# COS 316 Precept #6:
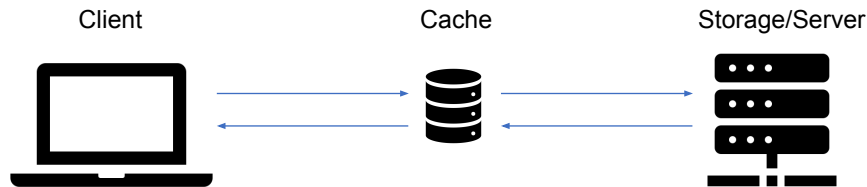## *Caching*
## *+*
## *Eviction (Replacement)*

Please visit ed discussion to download the
precept 6 exercise code ([here](#))

# What is caching?

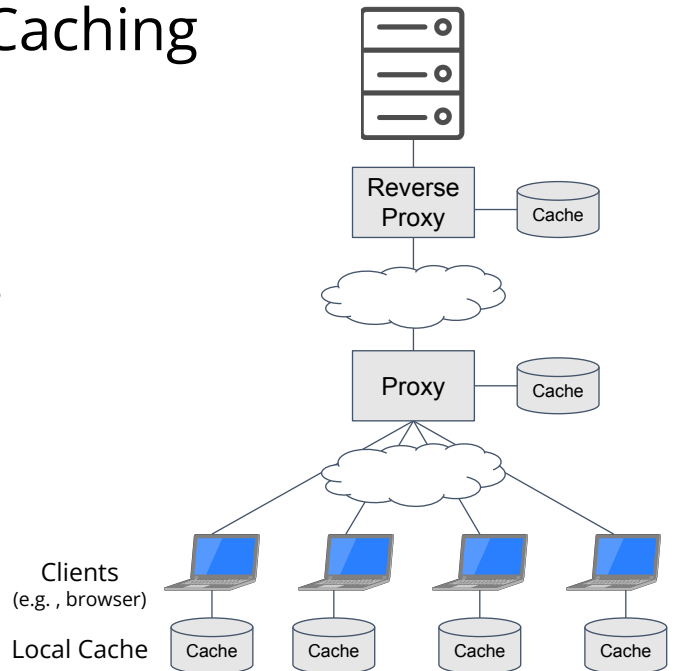Client                    Cache                 Storage/Server

- This week, you've learned a bit about web caching in lecture
- Let's talk about what a cache is and why'd you use one
- As a general computing term, a cache is a place where some objects from a larger dataset are stored for more convenient access
- You may have encountered a few caches in your computer science education
- There are caches in the CPU memory hierarchy that store a subset of main memory's contents
- The TLB, which you learned about earlier, is a cache that stores virtual to physical address translations
- And the list goes on
- So why would we want to store objects in two places?
- Really there are two reasons:
    - Accessing a cache is often faster than accessing the main datastore, because caches tend to be closer to the entity that's requesting the object
    - They also result in less traffic to the storage device or server that holds the entire dataset, which lowers the bandwidth and computing requirements of the storage server
- Caches generally take advantage of something called *locality*
- In this context, it refers to the tendency of data accesses to be correlated in some way
- Temporal locality is a type of locality that describes the tendency of things to be reused
- For example, if I access a particular object, I'm likely to use it again in the near

- future
- I've described this in a very vague and general sense so far. Let's dive a bit deeper with an example from lecture

# Overview of Web Caching

- Basic idea:
  - Bring objects "closer" to clients
- Three primary features:
  - Reduce network bandwidth
  - Reduce client-perceived delays
  - Reduce load on server
- Cache Replacement Strategy
  - When a cache becomes full, which object should be **evicted/replaced**?

Reverse Proxy — Cache

Proxy — Cache

Clients (e.g. , browser)

Local Cache — Cache   Cache   Cache   Cache

- Web caching is a specific instance of caching where objects are requested over the web and stored at one or more caches between the client and the server
- The cache may sit somewhere in the network between the client and the server
- Or the cache may be on the client itself, meaning the client has decided to store objects from previous requests
- This is done for several reasons:
  - First, it reduces network bandwidth usage
    - The network is a shared resource and if less data has to go over many of its connections, it can reduce congestion
  - Second, it reduces delays
    - Remember that caches are closer to clients than the servers are. This means that when a client needs to request an object, it will receive it faster from a cache
  - Third, it reduces load on the server
    - Simply put: If requests are handled by a cache, that means they don't need to be handled by a server.
- But caches require management
- Caches are necessarily limited in size for various reasons, principal among them being that storage is expensive
- What happens when we want to insert an object into a cache that's already full?
- We have to evict something!

- It turns out that what we choose to evict from a cache is important and there's an entire field of research based around this question

# Cache Eviction Algorithms

1. Client requests a new object
2. If object in cache
    1. return the object
3. Else:
    1. Get object from server/provider and return the object to client
    2. Attempt to insert the object into the cache:
        1. If Cache full:
            1. Identify an object in cache to evict
            2. Evict the object in the cache
            3. Replace with new object (insert new object)
    3. Admit the new object to the cache

# Cache Eviction Algorithms

- Least recently used (LRU): Evict the object from the cache whose last request is the oldest

- First-in, First-out (FIFO): Evict the object from the cache that has been in the cache the longest

- Many others...

---

- Here are a few important cache eviction algorithms that are commonly used and taught
- The first is least recently used, or LRU for short.
- LRU caches keep track of the last time an object was accessed.
- When an object needs to be evicted from the cache, it chooses the object whose most recent access is the furthest in the past
- LRU is a very intuitive algorithm; Just think about your closet.
- If you need to free up space in your closet, are you going to donate the shirt that you wore yesterday or the one that you can't even remember wearing?
- This kind of cache exploits temporal locality: Objects that were used recently are more likely to be used again
- Another algorithm that gets some use is FIFO. With this algorithm, you just evict the object that has been in the cache for the longest time
- There are many caching algorithms out there, each meant to maximize hit rate, which is the fraction of object requests that the cache can handle
- A high hit rate means that a cache is able to handle a large fraction of the requests
- A low hit rate means that most of the requests end up being handled by the server

# LRU

```
id:       8          ┌──────────┐
size:     10         │ Current  │
request: __:__       │  time:   │
admit:   __:__       │  16:00   │
                     └──────────┘
```

```
id:       6          id:       3
size:     2          size:     10
request: 13:00       request: 13:45
admit:   11:00       admit:   13:45

id:       1          id:       4
size:     3          size:     5
request: 15:01       request: 11:53
admit:   12:01       admit:   11:33

id:       11         id:       7
size:     8          size:     17
request: 11:30       request: 13:30
admit:   11:30       admit:   13:30
```

```
   Cache capacity = 50
     Cache size =  45
```

- Here we have an example of the LRU eviction algorithm in action
- This is a cache with capacity 50 units
- Each of the smaller rectangles is an object
- Notice the metadata associated with each object
- They each have an ID (just so we know what object we're referring to), size, most recent request time, and admission time
- In this example, the cache currently contains 45 units of objects when an object of size 10 comes in
- This would push the cache over its capacity, so something has to evicted
- Since this is an example of LRU, we have to look at the object that's currently in the cache with the oldest request time
- That is the object with id 11
- So we evict it!
- We then have enough space to place object 8

# LRU

| | |
|---|---|
| id: 8<br>size: 10<br>request: __:__<br>admit: __:__ | Current time: 16:00 |

**Cache capacity = 50**
**Cache size = 45**

| id: 6<br>size: 2<br>request: 13:00<br>admit: 11:00 | id: 3<br>size: 10<br>request: 13:45<br>admit: 13:45 |
|---|---|
| id: 1<br>size: 3<br>request: 15:01<br>admit: 12:01 | id: 4<br>size: 5<br>request: 11:53<br>admit: 11:33 |
| id: 11<br>size: 8<br>request: 11:30<br>admit: 11:30 | id: 7<br>size: 17<br>request: 13:30<br>admit: 13:30 |

| | |
|---|---|
| id: 8<br>size: 10<br>request: __:__<br>admit: __:__ | Current time: 16:00 |

**Cache capacity = 50**
**Cache size = 45**

| id: 6<br>size: 2<br>request: 13:00<br>admit: 11:00 | id: 3<br>size: 10<br>request: 13:45<br>admit: 13:45 |
|---|---|
| id: 1<br>size: 3<br>request: 15:01<br>admit: 12:01 | id: 4<br>size: 5<br>request: 11:53<br>admit: 11:33 |
| id: 11<br>size: 8<br>request: 11:30<br>admit: 11:30 | id: 7<br>size: 17<br>request: 13:30<br>admit: 13:30 |

**Cache capacity = 50**
**Cache size = 47**

| id: 6<br>size: 2<br>request: 13:00<br>admit: 11:00 | id: 3<br>size: 10<br>request: 13:45<br>admit: 13:45 |
|---|---|
| id: 1<br>size: 3<br>request: 15:01<br>admit: 12:01 | id: 4<br>size: 5<br>request: 11:53<br>admit: 11:33 |
| id: 8<br>size: 10<br>request: 16:00<br>admit: 16:00 | id: 7<br>size: 17<br>request: 13:30<br>admit: 13:30 |

- Here we have an example of the LRU eviction algorithm in action
- This is a cache with capacity 50 units
- Each of the smaller rectangles is an object
- Notice the metadata associated with each object
- They each have an ID (just so we know what object we're referring to), size, most recent request time, and admission time
- In this example, the cache currently contains 45 units of objects when an object of size 10 comes in
- This would push the cache over its capacity, so something has to evicted
- Since this is an example of LRU, we have to look at the object that's currently in the cache with the oldest request time
- That is the object with id 11
- So we evict it!
- We then have enough space to place object 8

# FIFO



```
id:       8
size:     10
request: __:__
admit:    __:__
```

```
Current
time:
16:00
```

```
id:       1
size:     3
request: 15:01
admit:    12:01
```

```
id:       3
size:     10
request: 13:45
admit:    13:45
```

```
id:       6
size:     8
request: 13:00
admit:    11:00
```

```
id:       4
size:     5
request: 11:53
admit:    11:33
```
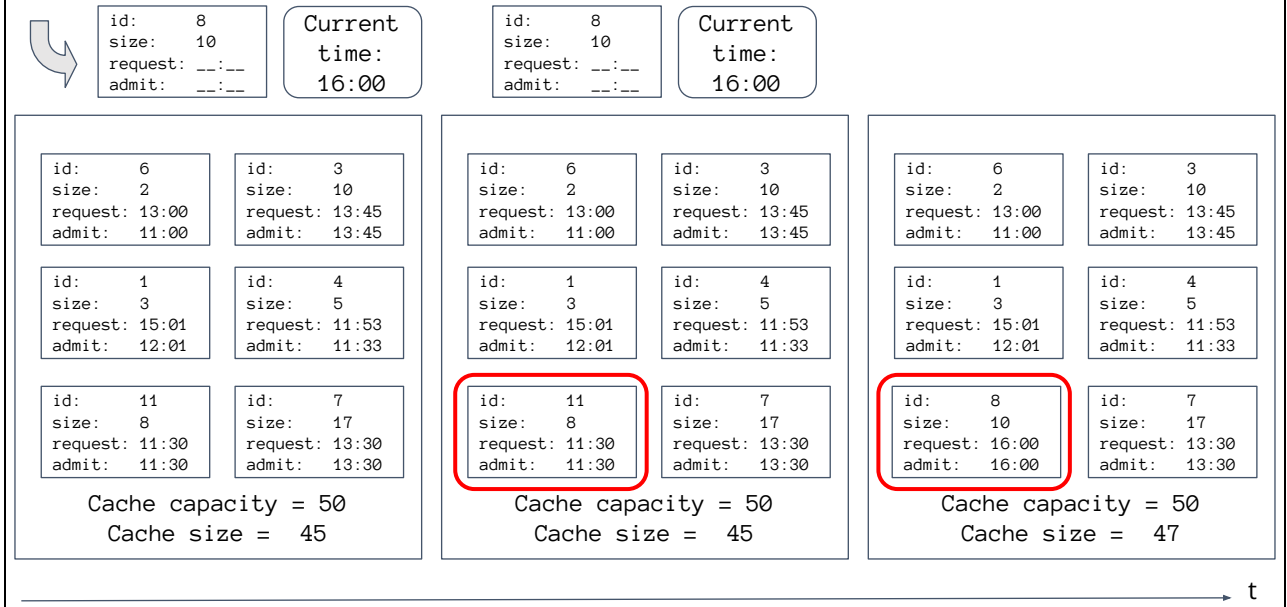
```
id:       11
size:     8
request: 11:30
admit:    11:30
```
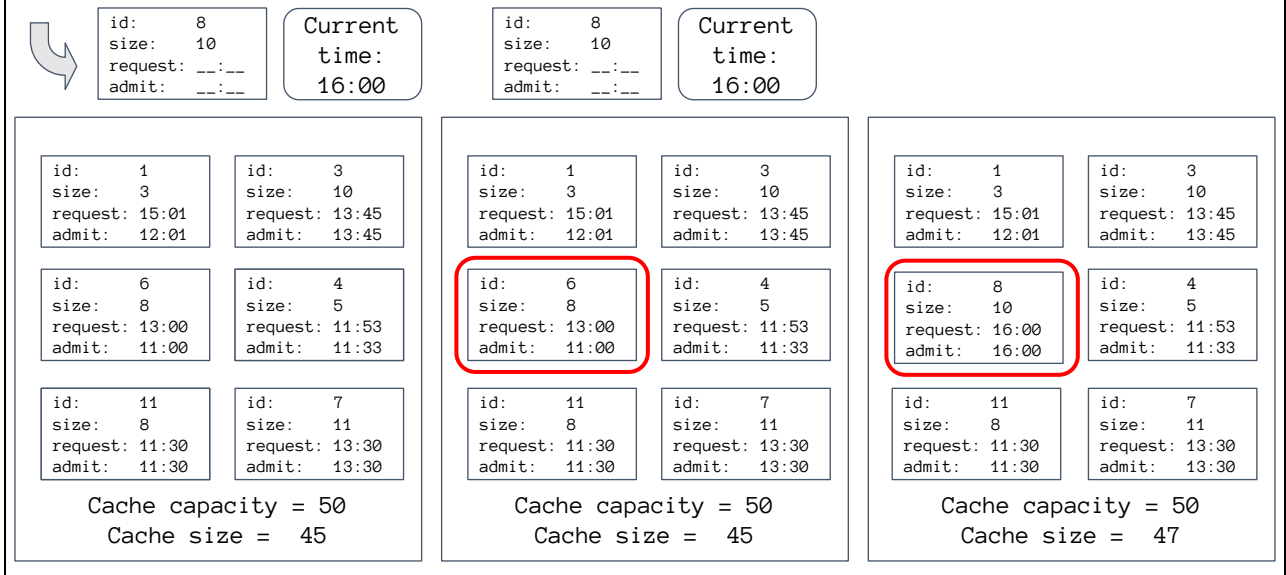
```
id:       7
size:     11
request: 13:30
admit:    13:30
```

Cache capacity = 50
Cache size =  45

- Next we have an example of the FIFO eviction algorithm
- Again, we have objects that keep track of all of the same metadata
- Again, an object comes in whose insertion would push the cache over its capacity
- We look for the object with the oldest admit time and evict it!

# FIFO

| id: 8 | | | |
| size: 10 | | | |
| request: \_\_:\_\_ | **Current time:** 16:00 | | |
| admit: \_\_:\_\_ | | | |

| id: 1 | id: 3 |
|---|---|
| size: 3 | size: 10 |
| request: 15:01 | request: 13:45 |
| admit: 12:01 | admit: 13:45 |

| id: 6 | id: 4 |
|---|---|
| size: 8 | size: 5 |
| request: 13:00 | request: 11:53 |
| admit: 11:00 | admit: 11:33 |

| id: 11 | id: 7 |
|---|---|
| size: 8 | size: 11 |
| request: 11:30 | request: 13:30 |
| admit: 11:30 | admit: 13:30 |

Cache capacity = 50
Cache size =  45

| id: 8 | | | |
| size: 10 | | | |
| request: \_\_:\_\_ | **Current time:** 16:00 | | |
| admit: \_\_:\_\_ | | | |

| id: 1 | id: 3 |
|---|---|
| size: 3 | size: 10 |
| request: 15:01 | request: 13:45 |
| admit: 12:01 | admit: 13:45 |

| id: 6 | id: 4 |
|---|---|
| size: 8 | size: 5 |
| request: 13:00 | request: 11:53 |
| admit: 11:00 | admit: 11:33 |

| id: 11 | id: 7 |
|---|---|
| size: 8 | size: 11 |
| request: 11:30 | request: 13:30 |
| admit: 11:30 | admit: 13:30 |

Cache capacity = 50
Cache size =  45

| id: 1 | id: 3 |
|---|---|
| size: 3 | size: 10 |
| request: 15:01 | request: 13:45 |
| admit: 12:01 | admit: 13:45 |

| id: 8 | id: 4 |
|---|---|
| size: 10 | size: 5 |
| request: 16:00 | request: 11:53 |
| admit: 16:00 | admit: 11:33 |

| id: 11 | id: 7 |
|---|---|
| size: 8 | size: 11 |
| request: 11:30 | request: 13:30 |
| admit: 11:30 | admit: 13:30 |

Cache capacity = 50
Cache size =  47

- Next we have an example of the FIFO eviction algorithm
- Again, we have objects that keep track of all of the same metadata
- Again, an object comes in whose insertion would push the cache over its capacity
- We look for the object with the oldest admit time and evict it!

# Experiments

> Download exercise code from [ed](ed)

> cd precept6/webcachesim-master

> make

- There's a web cache simulator on ed that we can use to experiment with different algorithms and cache sizes

# Trace File Form

- Request traces must be given in a space-separated format with three columns
- time -  long long int
- id -  long long int, used to uniquely identify objects
- size should be a long long int, object's size in bytes

•Example

| time | id | size |
|------|----|------|
| 1 | 1 | 120 |
| 2 | 2 | 64 |
| 3 | 1 | 120 |
| 4 | 3 | 14 |
| 4 | 1 | 120 |

•See test.tr

# Using the Simulator*

```
> ./webcachesim test.tr LRU 1000

LRU:1000 bytes, 10492 reqs, 8495 hits, 81 hits/reqs(%)

> ./webcachesim test.tr FIFO 1000

FIFO:1000 bytes, 10492 reqs, 8206 hits, 78 hits/reqs(%)
```

* Derived from https://github.com/dasebe/webcachesim

- Now we can see how these eviction algorithms perform by running the above commands

# Experiments

- LRU and FIFO
- Vary cache sizes
  - 80
  - 160
  - 320
  - 640
  - 1280
  - 2560
  - 5120

- Create a Google Sheet

- Three columns

- SIZE LRU FIFO

- Copy results accordingly

- Select three columns to create

  line chart

# Experiments

LRU and FIFO