



## (/explore)

엏 Tracks

(/tracks)

My Courses (/mycourses)



Edpresso (/edpresso)



Refer a Friend (/refer-afriend)

## **Grokking the System Design Interview**

(/collection/5668639101419520/56490502 G 87% completed

Search Course

### Glossary of System **Design Basics**

System Design Basics (/courses/grokking-the-systemdesign-interview/B892KY261z2)

Key Characteristics of Distributed Systems

(/courses/grokking-the-systemdesign-interview/YQWGjlZZVz9)

Load Balancing (/courses/grokking-the-systemdesign-interview/3jEwI04BL7Q)

Caching (/courses/grokking-the-systemdesign-interview/3j6NnJrpp5p)

Data Partitioning (/courses/grokking-the-systemdesign-interview/mEN8IJXV1LA)

Indexes (/courses/grokking-the-systemdesign-interview/gxkVE8NEvXj)

Proxies (/courses/grokking-the-system-

interview/N8G9MvM4OR2)

design-

Redundancy and Replication (/courses/grokking-the-systemdesign-interview/xV1qvj6PKkJ)

SQL vs. NoSQL (/courses/grokking-the-systemdesign-interview/YQIK1mDPgpK)

**CAP Theorem** (/courses/grokking-the-system-

design-interview/RMkqx1Egxqz)

Consistent Hashing (/courses/grokking-the-systemdesign-interview/B81vnyp0GpY)

Long-Polling vs WebSockets vs Server-Sent Events (/courses/grokking-the-systemdesign-interview/gx7wZzWn5Vj)

#### ırses/grokking-**Appendix**

Contact Us (/courses/grokking-the-systemdesign-interview/7nnxwPxAl 7il

## Caching

#### We'll cover the following

- Application server cache
- Content Distribution Network (CDN)
- Cache Invalidation
- Cache eviction policies

Load balancing helps you scale horizontally across an everincreasing number of servers, but caching will enable you to make vastly better use of the resources you already have as well as making otherwise unattainable product requirements feasible. Caches take advantage of the locality of reference principle: recently requested data is likely to be requested again. They are used in almost every layer of computing: hardware, operating systems, web browsers, web applications, and more. A cache is like short-term memory: it has a limited amount of space, but is typically faster than the original data source and contains the most recently accessed items. Caches can exist at all levels in architecture, but are often found at the level nearest to the front end where they are implemented to return data quickly without taxing downstream levels.

## Application server cache

#

Placing a cache directly on a request layer node enables the local storage of response data. Each time a request is made to the service, the node will quickly return local cached data if it exists. If it is not in the cache, the requesting node will query the data from disk. The cache on one request layer node could also be located both in memory (which is very fast) and on the node's local disk (faster than going to network storage).

What happens when you expand this to many nodes? If the request layer is expanded to multiple nodes, it's still quite possible to have each node host its own cache. However, if your load balancer randomly distributes requests across the nodes, the same request will go to different nodes, thus increasing cache misses. Two choices for overcoming this hurdle are global caches and distributed caches.

### Content Distribution Network (CDN)

#



Create

system-design view/3j6NnJrp

# 엏

Tracks (/tracks)

My Courses (/mycourses)



Edpresso (/edpresso)



Refer a Friend (/refer-afriend)

## **Grokking the System Design Interview**

(/collection/5668639101419520/56490502 87% completed G

Search Course

### Glossary of System **Design Basics**

System Design Basics (/courses/grokking-the-systemdesign-interview/B892KY261z2)

Key Characteristics of Distributed Systems

(/courses/grokking-the-systemdesign-interview/YQWGjlZZVz9)

Load Balancing (/courses/grokking-the-systemdesign-interview/3jEwI04BL7Q)

Caching (/courses/grokking-the-systemdesign-interview/3j6NnJrpp5p)

Data Partitioning (/courses/grokking-the-systemdesign-interview/mEN8IJXV1LA)

Indexes (/courses/grokking-the-systemdesign-interview/gxkVE8NEvXj)

Proxies (/courses/grokking-the-systemdesign-

interview/N8G9MvM4OR2)

Redundancy and Replication (/courses/grokking-the-systemdesign-interview/xV1qvj6PKkJ)

SQL vs. NoSQL (/courses/grokking-the-systemdesign-interview/YQIK1mDPgpK)

**CAP Theorem** (/courses/grokking-the-systemdesign-interview/RMkqx1Egxqz)

Consistent Hashing (/courses/grokking-the-systemdesign-interview/B81vnyp0GpY)

Long-Polling vs WebSockets vs Server-Sent Events (/courses/grokking-the-systemdesign-interview/gx7wZzWn5Vj)

## **Appendix**

Contact Us design-interview/7nnxwPxAl 7il

CDNs are a kind of cache that comes into play for sites serving large amounts of static media. In a typical CDN setup, a request will first ask the CDN for a piece of static media; the CDN will serve that content if it has it locally available. If it isn't available, the CDN will query the back-end servers for the file, cache it locally, and serve it to the requesting user.

If the system we are building isn't yet large enough to have its own CDN, we can ease a future transition by serving the static media off a separate subdomain (e.g. static.yourservice.com (http://static.yourservice.com)) using a lightweight HTTP server like Nginx, and cut-over the DNS from your servers to a CDN later.

#### Cache Invalidation

#

While caching is fantastic, it does require some maintenance for keeping cache coherent with the source of truth (e.g., database). If the data is modified in the database, it should be invalidated in the cache; if not, this can cause inconsistent application behavior.

Solving this problem is known as cache invalidation; there are three main schemes that are used:

**Write-through cache:** Under this scheme, data is written into the cache and the corresponding database at the same time. The cached data allows for fast retrieval and, since the same data gets written in the permanent storage, we will have complete data consistency between the cache and the storage. Also, this scheme ensures that nothing will get lost in case of a crash, power failure, or other system disruptions.

Although, write through minimizes the risk of data loss, since every write operation must be done twice before returning success to the client, this scheme has the disadvantage of higher latency for write operations.

**Write-around cache:** This technique is similar to write through cache, but data is written directly to permanent storage, bypassing the cache. This can reduce the cache being flooded with write operations that will not subsequently be re-read, but has the disadvantage that a read request for recently written data will create a "cache miss" and must be read from slower back-end storage and experience higher latency.

Write-back cache: Under this scheme, data is written to cache alone and completion is immediately confirmed to the client. The write to the permanent storage is done after specified intervals or under certain conditions. This results in low latency and high throughput for write-intensive applications, however, this speed comes with the risk of data loss in case of a crash or other adverse

Create ırses/grokkingsystem-design

view/3j6NnJrp

(/courses/grokking-the-system-



# ₩

Explore (/explore)

### 엏 Tracks

(/tracks)





Edpresso (/edpresso)



Refer a Friend (/refer-afriend)

## **Grokking the System Design Interview**

(/collection/5668639101419520/56490502

87% completed

G

Search Course

### Glossary of System **Design Basics**

System Design Basics (/courses/grokking-the-systemdesign-interview/B892KY261z2)

Key Characteristics of Distributed Systems

(/courses/grokking-the-systemdesign-interview/YQWGjlZZVz9)

Load Balancing (/courses/grokking-the-systemdesign-interview/3jEwl04BL7Q)

Caching (/courses/grokking-the-system-

design-interview/3j6NnJrpp5p)

Data Partitioning (/courses/grokking-the-system-

design-interview/mEN8IJXV1LA)

Indexes (/courses/grokking-the-system-

design-interview/gxkVE8NEvXj)

Proxies (/courses/grokking-the-system-

interview/N8G9MvM4OR2)

design-

Redundancy and Replication (/courses/grokking-the-systemdesign-interview/xV1qvj6PKkJ)

SQL vs. NoSQL (/courses/grokking-the-systemdesign-interview/YQIK1mDPgpK)

**CAP Theorem** (/courses/grokking-the-systemdesign-interview/RMkqx1Egxqz)

Consistent Hashing (/courses/grokking-the-systemdesign-interview/B81vnyp0GpY)

Long-Polling vs WebSockets vs Server-Sent Events (/courses/grokking-the-systemdesign-interview/gx7wZzWn5Vi)

#### ırses/grokking-**Appendix** system-design

Create

view/3j6NnJrp

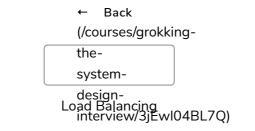
Contact Us (/courses/grokking-the-systemdesign-interview/7nnxwPxAl 7il

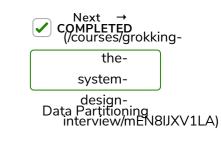
## Cache eviction policies

Following are some of the most common cache eviction policies:

- 1. First In First Out (FIFO): The cache evicts the first block accessed first without any regard to how often or how many times it was accessed before.
- 2. Last In First Out (LIFO): The cache evicts the block accessed most recently first without any regard to how often or how many times it was accessed before.
- 3. Least Recently Used (LRU): Discards the least recently used items first.
- 4. Most Recently Used (MRU): Discards, in contrast to LRU, the most recently used items first.
- 5. Least Frequently Used (LFU): Counts how often an item is needed. Those that are used least often are discarded first.
- 6. Random Replacement (RR): Randomly selects a candidate item and discards it to make space when necessary.

Following links have some good discussion about caching: [1] Cache (https://en.wikipedia.org/wiki/Cache\_(computing)) [2] Introduction to architecting systems (https://lethain.com/introduction-to-architecting-systems-for-scale/)







Issue

? Ask a Question

(https://discuss.educative.io/c/grokking-the-system-designinterview-design-gurus/glossary-of-system-design-basics-caching)