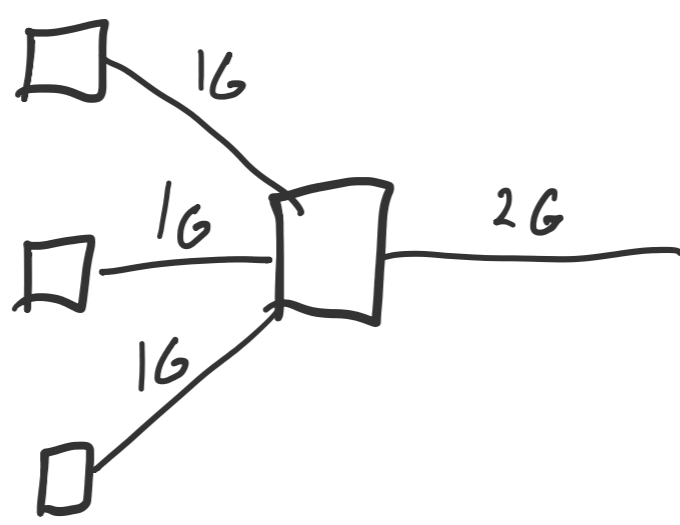
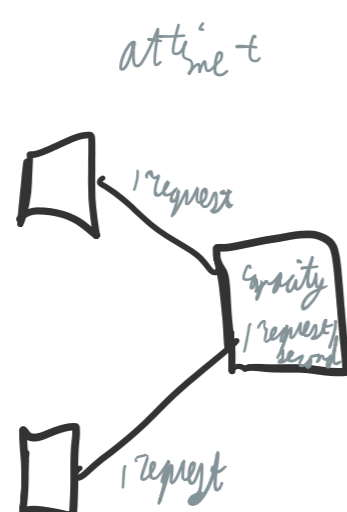


over-subscribe



overload



no uniform definition of scalability in literature

definitions:

1. ability to handle increasing workload *without adding resources*
2. ability to handle increased workload by repeatedly applying a cost-effective strategy for setting a system's capacity

load scalability

deal with increased load

space scalability

deal with increased memory demand *generally not storage*

space-time scalability

deal with increased demand for number of objects

structural scalability

deal with expanding system without major modifications to its architecture

slide 15, middle figure: network limit = $\frac{16 \text{ servers} \cdot 12.5 \text{ MB/s}}{3 \text{ replicas}} \approx 67 \text{ MB/s}$

scaling has no effect on the sequential part of the system, and a linear effect on the parallel part

Amdahl's law keeps the problem size; then

$$\text{speedup} = \frac{1}{1-p + \frac{p}{N}}$$

with p the portion of the program which can run in parallel and N the number of processors

Gustafson's law observes that parallelism increases in an application when the problem size increases

$$\text{speedup} = S + p \cdot N$$

with S the execution time of the serial portion, p the parallelizable portion of the program, and N the number of processors

Amdahl's law: if tuning one variable does not work, tune another one
 Gustafson's law: use more resources to solve larger problems

tactics for scalability:

- limit dependencies
- no machine has/needs complete information
- decisions are based on local information
- failure of one machine does not ruin results
- no assumption of global/shared clock

techniques

- hide latency *do something useful while waiting*
- introduce concurrency
- partition
- limit communication

scalability is a prerequisite for elasticity

elasticity is the ability of a system to adapt to workload changes by provisioning and deprovisioning in an automatic manner

adaptation mechanism: auto scalers

aspects of elasticity:

- timing: % of under/overprovisioned time
- frequency of adaptation
- latency: time to bring up/down resources
- accuracy: relative deviation of allocated resources from actual demand



jitter $y = \frac{E_s - E_d}{T}$, where E_s is the number of supply changes during interval T and similarly for E_d