# Housekeeping items and feedback

1. Start recording

# Recap of our last lecture

# Module overview

## Sections

1. Architectural need

2. Overview of caching

3. Edge caching ← We stopped here in the last lecture

4. Caching web sessions

5. Caching databases

## Lab

- Guided Lab: Streaming Dynamic Content Using Amazon CloudFront

# Section 4: Caching web sessions

aws academy

# Session management: Sticky sessions
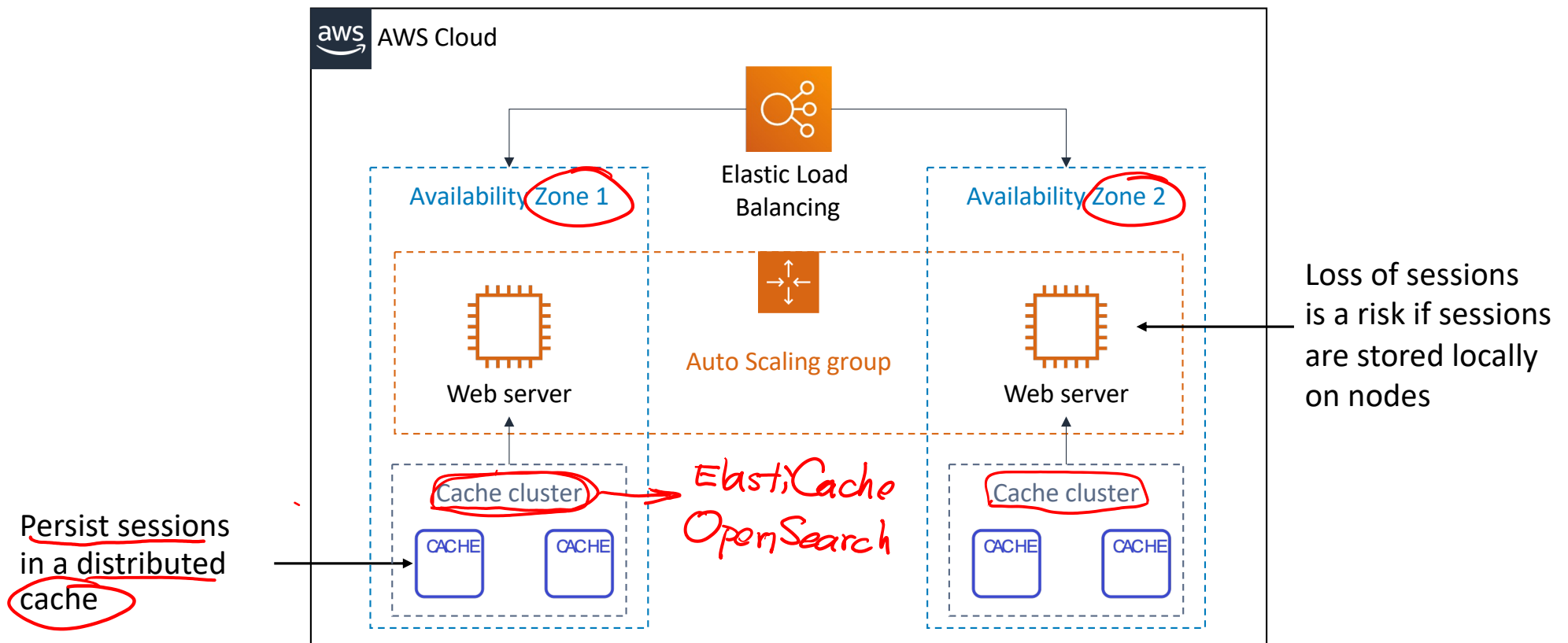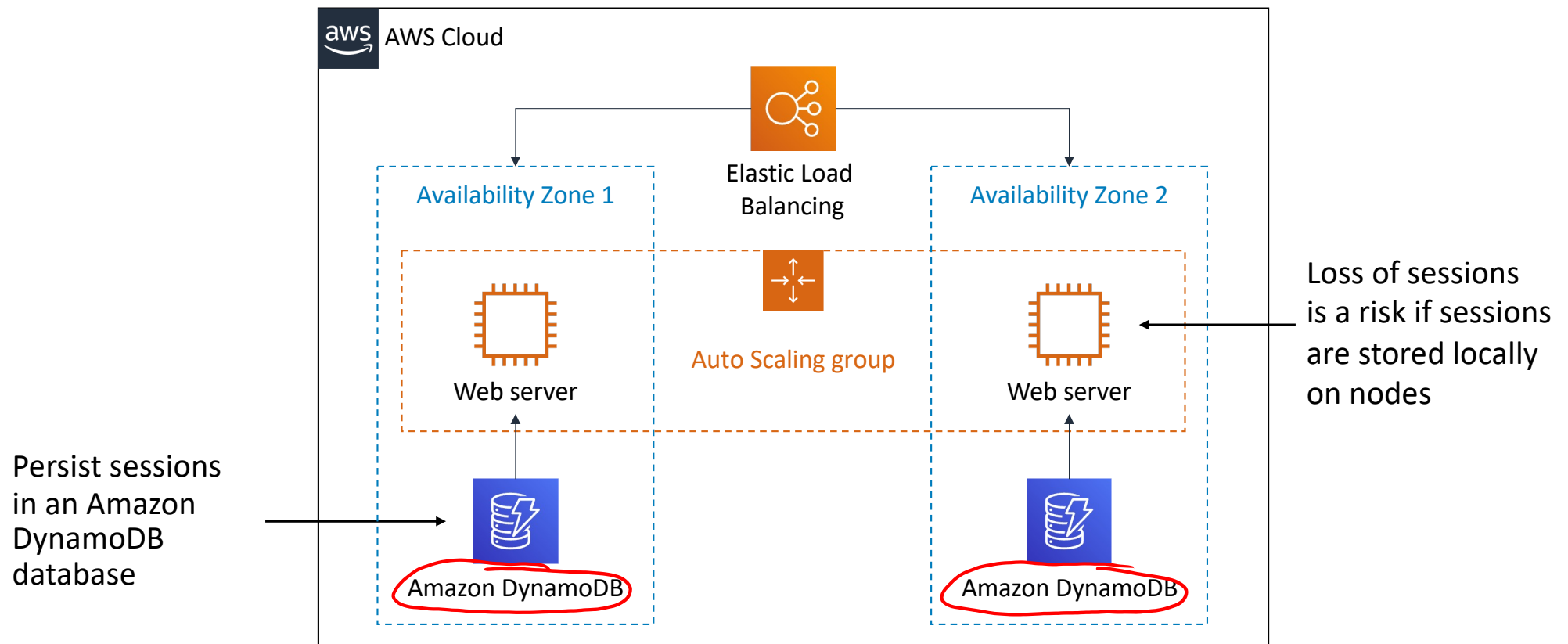
Elastic Load Balancing



Sticky sessions

Feature that enables a load balancer to route a request to the specific server that manages the user's session.

- Use client-side cookies
- Are cost-effective
- Speed up retrieval of sessions
- Have disadvantages –
  - Loss of sessions when you have an instance failure
  - Limit scalability: Uneven load distribution and increased latency

# Instead of sticky sessions:
# Persist sessions inside a distributed cache

AWS Cloud

Elastic Load Balancing

Availability Zone 1

Availability Zone 2

Auto Scaling group

Web server

Web server

Loss of sessions is a risk if sessions are stored locally on nodes
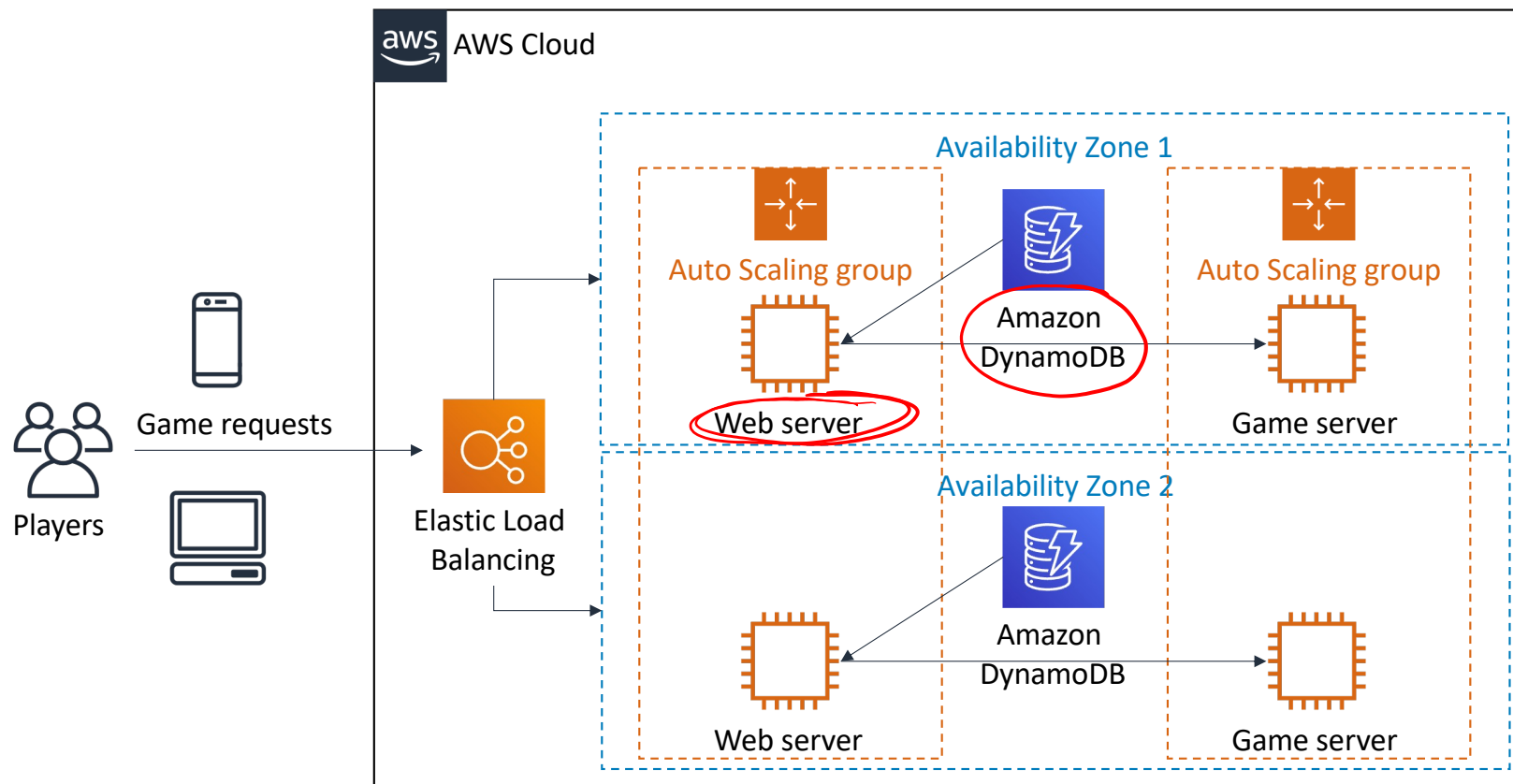
Cache cluster

Cache cluster

ElastiCache OpenSearch

CACHE    CACHE

CACHE    CACHE

Persist sessions in a distributed cache

# Instead of sticky sessions:
# Persist sessions inside a DynamoDB table



Persist sessions in an Amazon DynamoDB database

Loss of sessions is a risk if sessions are stored locally on nodes

AWS Cloud

Elastic Load Balancing

Availability Zone 1

Availability Zone 2

Auto Scaling group

Web server

Web server

Amazon DynamoDB

Amazon DynamoDB

# Example: Storing session states for an online gaming application
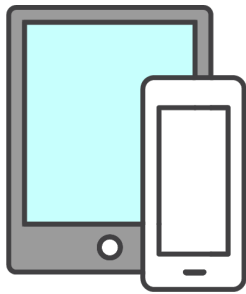
# Section 4 key takeaways



- Sessions are used to manage user authentication and store user data while the user interacts with the application.

- You can manage sessions with sticky sessions, which is a feature of Elastic Load Balancing load balancers. Sticky sessions route requests to the specific server that's managing the user's session.

- You can also manage sessions by persisting session data outside the web server instance—for example, in a distributed cache or DynamoDB table.

# Section 5: Caching databases

aws academy

# When should you cache your database?

You are concerned about response times for your customer.

*latency*

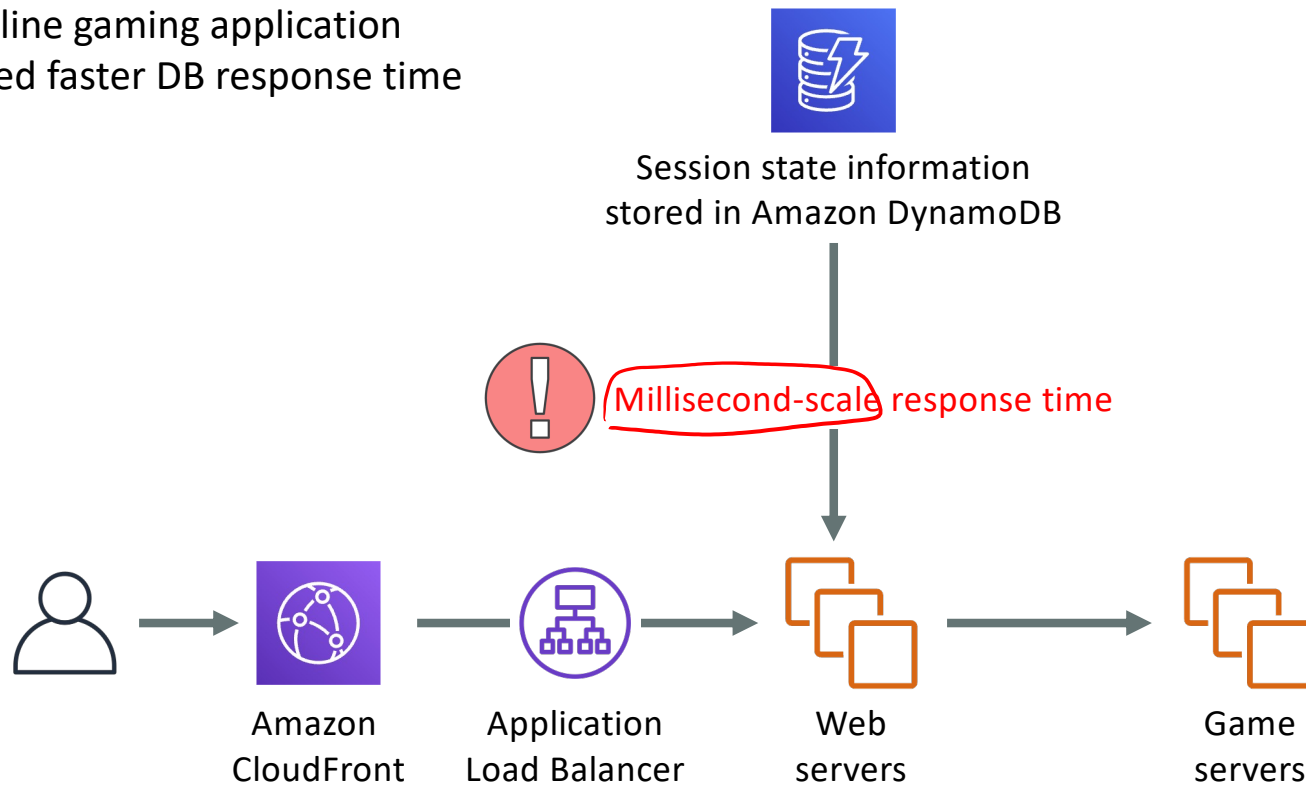You have a high volume of requests that are inundating your database.

You would like to reduce your database costs.

*read replicas*

*Cache*

# Using DynamoDB for state information

Wait, I should include the aws academy text in header area but it's a logo. Transcribe minimal.

Use case: Online gaming application
Problem: Need faster DB response time

Session state information
stored in Amazon DynamoDB

Millisecond-scale response time

Amazon
CloudFront

Application
Load Balancer

Web
servers

Game
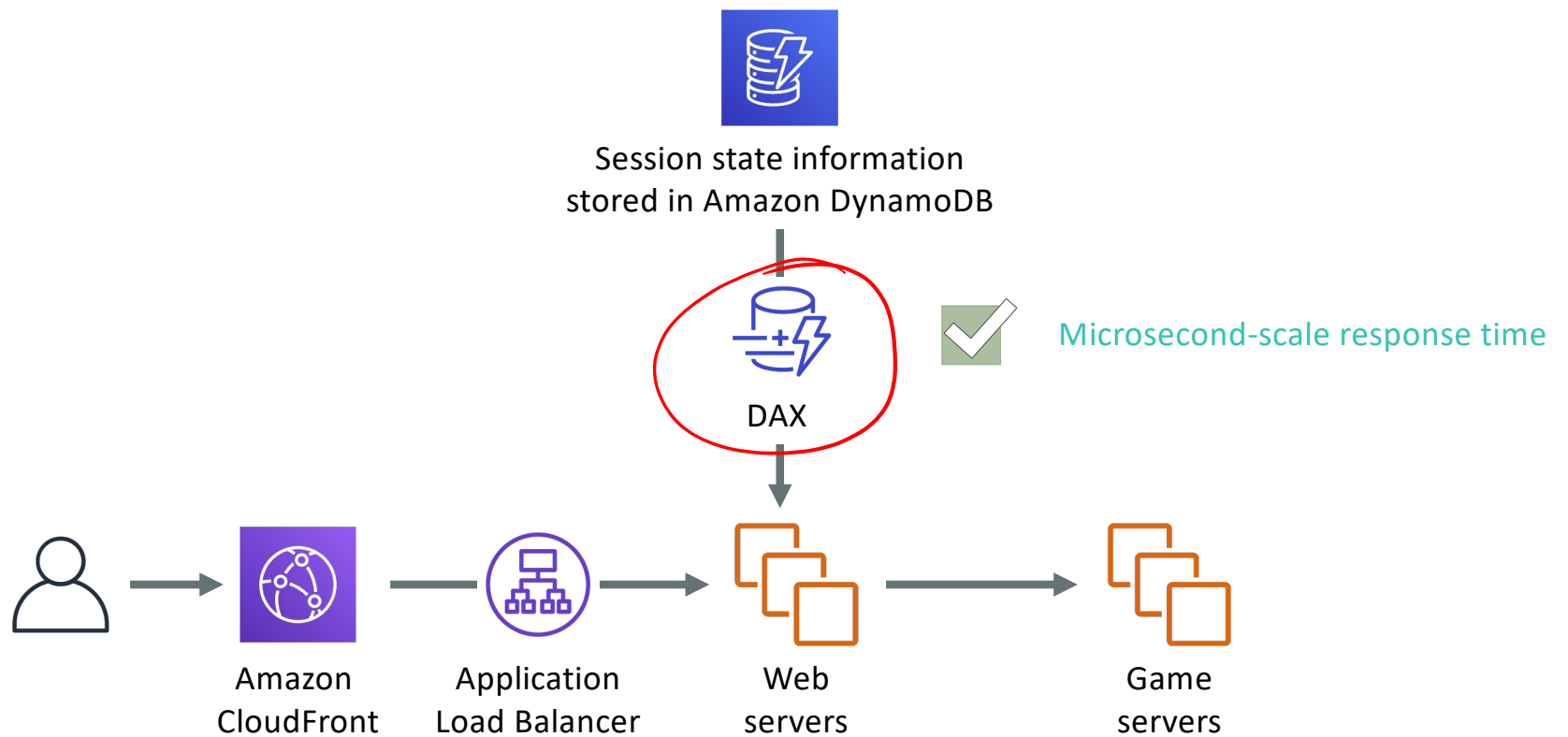servers

# Amazon DynamoDB Accelerator

DAX

Fully managed, highly available, <u>in-memory cache</u> for DynamoDB

- Extreme <u>performance</u> (<u>microsecond</u>-scale response time)
- Highly scalable
- Fully managed
- Integrated with DynamoDB
- Flexible
- <u>Secure</u>
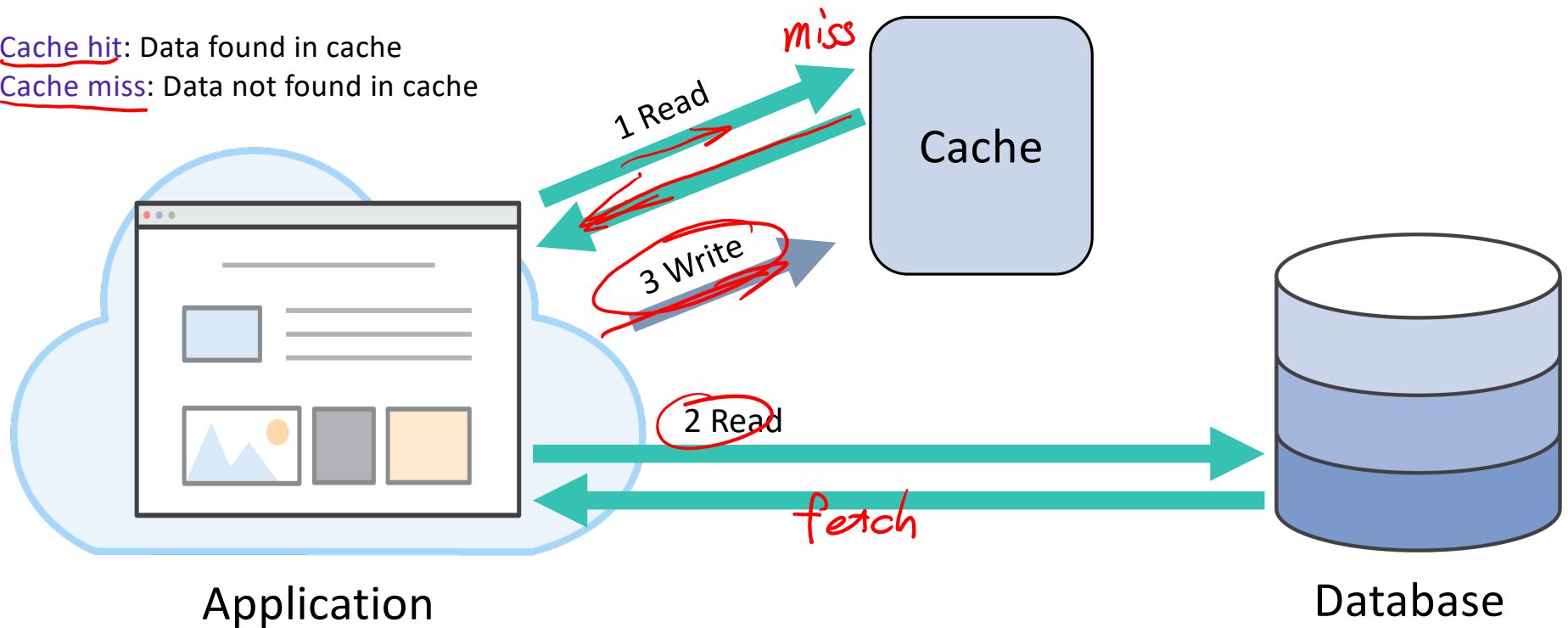
Amazon DynamoDB Accelerator

# Using DynamoDB with DAX to accelerate response time

Session state information
stored in Amazon DynamoDB

DAX

Microsecond-scale response time

Amazon
CloudFront

Application
Load Balancer

Web
servers

Game
servers

# Remote or side caches

Cache hit: Data found in cache
Cache miss: Data not found in cache

*miss*

Cache

1 Read

3 Write

2 Read

*fetch*

Application

Database

# Amazon ElastiCache
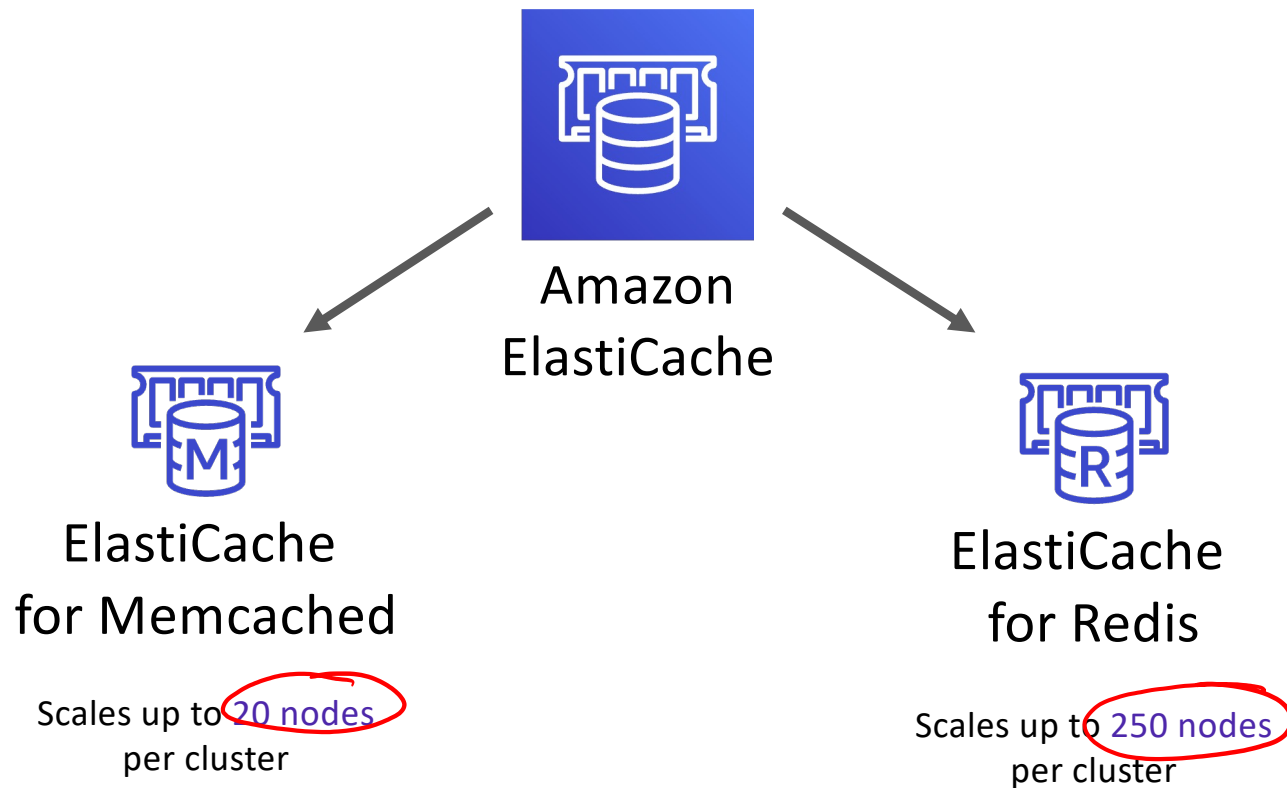
Amazon
ElastiCache

ElastiCache provides <u>web applications</u> with an <u>in-memory data store</u> in the cloud.

- Works an <u>in-memory data store</u> and cache
- Offers <u>high performance</u>  *microsecond*
- Is fully managed
- Is scalable
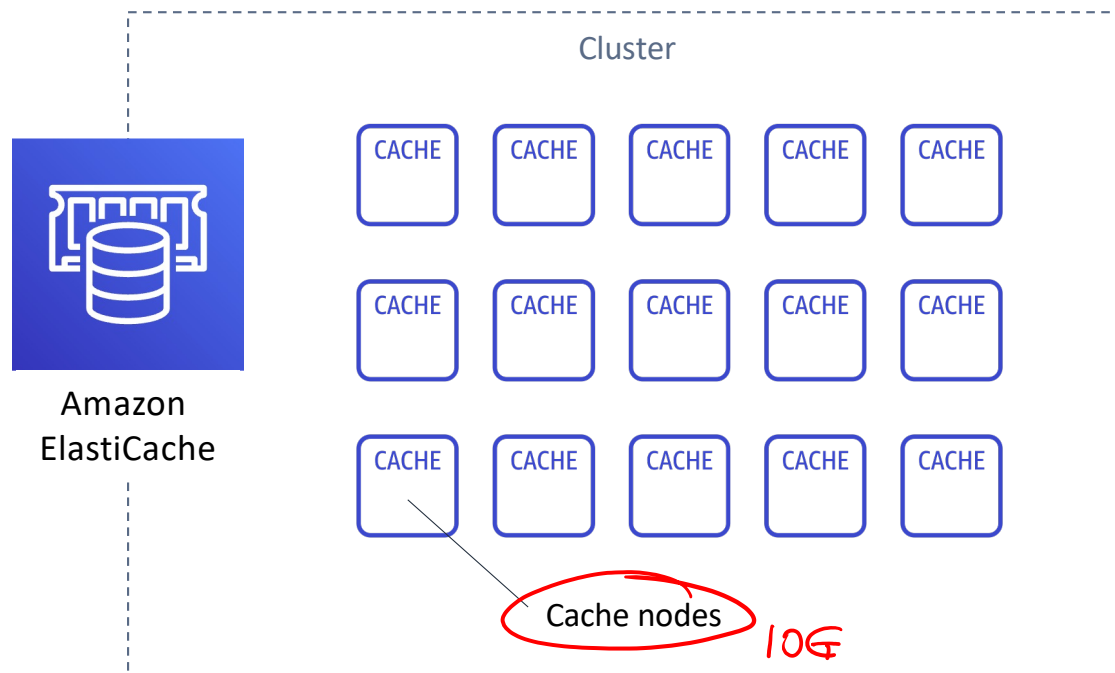- Supports <u>Redis</u> and <u>Memcached</u>

# Redis and Memcached

Amazon
ElastiCache

ElastiCache
for Memcached

Scales up to 20 nodes
per cluster

ElastiCache
for Redis

Scales up to 250 nodes
per cluster

# Memcached versus Redis comparison

| Feature | Memcached | Redis |
|---|---|---|
| Sub-millisecond latency | Yes | Yes |
| Ability to scale horizontally for writes and storage | Yes | No |
| Multi-threaded performance | Yes | No |
| Advanced data structures | No | Yes |
| Sorting and ranking datasets | No | Yes |
| Publish/subscribe messaging | No | Yes |
| Multi-AZ deployments with automatic failover | No | Yes |
| Persistence | No | Yes |

# ElastiCache components

Cluster

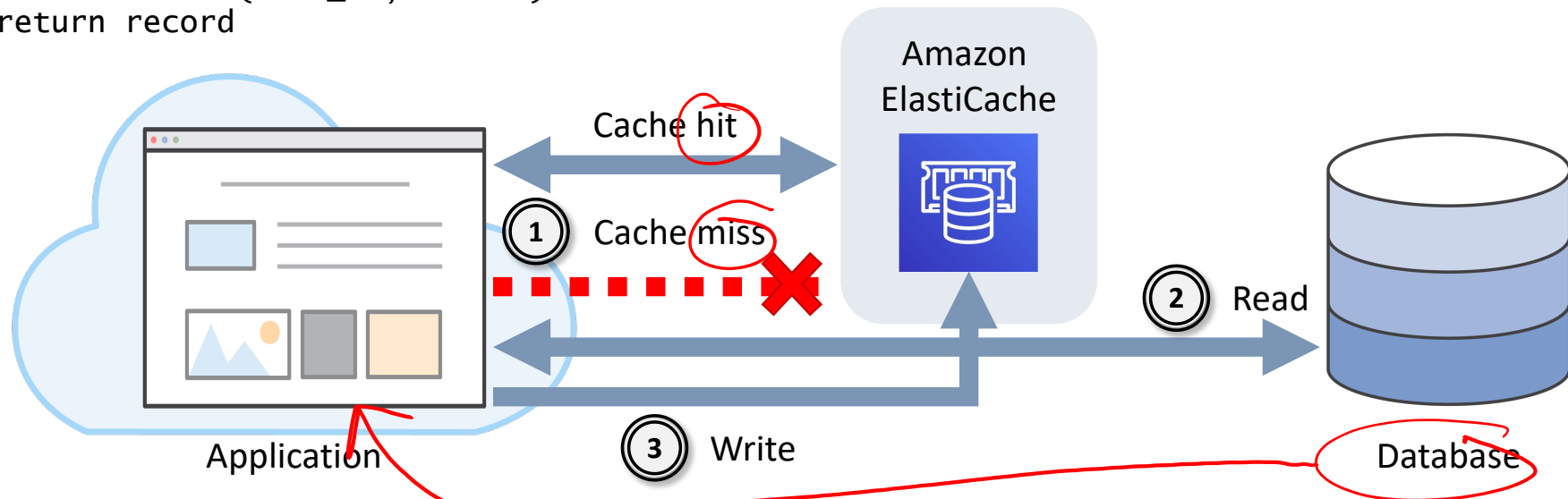| CACHE | CACHE | CACHE | CACHE | CACHE |
|-------|-------|-------|-------|-------|
| CACHE | CACHE | CACHE | CACHE | CACHE |
| CACHE | CACHE | CACHE | CACHE | CACHE |

Amazon ElastiCache

Cache nodes *10G*

- A node is the smallest block of an ElastiCache deployment

- Each node has its own DNS name and port

- A cluster is a logical grouping of one or more nodes
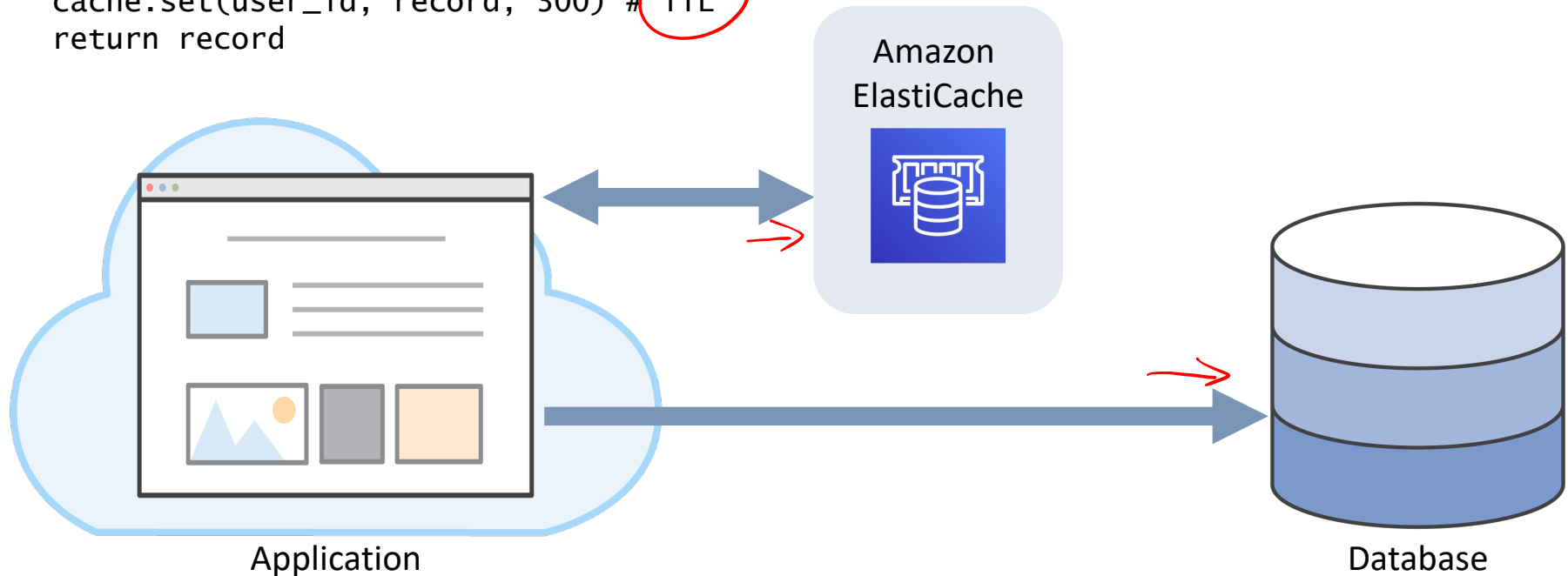
# Caching strategies: Lazy loading

```
def get_user(user_id):
    # Check the cache
    record = cache.get(user_id)
    if record is None:
        # Run a DB query
        record = db.query("select * from users where id = ?", user_id)
        # Populate the cache
        cache.set(user_id, record)
    return record
```
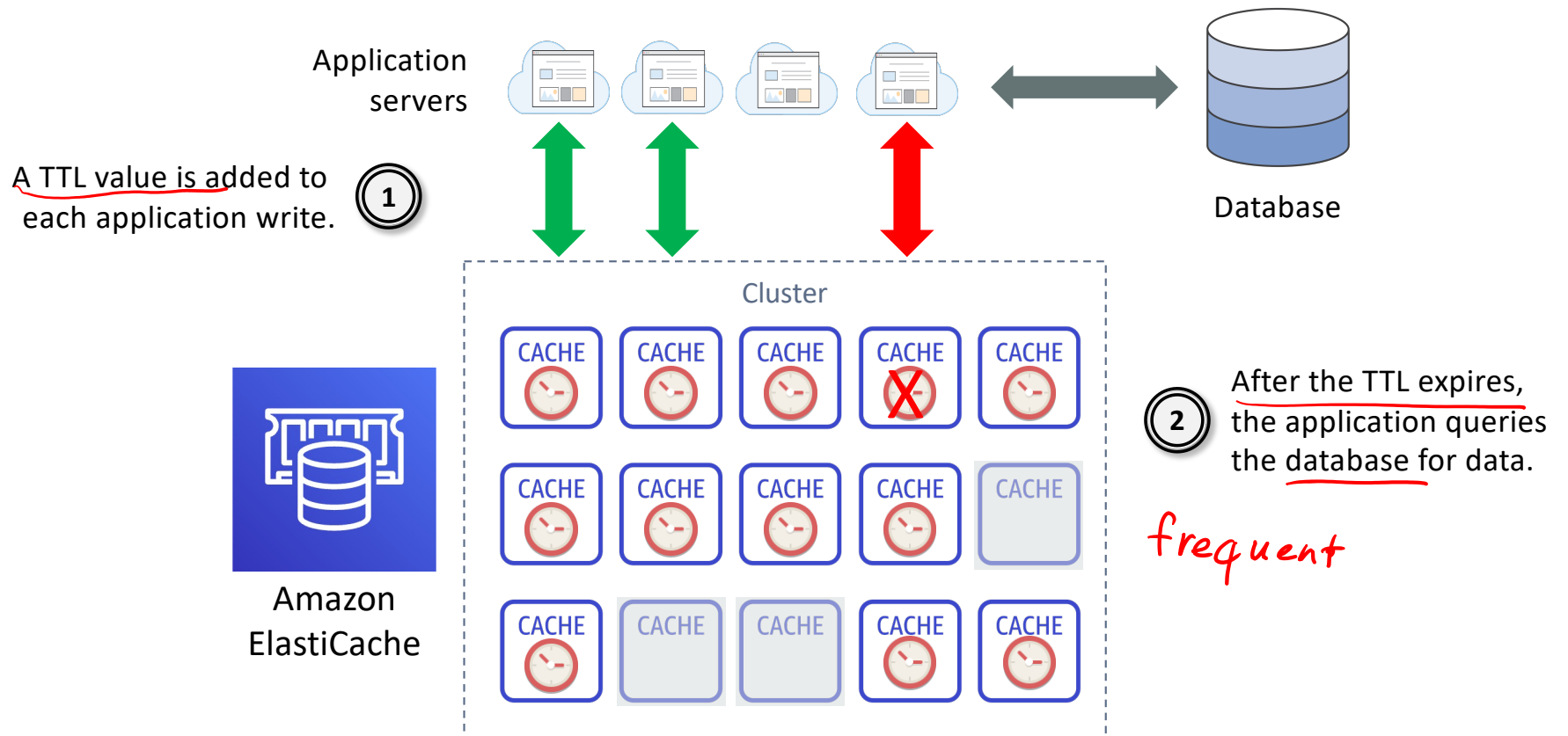
Amazon ElastiCache

Cache hit

① Cache miss

② Read

③ Write

Application

Database

# Caching strategies: Write-through
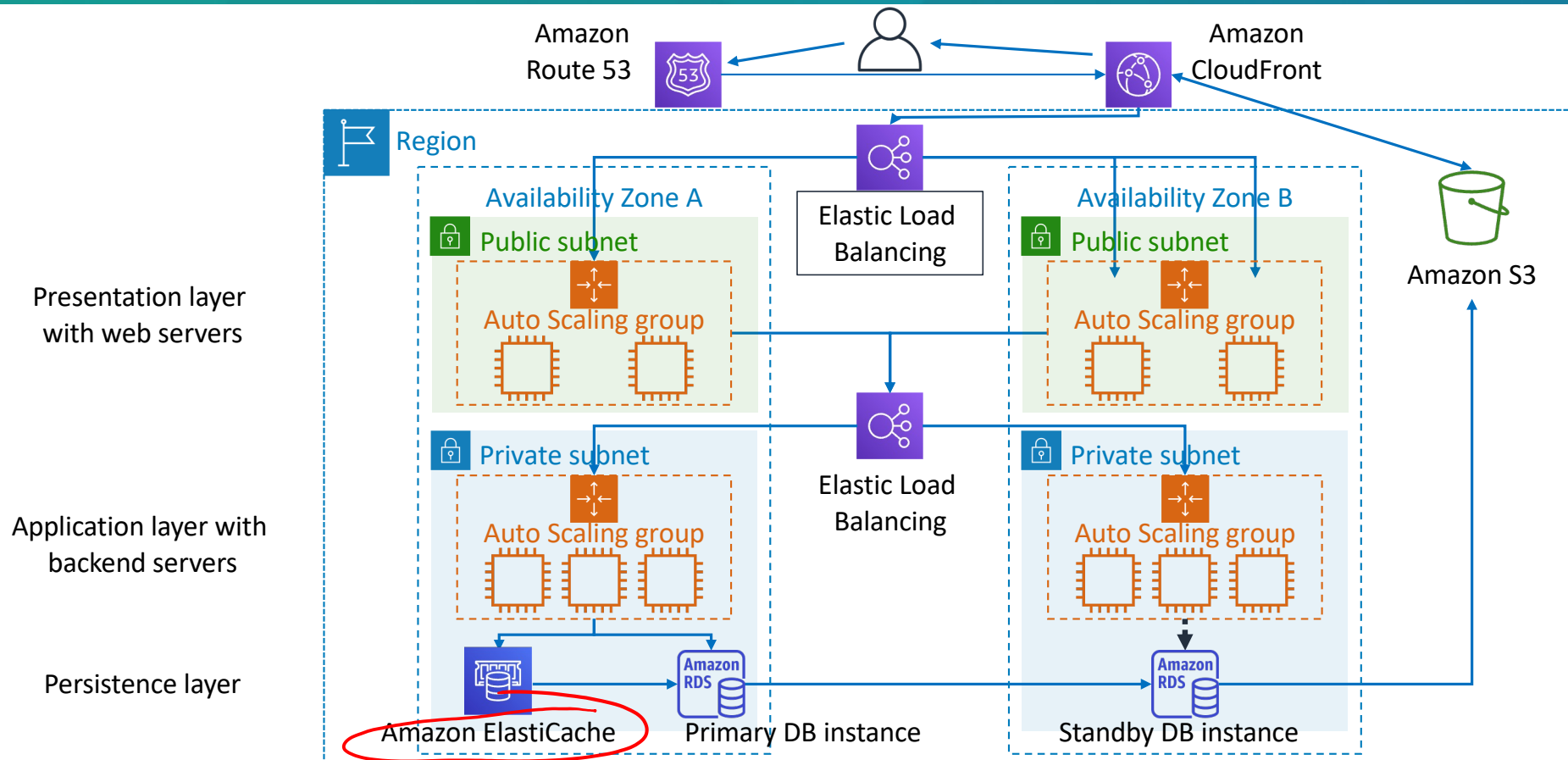
```
def save_user(user_id, values):
    # Save to DB
    record = db.query("update users…where id = ?", user_id, values)
    # Push into cache
    cache.set(user_id, record, 300) # TTL
    return record
```

Amazon
ElastiCache

Application

Database

# Adding TTL

Application servers

Database

**A TTL value is added to each application write.** ①

Amazon ElastiCache

## Cluster

| CACHE | CACHE | CACHE | CACHE | CACHE |
| CACHE | CACHE | CACHE | CACHE | CACHE |
| CACHE | CACHE | CACHE | CACHE | CACHE |

② After the TTL expires, the application queries the database for data.

frequent

# Three-tier web hosting architecture



Amazon Route 53

Amazon CloudFront

Region

**Availability Zone A**

Public subnet

Auto Scaling group

Elastic Load Balancing

**Availability Zone B**

Public subnet

Auto Scaling group

Amazon S3

Presentation layer with web servers

Private subnet

Auto Scaling group

Elastic Load Balancing

Private subnet

Auto Scaling group

Application layer with backend servers

Persistence layer

Amazon ElastiCache

Amazon RDS
Primary DB instance

Amazon RDS
Standby DB instance

# Section 5 key takeaways

- A database cache supplements your primary database by removing unnecessary pressure on it, typically in the form of frequently accessed read data

- DAX is a fully managed, highly available, in-memory cache for DynamoDB that delivers a performance improvement of up to 10 times—from milliseconds to microseconds

- Amazon ElastiCache is a side cache that works as an in-memory data store to support the most demanding applications that require microsecond response times

Module 13: Caching Content

# Module wrap-up

aws academy

# Module summary

In summary, in this module, you learned how to:

- Identify how caching content can improve application performance and reduce latency
- Create architectures that use Amazon CloudFront to cache content
- Identify how to design architectures that use edge locations for distribution and distributed denial of service (DDoS) protection
- Recognize how session management relates to caching
- Describe how to design architectures that use Amazon ElastiCache

# Thank you

aws academy