

Distribuované systémy

Obsah

- Proč distribuované systémy
 - Sjednocení pojmů
 - CAP theorem
 - Replikace
-
- <http://book.mixu.net/distsys/single-page.html>

Když server nestíhá

- Upgrade HW
- Co pak?

Škálovatelnost

Schopnost systému zvládnout zvětšující se množství práce rozumným způsobem

- Velikostní
- Geografická
- Administrativní

Výkon

- Nízká latence
- Vysoká propustnost
- Nízké využití výpočetního výkonu

Dostupnost

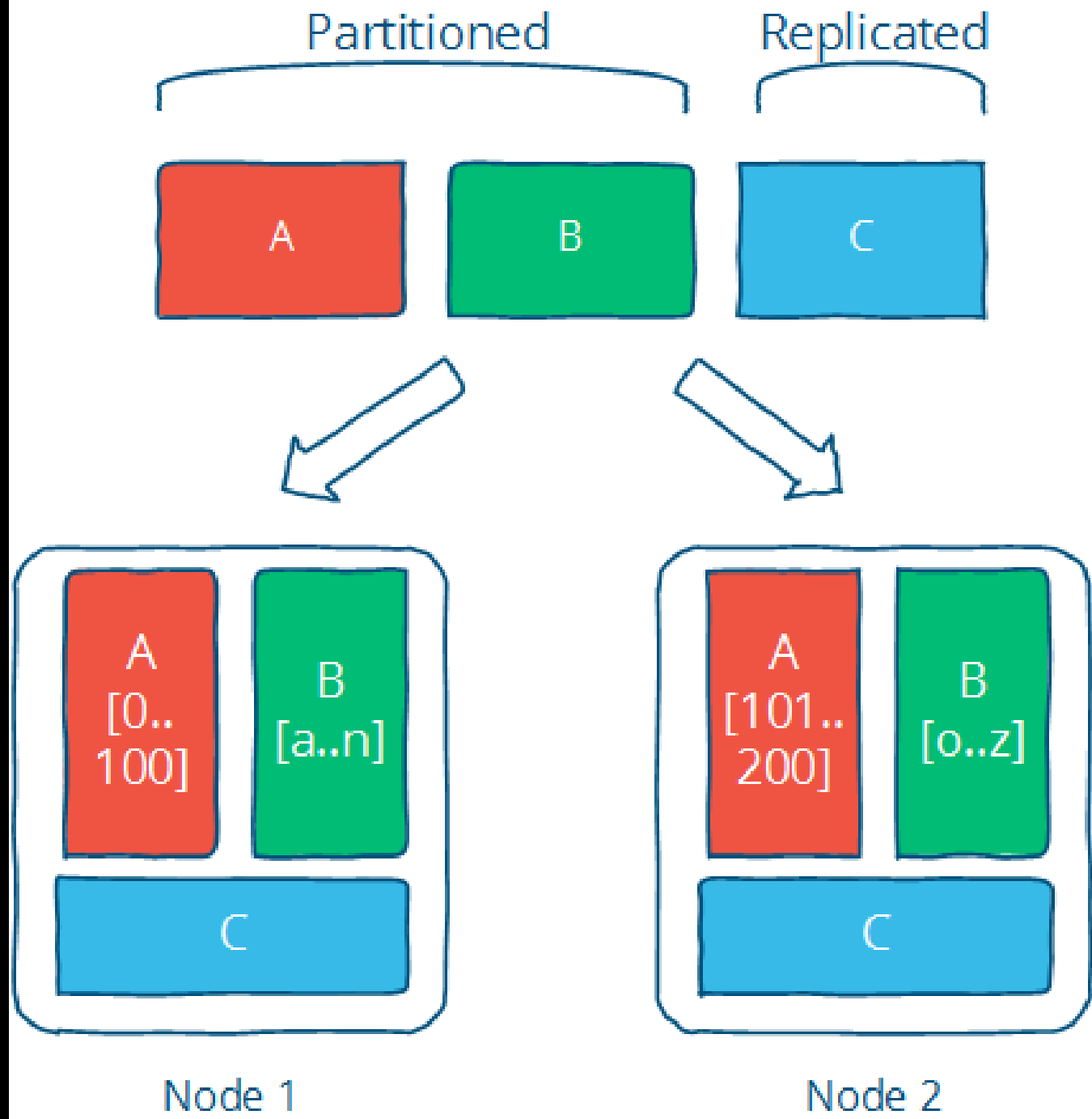
- Procento času, po který byl systém „funkční“
- Jeden stroj nemůže být fault tolerant

Fault tolerance

- Schopnost systému chovat se „nějak“ při výpadku
- Výpadku čeho?

Důsledky distribuovanosti

- Více uzlů => větší pravděpodobnost selhání
- Více uzlů => více komunikace mezi uzly
- Geografická vzdálenost => latence



Partitioning

- + výkon (divide & conquer)
- + dostupnost celku (fail independently)

Replikace

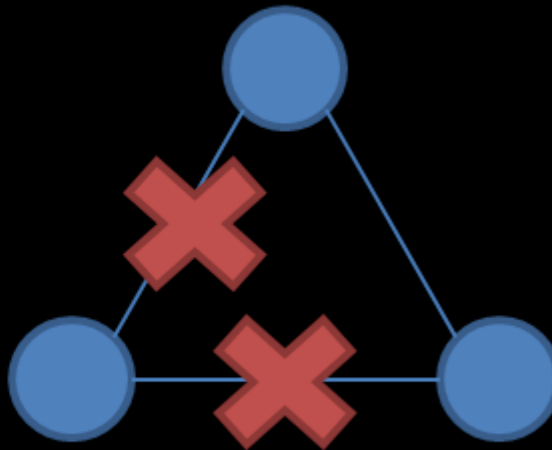
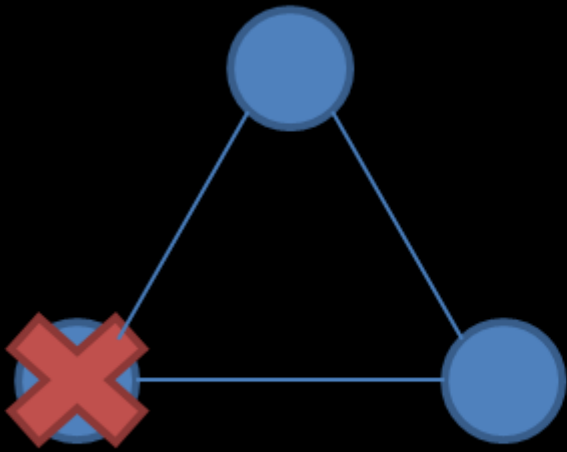
- + výkon (více nodů může řešit stejný problém)
- + dostupnost ($n+1$ uzlů)
- konzistence ?

Model systému

Sada předpokladů o systému

- Co vše může selhat
- Pořadí operací
- Lokální vs vzdálený stav
- Synchronizace času / dat
- Pád sítě / pád uzlu

Network partition



Synchronní systém

- Lock, write, unlock
- Omezený výkon

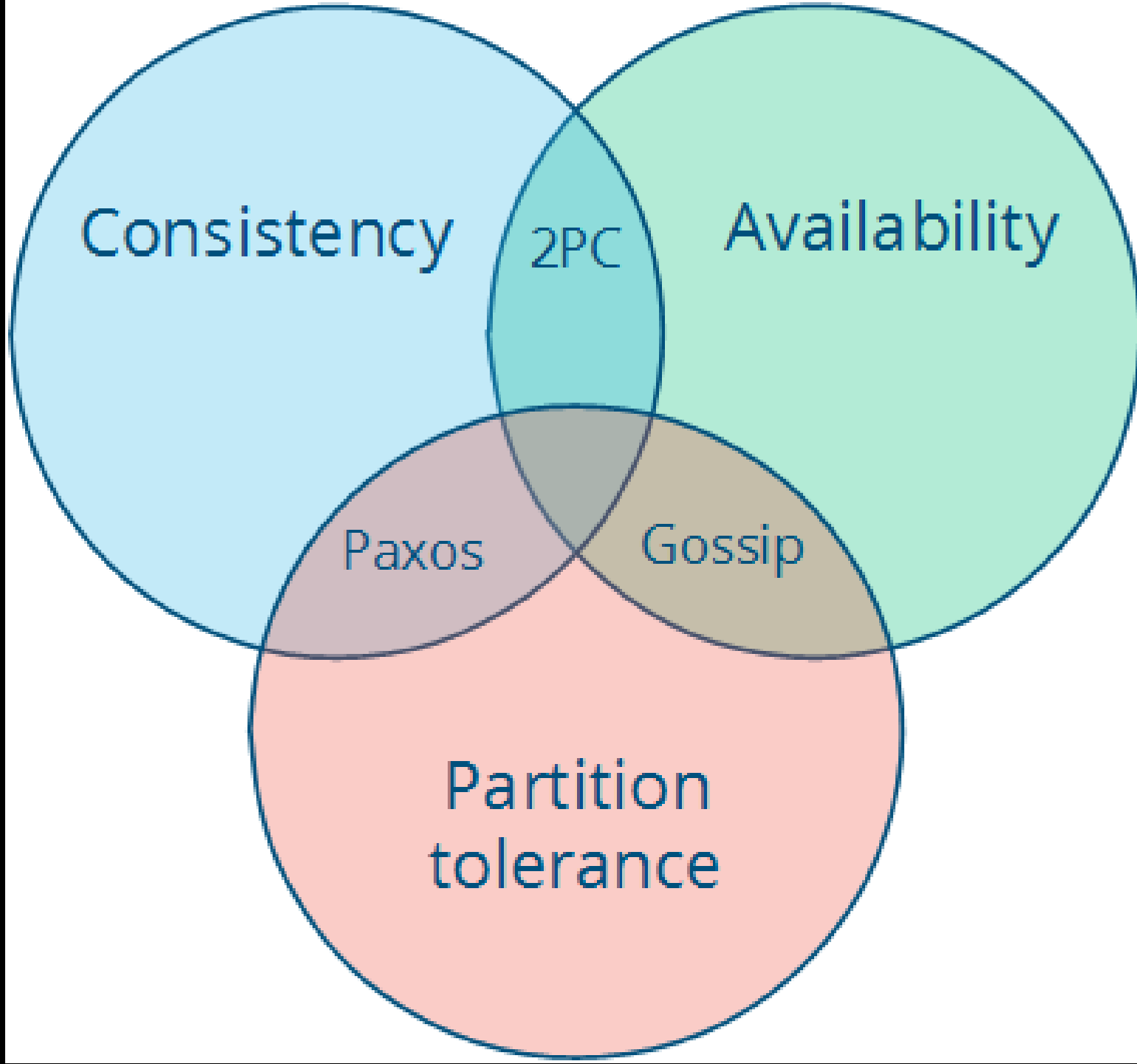
Asynchronní systém

- Žádný předpoklad o časování/pořadí

Problém konsenzu

Uzly dosáhnou konsenzu, pokud se dohodnou na nějaké jedné hodnotě.

- Dohoda
- Integrita
- Ukončení
- Validita

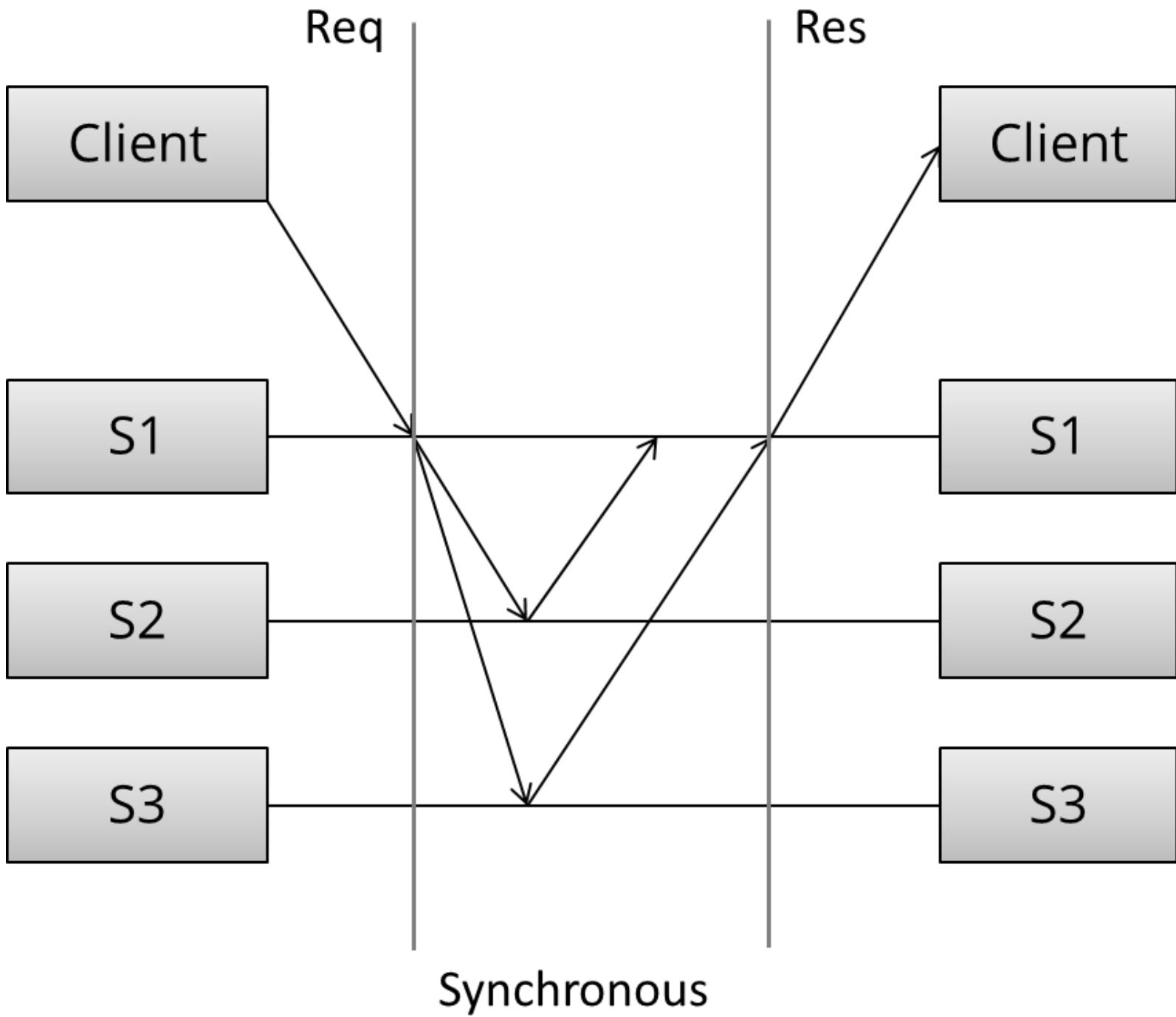


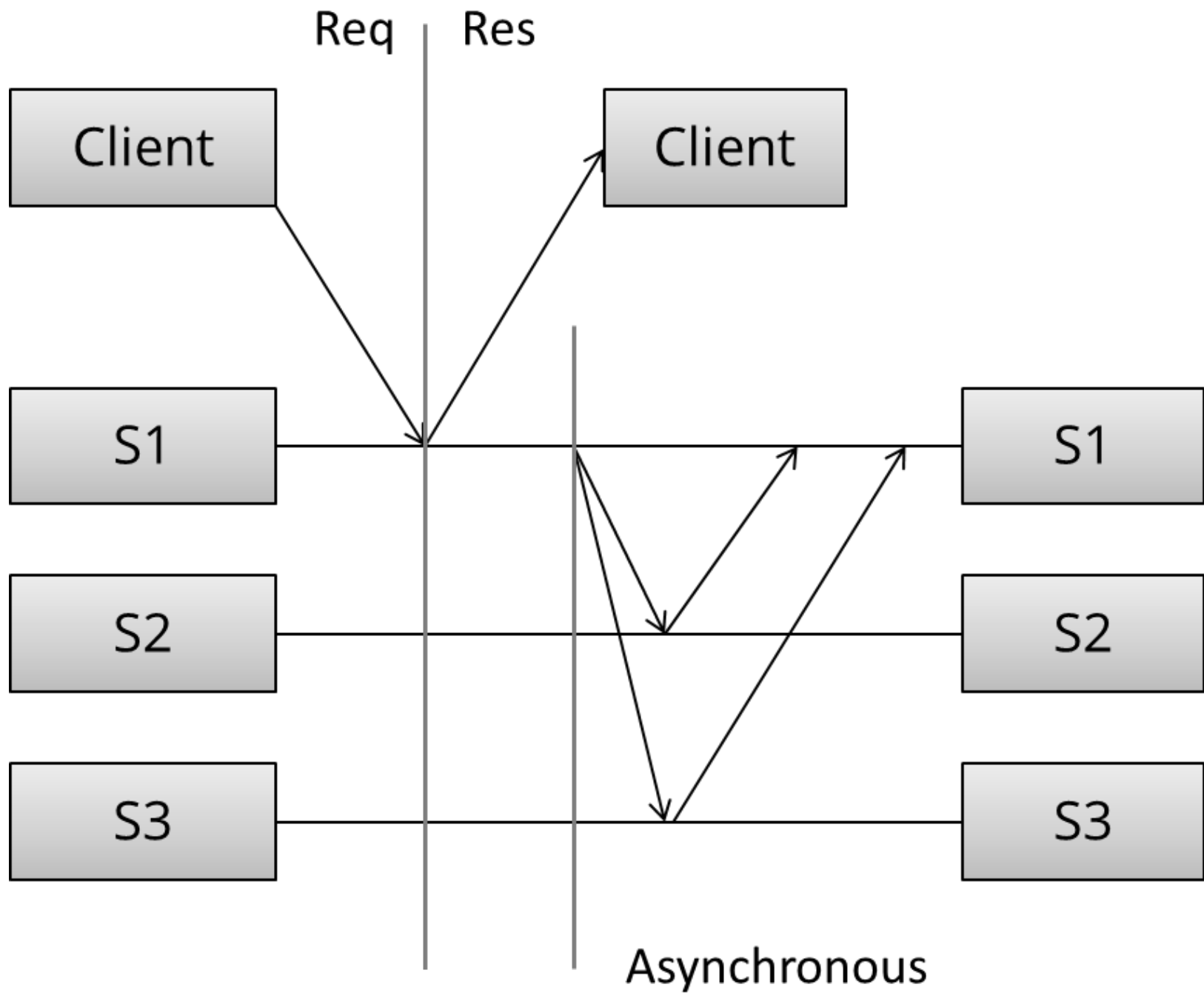
Pozorování

- Mnoho starších databázových systémů je CA
- Musíme si vybrat A nebo C při network partition
- Silná konzistence vs výkon
- Pokud chceme A+P, musíme snížit C

Pořadí a čas

- Úplné seřazení (total order)
- Částečné (partial order)
- Čas => pořadí
- Globální / lokální čas
- Serializovatelnost
- Detekce selhání (failure)





2-phase commit

[Coordinator] -> OK to commit? [Peers]
 <- Yes / No

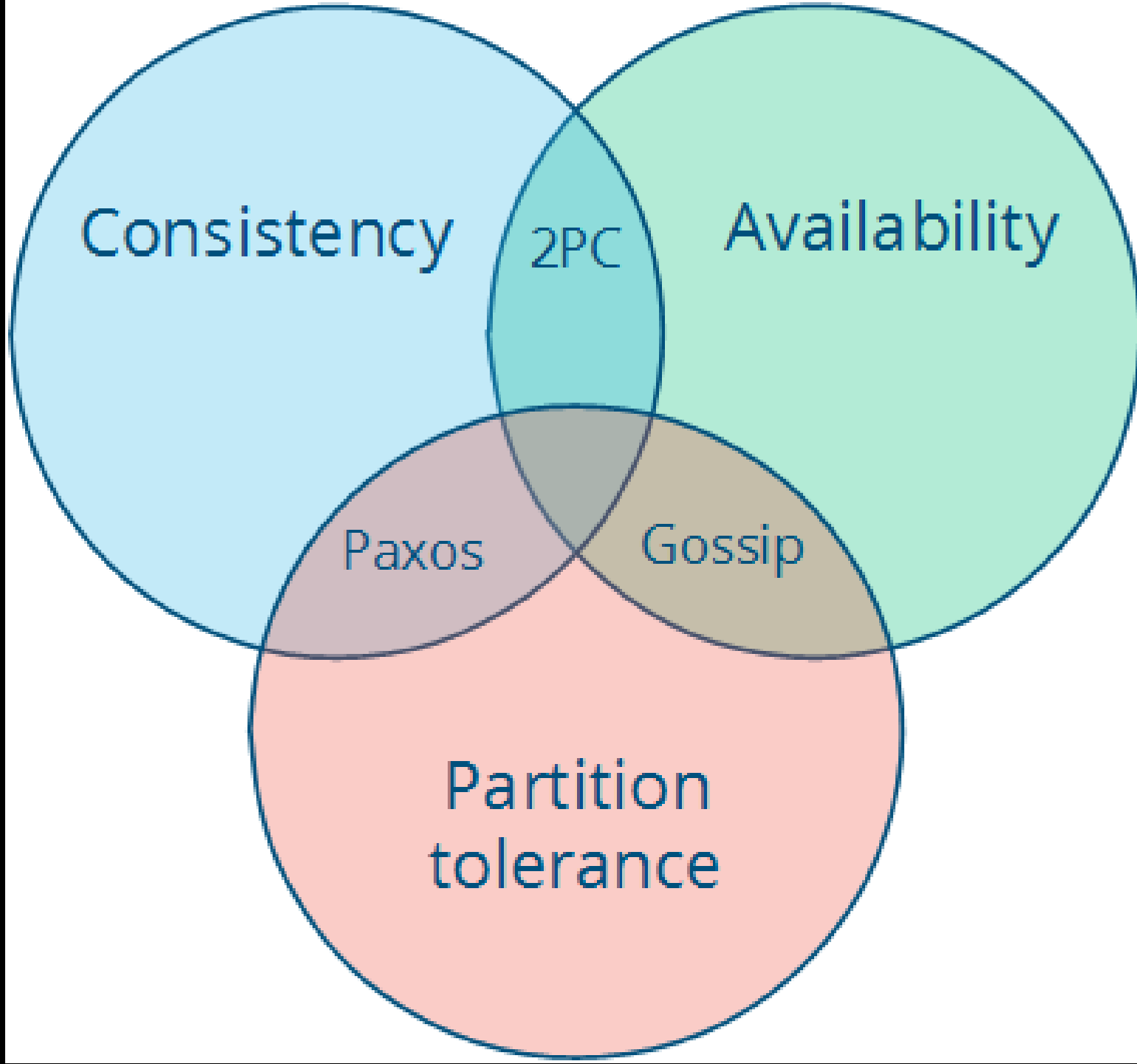
[Coordinator] -> Commit / Rollback [Peers]
 <- ACK

Paxos / Raft

- Hlasování většiny
- Dynamický master
- $n/2 - 1$ nedostupných nodů je OK

CP systémy

- Raft (etcd)
- Paxos (ZooKeeper)

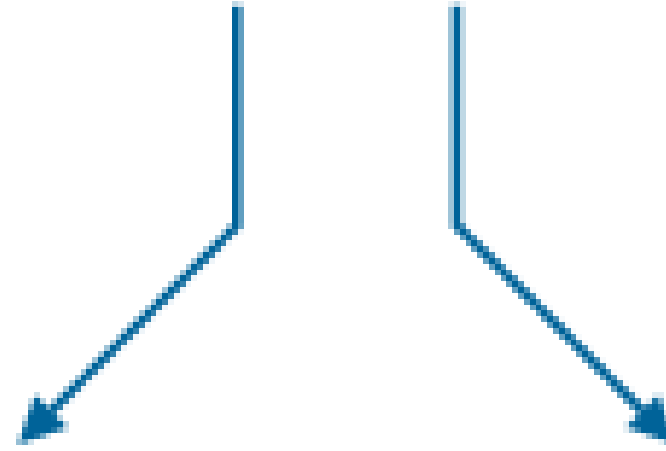


Partition
tolerance

Consistency

?

Availability

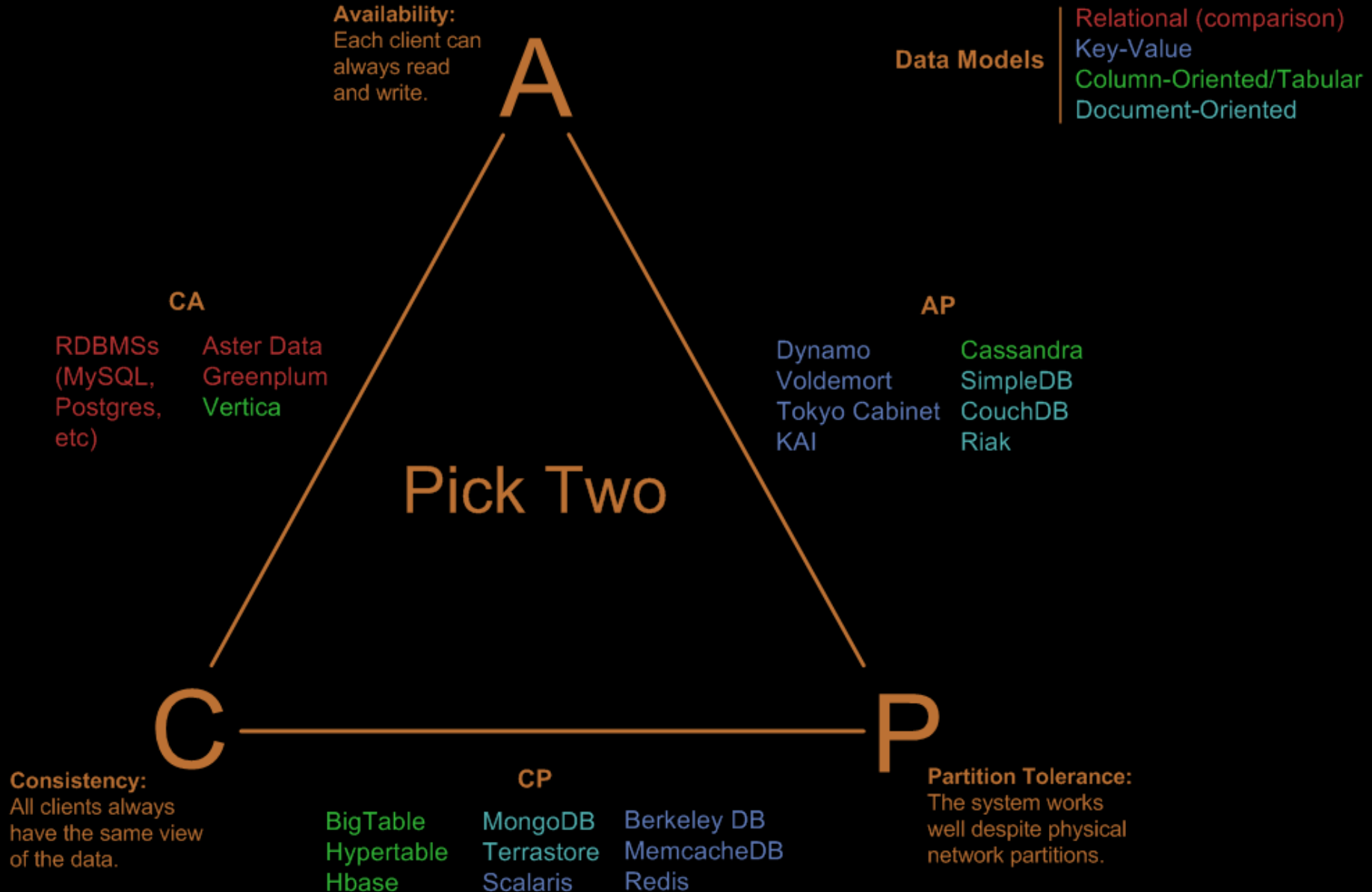


Děkuju za pozornost.

Štefan Šafár

<http://book.mixu.net/distsys/single-page.html>

Visual Guide to NoSQL Systems



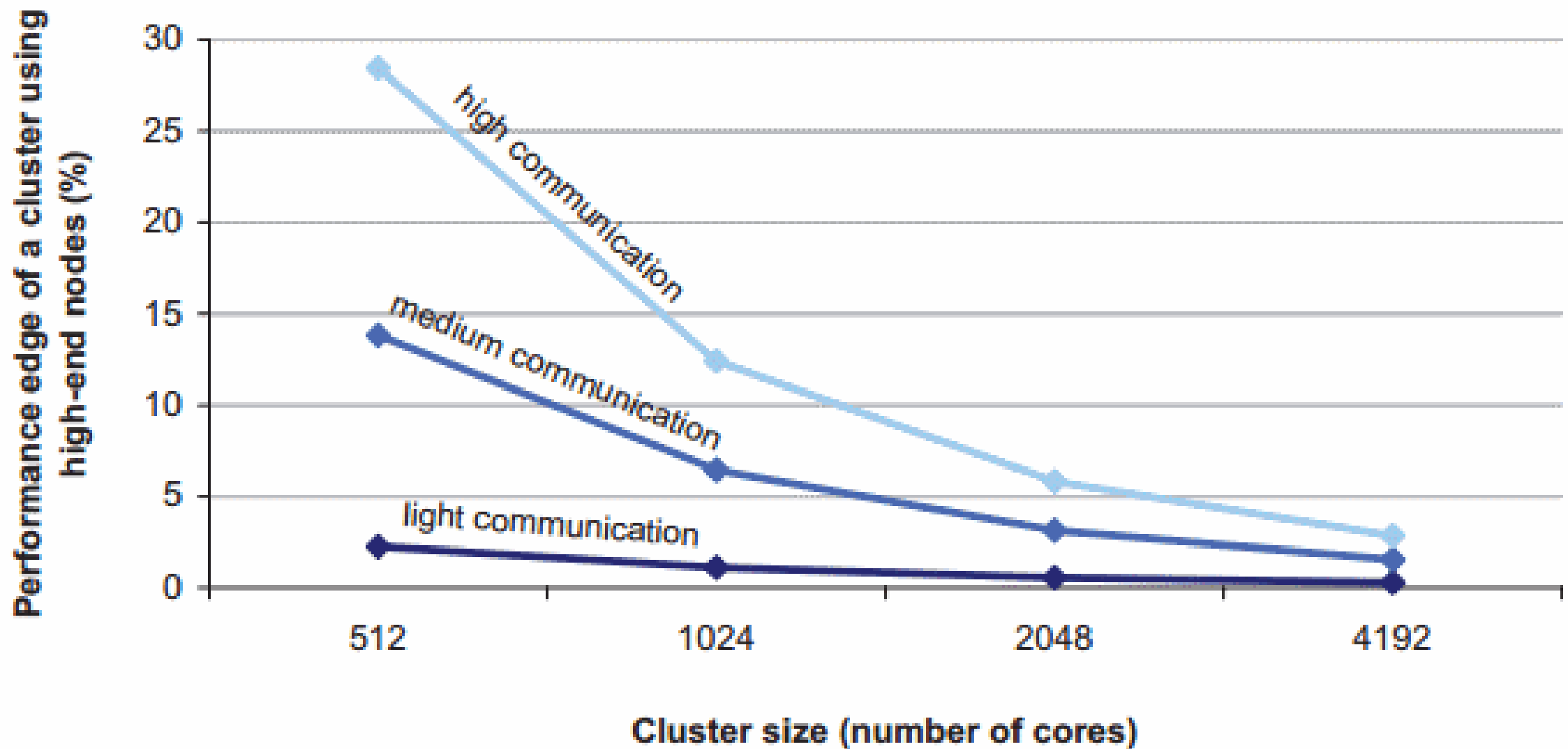


FIGURE 3.2: Performance advantage of a cluster built with high-end server nodes (128-core SMP) over a cluster with the same number of processor cores built with low-end server nodes (four-core SMP), for clusters of varying size.