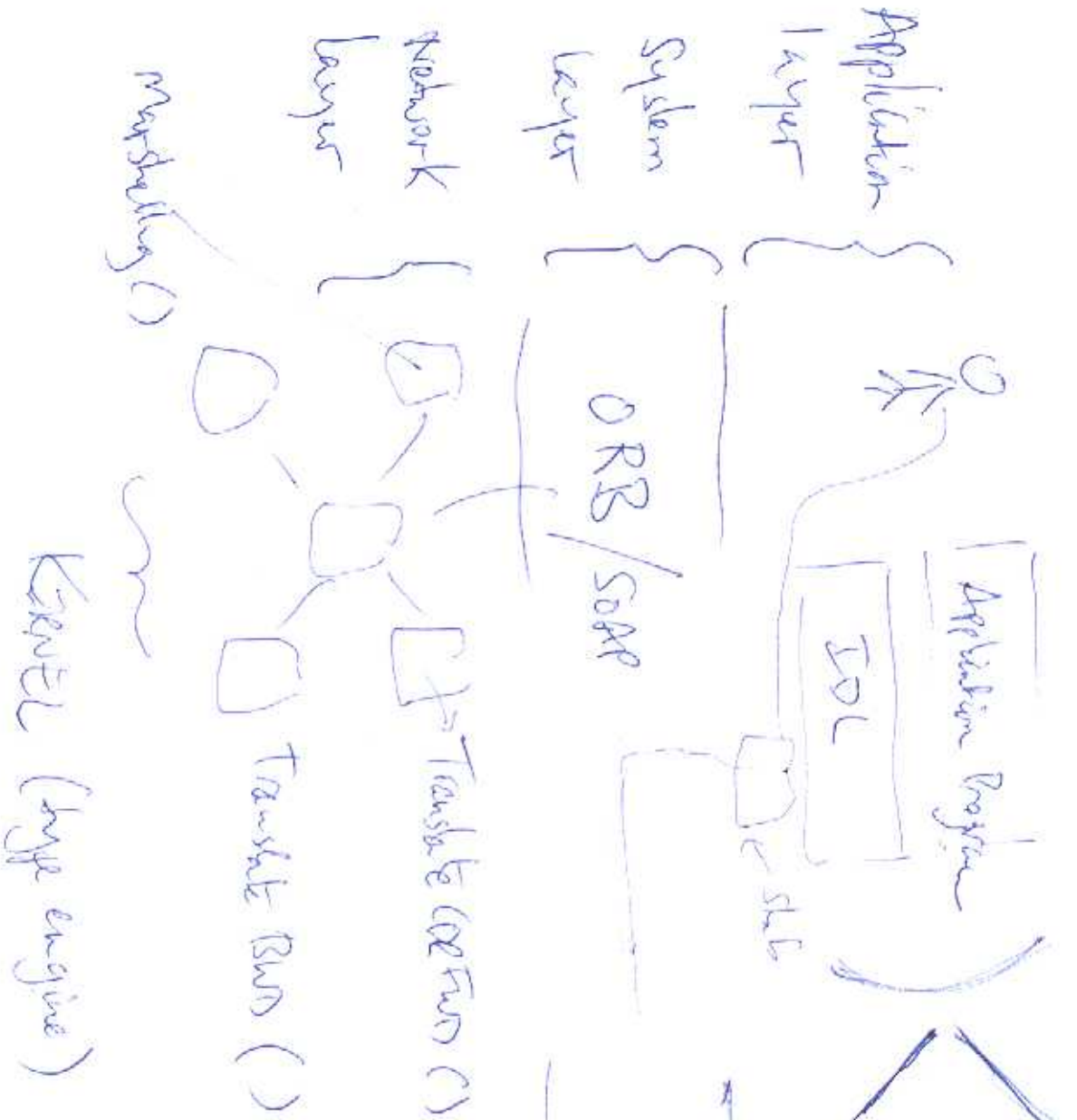


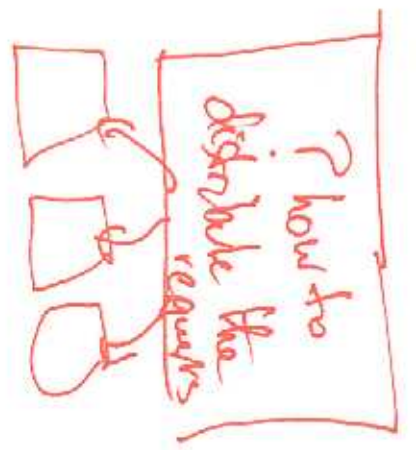
Performance (QoS) ⁰ improve response time ¹¹



Why IDL/API are designed?
 (e.g. changing the sets of refs
 into sets of structs)

Reduce the # invocations
 (independently on the applications)

implementation
 dependant

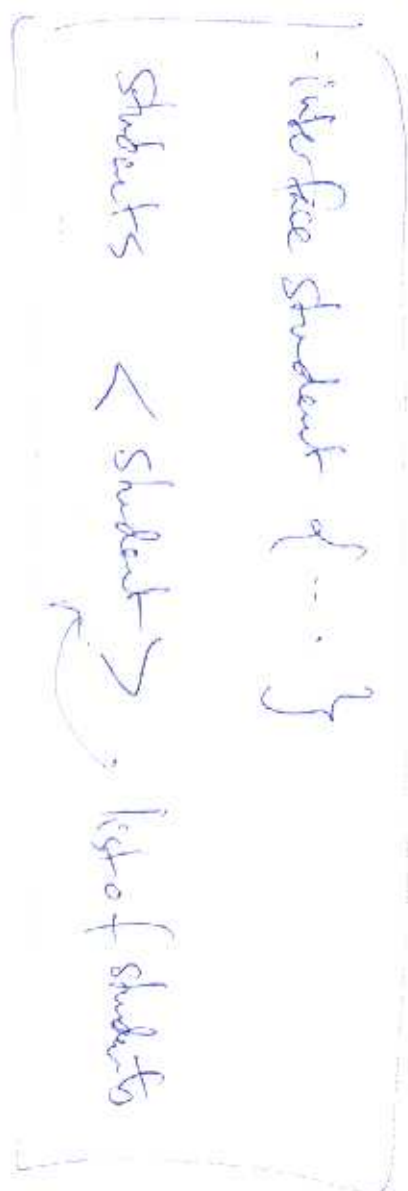


Server level
 Very high # of requests

LOAD BALANCING

LAM (Local Area Network)

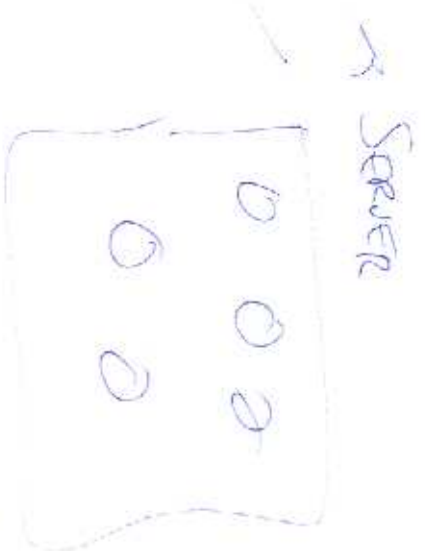
JPL



① instance of students

<ref1, ref2, ... ref 1000 >

② on each element of the set, you apply a certain operation (e.g. getNAME(), getMARK()).



① + ② numbers of calls to server \equiv $1000 + 1 \equiv 1001$

if N is the size of the set \Rightarrow $N+1$ calls

(5-2)

A possible solution is to define Student as
struct (structure)

struct Student { - }

~~Students~~
Students < Student >

only

How many calls to the server do we need to
get all the information about students? 1

< ^{name} [name, name], [], [] >

all values of the various attributes of
the different elements of struct :

(53)

CACHE

improve the locality of objects (not all objects, only those that are "critical". e.g. frequent objects, big objects, etc)



When the cache is full

eviction: which object to evict to allow a new object to be included

Consideration

↳ increase

Hit Rate \equiv High probability to find the object requested in the cache (15-4)

LOAD BALANCING is about IMPROVING ~~THE~~ RESPONSE TIME
(Performance)

by replicating the Server

REPLICATION is also about replicating the server, but
the focus is on improving AVAILABILITY

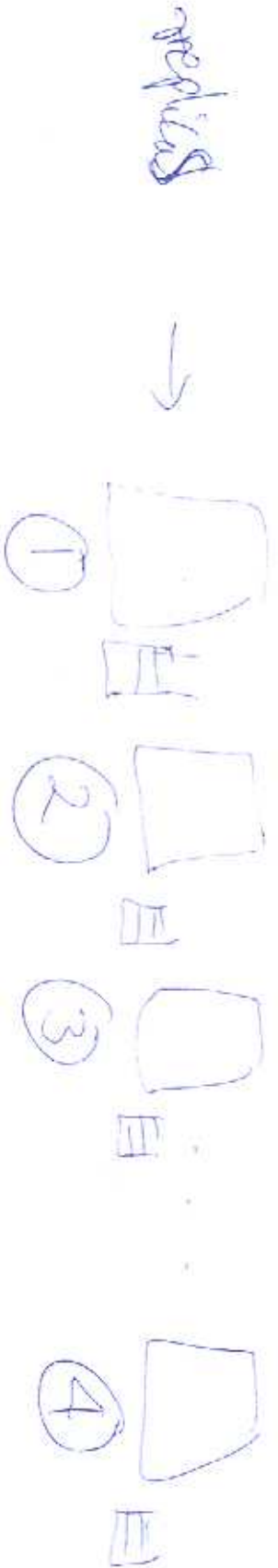
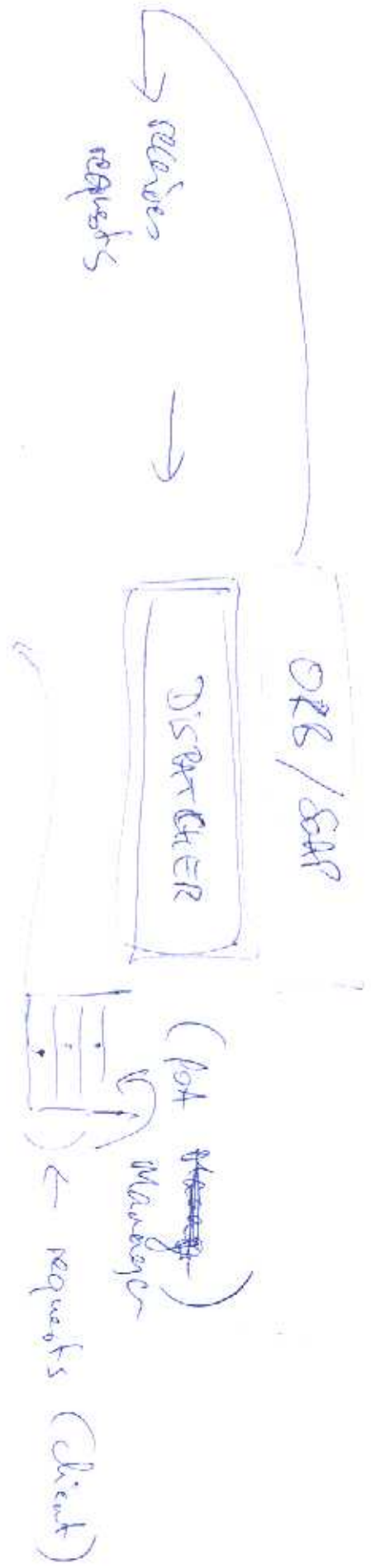
if you want BOTH (performance + availability) \Rightarrow

Sacrifice ~~to~~ the degree of performance (ie
not the best response time) and availability

(not very high \approx 100% available servers)

LS-5

Load Balancing (Server Side)



Dispatcher (Issues)

- Load INDEX (computes the load of each server)
- Task assignment (making a decision about which server to assign the task)
- Task migration

15-6

When you know the load index, normally speaking the best assignment policy is straight-forward (ie because you assign tasks to the least load first (LWF) server as this server has the smallest # of tasks in its queue)

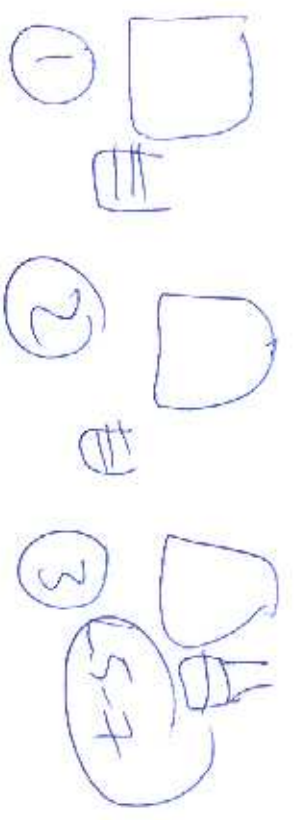
TRUE: when the task size distribution

is EXPONENTIAL

LWF is the best strategy

$$P(x \leq x_0) \approx e^{-\lambda x_0}$$

$$x \uparrow P(\text{find a big task}) \approx 0$$



High variability
in the task size
(eg HTTP traffic)



LLF (is based on
the shortest queues)
is NOT optimal



FALSE :
the task distribution is
not exponential
(PARETO distribution)

$$P(X > x) \neq 0$$

probability that a task
goes beyond a certain
threshold is not nil

(L5-8)

Slows Down (SD)

because you have high variability in the task size distribution (ie it is not exponential)

$P(X > x_0) \neq$ = probability that the task size is

larger than x_0

$P(X > x_0)$ when $x_0 \uparrow$ (high) \neq

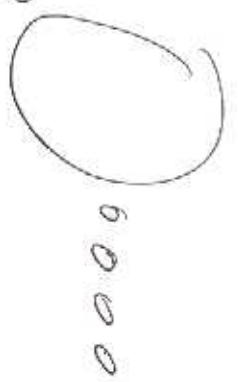
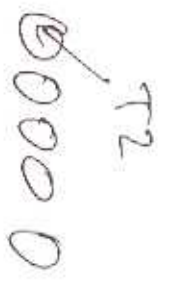
$P(X > x_0)$ when $x_0 \uparrow$ (high) \neq

(~~HTTP~~ traffic) Rareho

Exponential

80

80% of tasks are small



Lot of small tasks

in order to compare T_1 and T_2 (you need to use a different metric (than waiting time))

Aim

$$\frac{\text{waiting time } (T_1)}{\text{Task Size } (T_1)} = \frac{\text{waiting time } (T_2)}{\text{Task Size } (T_2)}$$

fairness

↳ fairness \equiv treat the tasks in the same way i.e. the ratio of waiting time by the task size is similar for all tasks

T_1 / Task size = 100
 / Waiting Time = 10

T_2 / task size = 10
 / waiting Time = 1

15-10