

# Introduction

Caching is a critical optimization technique in software systems designed to reduce data retrieval time and improve overall application performance. By temporarily storing frequently accessed or computationally expensive data, caching can reduce latency, backend load, and costs. However, mismanagement can lead to performance bottlenecks, stale data issues, and increased complexity.

Caching exists in various forms:

## 1. Levels of Caching:

- **Client-side Caching:** Browser or local storage (e.g., HTTP headers like ETag, Cache-Control).
- **Server-side Caching:** Application-level caching (e.g., memory, disk).
- **Distributed Caching:** Clusters of cache servers (e.g., Redis, Memcached).
- **Database Caching:** Query results caching or materialized views.
- **CDNs (Content Delivery Networks):** Edge caching for static assets.

## 2. Types of Cached Data:

- **Data Caches:** Storing frequently requested database query results.
- **Content Caches:** Static assets like CSS, JavaScript, or images.
- **Computation Caches:** Results of expensive computations (e.g., recommendation engines).

---

## Positive Impacts of Caching

### 1. Reduced Latency

Caching stores data closer to the requesting system, minimizing retrieval time. For example:

- In-memory caching via Redis/Memcached can fetch data within **sub-millisecond latency**, compared to **tens or hundreds of milliseconds** for database queries.

**Scenario:** An e-commerce platform caches product details to serve user requests rapidly without hitting the database.

### Configuration Example (Redis):

- Install Redis:

```
sudo apt-get update
sudo apt-get install redis-server
```

- Configure Redis (/etc/redis/redis.conf):

```
maxmemory 512mb
maxmemory-policy allkeys-lru
```

### Code Example (Node.js with Redis):

```
const Redis = require('ioredis');
const redis = new Redis();

async function getProductDetails(productId) {
  const cacheKey = `product:${productId}`;
  const cachedData = await redis.get(cacheKey);

  if (cachedData) {
    console.log('Cache Hit');
    return JSON.parse(cachedData);
  }

  console.log('Cache Miss');
  const productDetails = await fetchFromDatabase(productId); // Simulated DB call
  await redis.set(cacheKey, JSON.stringify(productDetails), 'EX', 3600); // 1-hour TTL

  return productDetails;
}
```

---

## 2. Decreased Backend Load

By offloading repetitive queries and requests from the backend, caching reduces database or API stress.

**Scenario:** A high-traffic website that repeatedly queries user profile details from the database caches this data in Redis. Only changes (updates) trigger database access.

### Performance Metrics:

- Reduced Database Queries: 10,000/sec → 2,000/sec
- Cache Hit Ratio: >80%

**Configuration (MySQL Query Cache):** Enable query\_cache\_size in MySQL:

```
SET GLOBAL query_cache_size = 1048576; -- 1 MB
SET GLOBAL query_cache_type = ON;
```

---

### 3. Improved Scalability

Caching enhances horizontal scaling by decoupling data retrieval from backend systems.

**Scenario:** A video streaming platform uses a multi-level cache:

1. **Edge Cache (CDN):** Stores static content.
2. **Application Cache (Redis):** Caches metadata (e.g., video recommendations).
3. **In-Memory Cache (Local):** Stores transient session data.

#### Edge Caching with NGINX Configuration:

```
proxy_cache_path /data/nginx/cache levels=1:2 keys_zone=STATIC:10m max_size=10g;
server{
    location /videos/{
        proxy_cache STATIC;
        proxy_cache_key $uri;
        proxy_pass http://backend;
        add_header X-Cache-Status $upstream_cache_status;
    }
}
```

---

### 4. Enhanced User Experience

Faster response times lead to better user engagement. CDNs like Cloudflare or Akamai can cache static assets globally, reducing latency.

**Scenario:** A news website serves static assets (CSS/JS) via CDN with a cache lifetime of 24 hours.

#### Configuration (HTTP Headers):

```
Cache-Control: public, max-age=86400
ETag: "5d8c72a6-3e2a"
```

---

### Case Study: Netflix

Netflix optimizes its performance with a three-tier caching strategy:

1. **Client-Side Caching:** Smart caching within apps.
  2. **Edge Caching (CDNs):** Reduces streaming latency.
  3. **Metadata Caching (EVCache):** Uses Redis-based architecture for session data.
-

## Negative Impacts of Caching

### 1. Cache Invalidation Complexity

Keeping the cache consistent with the source of truth is a well-known challenge.

**Scenario:** An inventory system updates product stock frequently. Cache invalidation delays can lead to discrepancies, showing outdated stock levels.

**Technical Solution:** Use event-driven invalidation with tools like Kafka.

#### Configuration Example (Redis Keyspace Notifications):

Enable notifications in redis.conf:

```
notify-keyspace-events Ex
```

Consume invalidation events:

```
import redis

def handle_event(event):
    print(f"Cache invalidation triggered for: {event['key']}")

r = redis.StrictRedis()
pubsub = r.pubsub()
pubsub.subscribe('__keyevent@0__:expired')

for message in pubsub.listen():
    handle_event(message)
```

---

### 2. Stale Data

If cached data is not refreshed in time, users may see outdated information.

**Scenario:** A stock trading platform shows old stock prices due to delayed cache refresh.

**Mitigation:** Use a combination of:

- **Short TTLs**
- **Write-through cache:** Immediately update cache during writes.

## Example: Spring Boot + Redis

```
@Cacheable(value = "stockPrices", key = "#symbol")
public StockPrice getStockPrice(String symbol) {
    return stockRepository.findPriceBySymbol(symbol);
}

@CacheEvict(value = "stockPrices", key = "#symbol")
public void updateStockPrice(String symbol, double price) {
    stockRepository.updatePrice(symbol, price);
}
```

---

### 3. Cache Misses and Thundering Herd

When many requests miss the cache, backend systems can get overwhelmed.

**Scenario:** A cache server crash leads to all requests hitting the database, causing a bottleneck.

#### Mitigation:

- **Cache Prewarming:** Load frequently accessed data during startup.
- **Request Coalescing:** Deduplicate simultaneous requests for the same data.

#### Example (Redis Lua Script for Request Coalescing):

```
-- Redis Lua Script: Return cached value or set a lock
local key = KEYS[1]
local lockKey = key .. ":lock"

if redis.call("EXISTS", lockKey) == 1 then
    return "LOCKED"
end

local value = redis.call("GET", key)
if value then
    return value
else
    redis.call("SET", lockKey, 1, "EX", 10)
    return nil
end
```

---

## 4. Increased Memory Consumption

Caches consume memory, and inefficient eviction policies can cause performance degradation.

**Scenario:** A cache using the default FIFO eviction policy stores low-priority data longer than critical data.

**Mitigation:** Use eviction policies like **LRU** or **LFU** based on access patterns.

### Redis Configuration:

```
maxmemory 512mb  
maxmemory-policy allkeys-lru
```

---

## Case Study: Twitter

**Problem:** Twitter faced issues with stale tweets due to delayed cache invalidation, especially during high-traffic events.

**Solution:** Implemented Kafka-based event-driven cache updates and tuned TTLs.

---

## Performance Engineering Metrics

1. **Latency:** Measure the time to serve requests with and without caching.
  2. **Hit/Miss Ratios:** High hit ratios indicate effective caching.
  3. **Evictions:** Monitor for unnecessary evictions due to memory constraints.
  4. **Resource Usage:** Monitor memory/CPU consumption of cache servers.
- 

## Conclusion

Caching is a powerful technique for improving performance, scalability, and user experience. However, it must be carefully designed and monitored. Performance engineers must balance caching benefits with potential trade-offs like consistency challenges, stale data, and resource consumption. By using robust configurations, monitoring tools, and advanced strategies like request coalescing and cache invalidation frameworks, applications can leverage caching effectively while minimizing its drawbacks.