

The logo for Quorade features the word "Quorade" in a blue, cursive script font. A yellow ring, resembling a planet's ring system, is positioned around the letter 'Q'.

Quorade

Michael Whittaker, Aleksey
Charapko, Joseph M. Hellerstein,
Heidi Howard, Ion Stoica

Replica 1

Key	value
x	0
y	0

Replica 2

Key	value
x	0
y	0

Replica 3

Key	value
x	0
y	0



Client 1



Client 2

Replica 1

Key	value
x	0 1
y	0

Replica 2

Key	value
x	0
y	0

Replica 3

Key	value
x	0
y	0

OK Write(x,1)

Client 1

Client 2

Replica 1

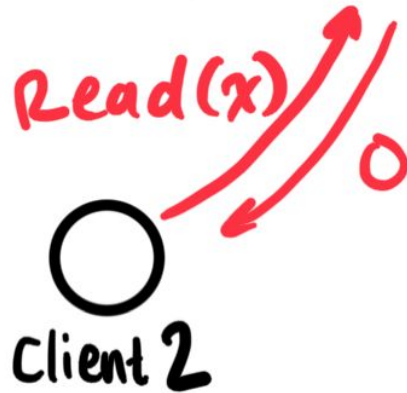
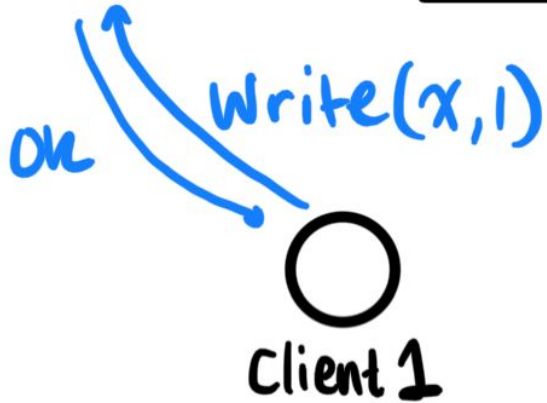
Key	value
x	0 1
y	0

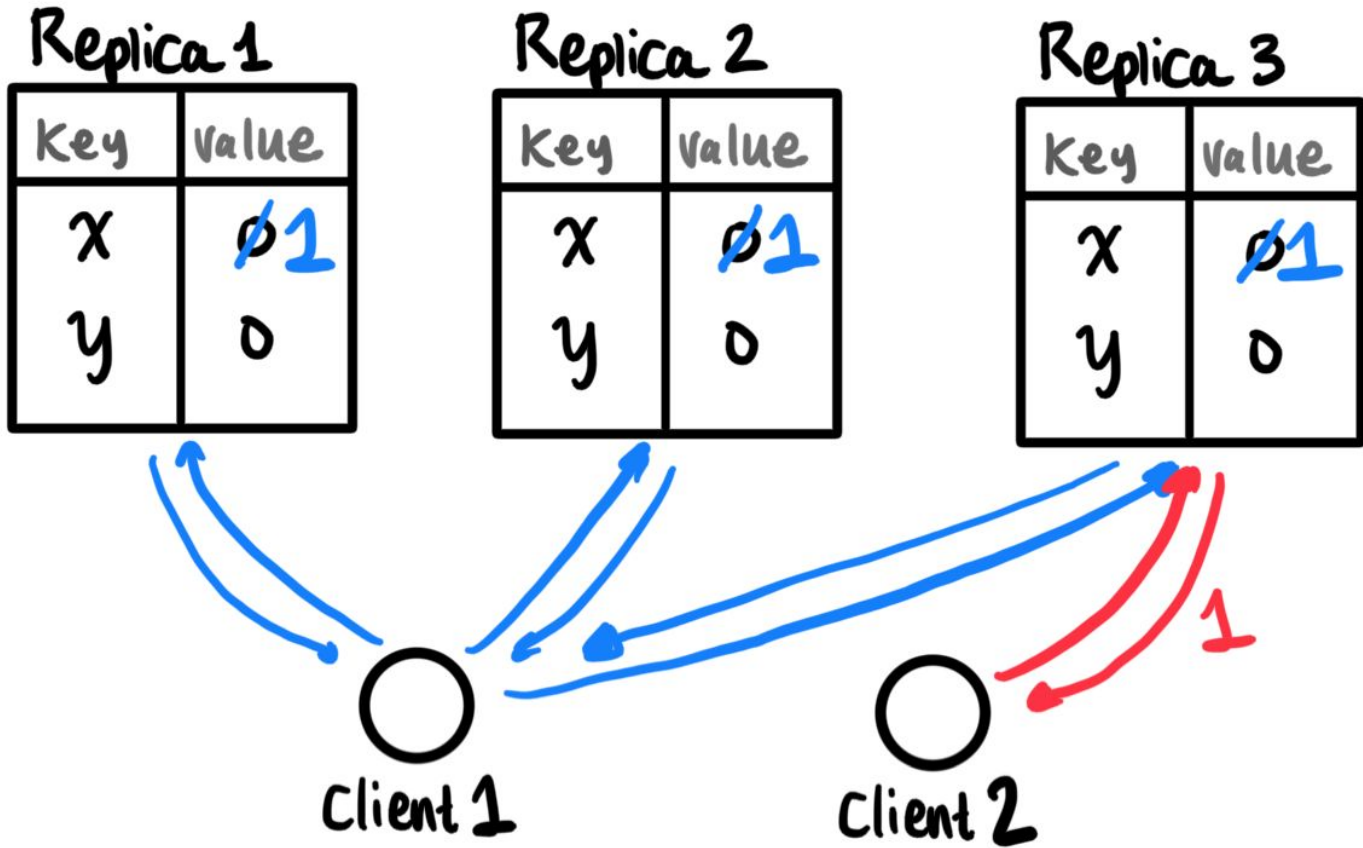
Replica 2

Key	value
x	0
y	0

Replica 3

Key	value
x	0
y	0





Given a set $X = \{x_1, \dots, x_n\}$, a **read-write quorum system** over X is a pair $Q = (R, W)$ where

Given a set $X = \{x_1, \dots, x_n\}$, a **read-write quorum system** over X is a pair $Q = (R, W)$ where

- ① R is a set of subsets of X called **read quorums**.

Given a set $X = \{x_1, \dots, x_n\}$, a **read-write quorum system** over X is a pair $Q = (R, W)$ where

- ① R is a set of subsets of X called **read quorums**.
- ② W is a set of subsets of X called **write quorums**.

Given a set $X = \{x_1, \dots, x_n\}$, a **read-write quorum system** over X is a pair $Q = (R, W)$ where

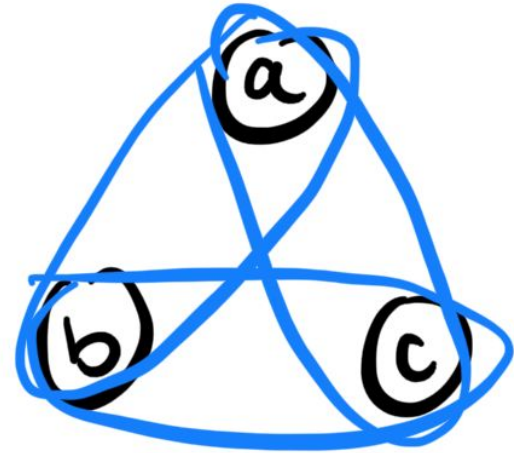
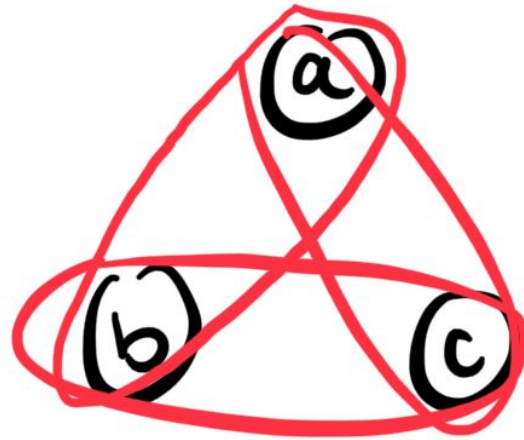
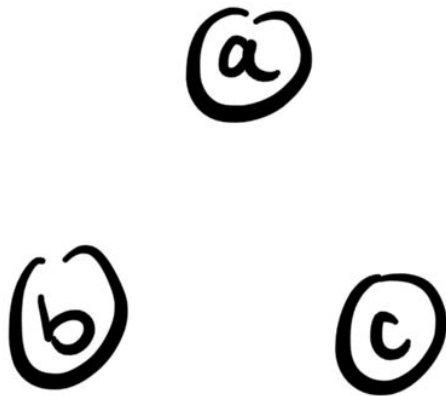
- ① R is a set of subsets of X called **read quorums**.
- ② W is a set of subsets of X called **write quorums**.
- ③ Every read quorum intersects every write quorum: $\forall r \in R. \forall w \in W. r \cap w \neq \emptyset$

The majority quorum system.

$$X = \{a, b, c\}$$

$$R = \{ab, bc, ac\}$$

$$W = \{ab, bc, ac\}$$



The grid quorum system.

$X = \{a, b, c, d\}$

Read quorums
 $R = \{ab, cd\}$

Write quorums
 $W = \{ac, bd\}$

(a) (b)

(a) (b)

(a)

(b)

(c) (d)

(c) (d)

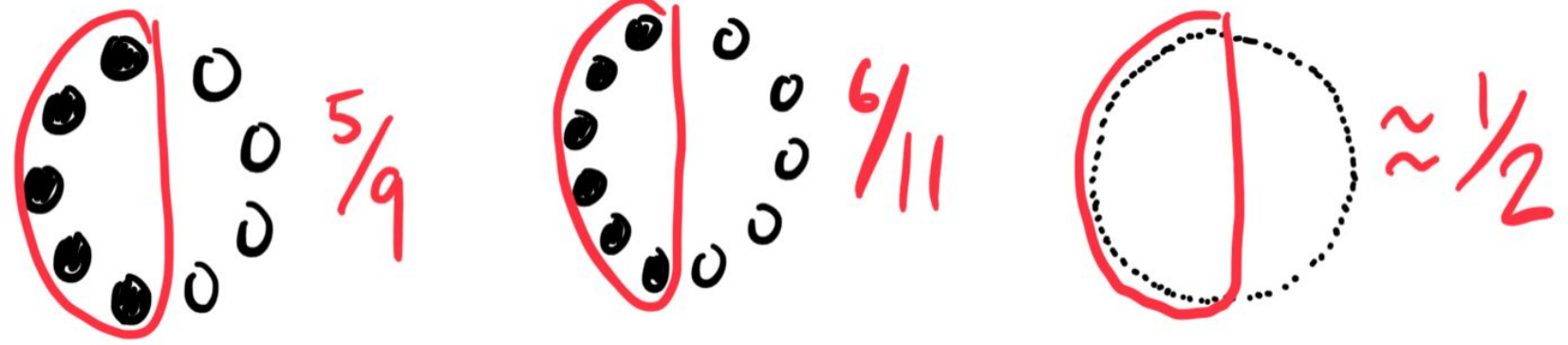
(c)

(d)

3 nodes... → 5 nodes... → 7 nodes...



9 nodes... → 11 nodes... → ∞ nodes!



4 nodes...



6 nodes...



8 nodes...



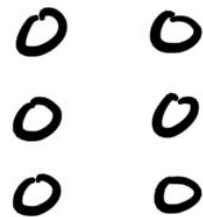
$1/2$



$1/3$



$1/4$



10 nodes...



12 nodes...



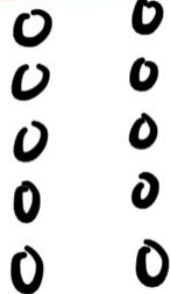
∞ nodes!



$1/5$



$1/6$

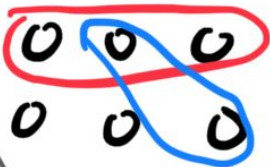


≈ 0

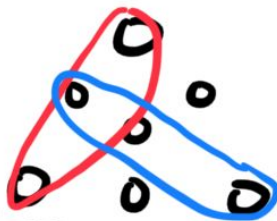
Majority Quorums [1]



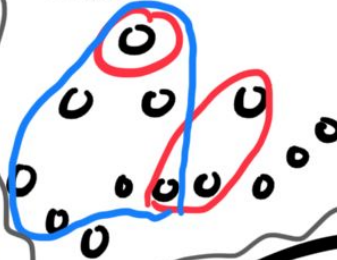
Grid Quorums [1]



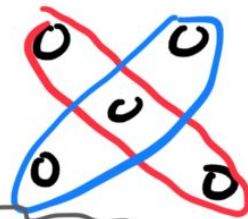
Finite Projective Planes



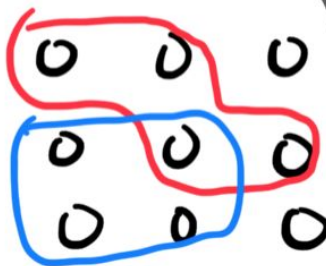
Tree Quorums [2]



Path Quorums [3]



WPaxos Quorums [4]



Which one do I pick?



[1]: The Origin of Quorum Systems

[2]: The Tree Quorum Protocol: An Efficient Approach for Managing Replicated Data

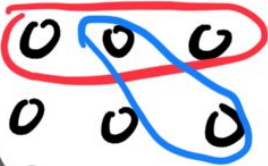
[3]: The Load, Capacity, and Availability of Quorum Systems

[4]: WPaxos: Wide Area Network Flexible Consensus

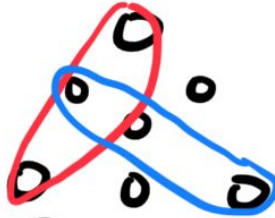
Majority Quorums [1]



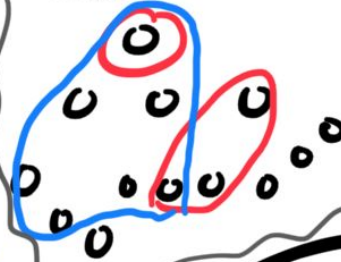
Grid Quorums [1]



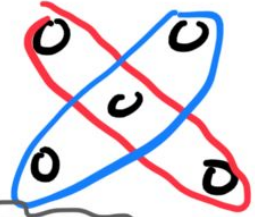
Finite Projective Planes



Tree Quorums [2]



Path Quorums [3]

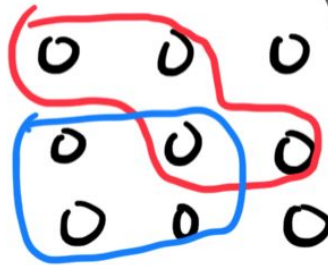


I can help with that!



Quoracle

WPaxos [4] Quorums



Which one do I pick?



[1]: The Origin of Quorum Systems
[2]: The Tree Quorum Protocol: An Efficient Approach for Managing Replicated Data
[3]: The Load, Capacity, and Availability of Quorum Systems
[4]: WPaxos: Wide Area Network Flexible Consensus

New quorum
system theory



Python
library/tool

$$f_r \left(\sum_{\{r \in R \mid x \in R\}} P_r \right) +$$

$$(1 - f_r) \left(\sum_{\{w \in W \mid x \in W\}} P_w \right)$$

```
from quoraele import *
```

```
a = Node('a')
```

```
v = Node('v')
```

```
t = Node('t')
```

```
n = Node('n')
```

```
qs = QuorumSystem(a*v*a + t*a*n)
```



```
from quoracle import *

a = Node('a')
b = Node('b')
c = Node('c')
majority = QuorumSystem(reads=a*b + b*c + a*c)

print(majority.fault_tolerance()) # 1
print(majority.load(read_fraction=1)) # 2/3
print(majority.capacity(read_fraction=1)) # 3/2
```

```
from quoracle import *

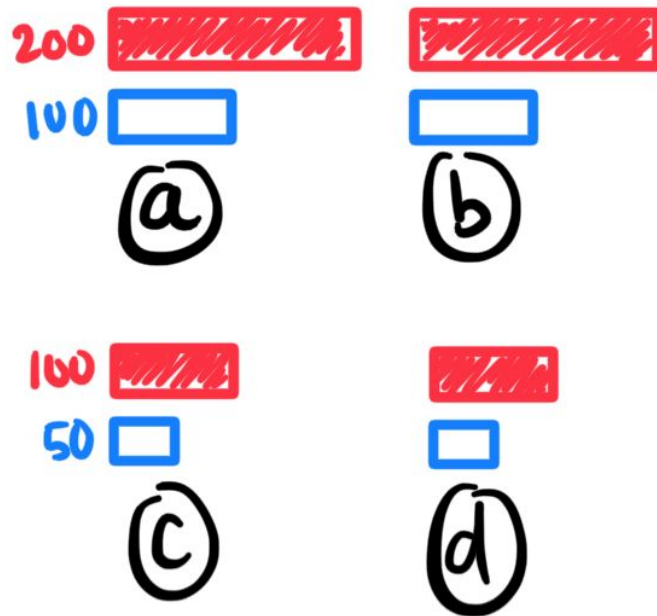
a = Node('a', write_cap=100, read_cap=200)
b = Node('b', write_cap=100, read_cap=200)
c = Node('c', write_cap=50, read_cap=100)
d = Node('d', write_cap=50, read_cap=100)
grid = QuorumSystem(reads=a*b + c*d)

print(grid.capacity(read_fraction=1)) # 300
print(grid.capacity(read_fraction=0.5)) # 200
print(grid.capacity(read_fraction=0)) # 100
```

```
from quorace import *
```

```
a = Node('a', write_cap=100, read_cap=200)  
b = Node('b', write_cap=100, read_cap=200)  
c = Node('c', write_cap=50, read_cap=100)  
d = Node('d', write_cap=50, read_cap=100)  
grid = QuorumSystem(reads=a*b + c*d)
```

```
print(grid.capacity(read_fraction=1)) # 300  
print(grid.capacity(read_fraction=0.5)) # 200  
print(grid.capacity(read_fraction=0)) # 100
```

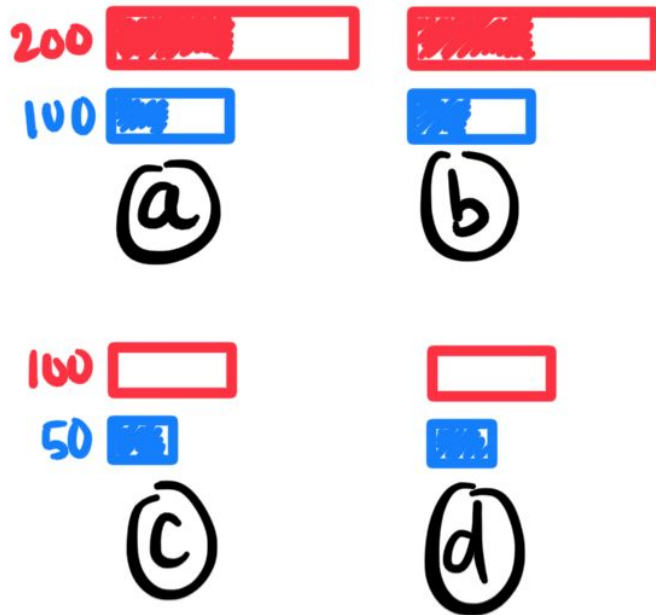


100% reads

```
from quorace import *
```

```
a = Node('a', write_cap=100, read_cap=200)  
b = Node('b', write_cap=100, read_cap=200)  
c = Node('c', write_cap=50, read_cap=100)  
d = Node('d', write_cap=50, read_cap=100)  
grid = QuorumSystem(reads=a*b + c*d)
```

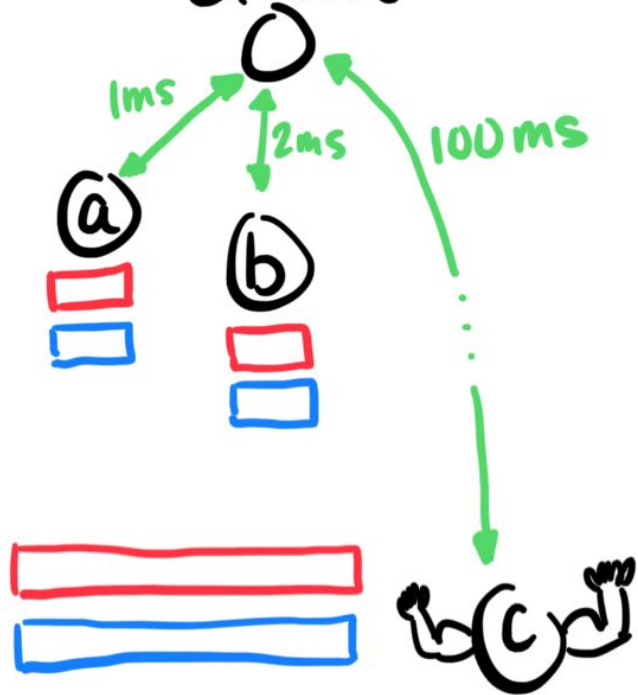
```
print(grid.capacity(read_fraction=1)) # 300  
print(grid.capacity(read_fraction=0.5)) # 200  
print(grid.capacity(read_fraction=0)) # 100
```



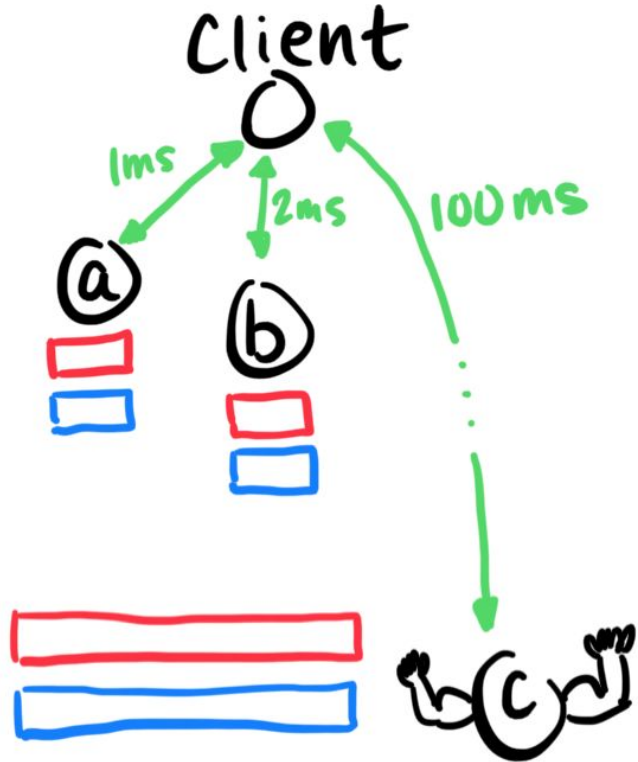
50% reads
50% writes

Latency

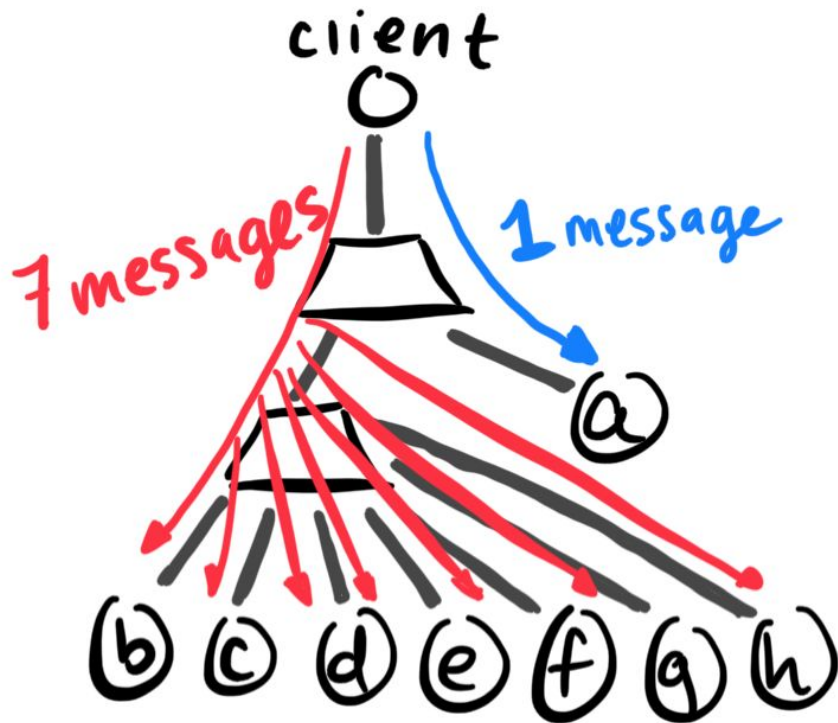
Client



Latency



Network Load



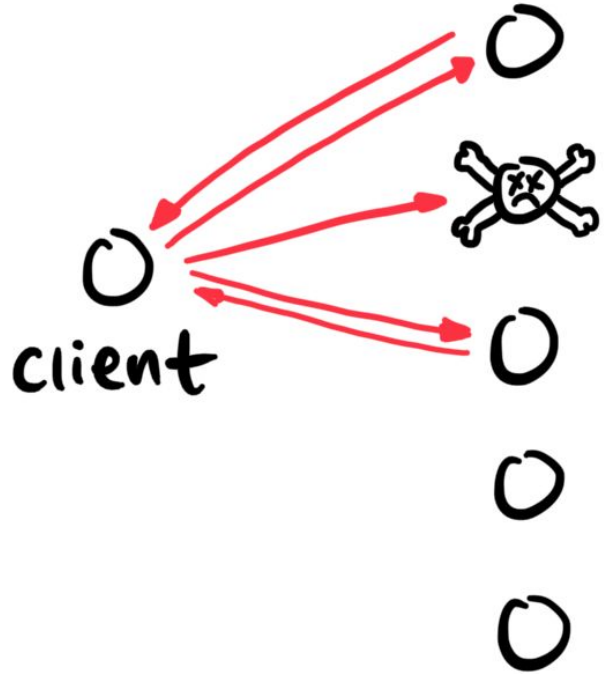
```
import datetime

def seconds(x: int) -> datetime.timedelta:
    return datetime.timedelta(seconds=x)

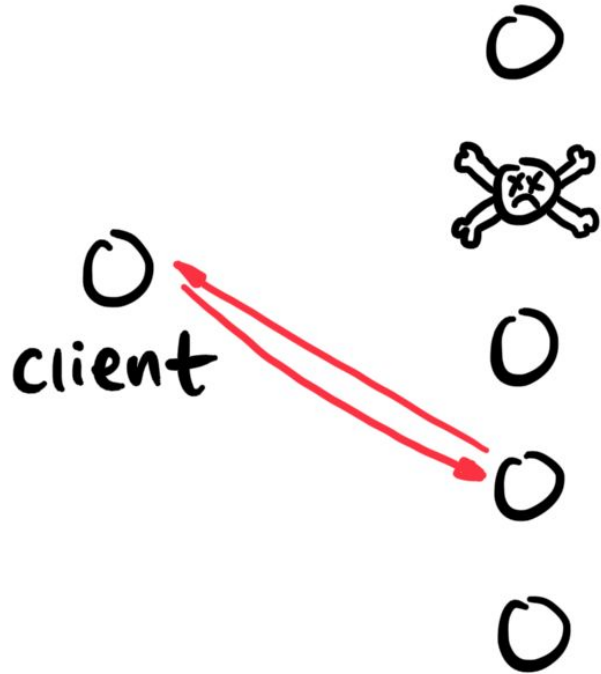
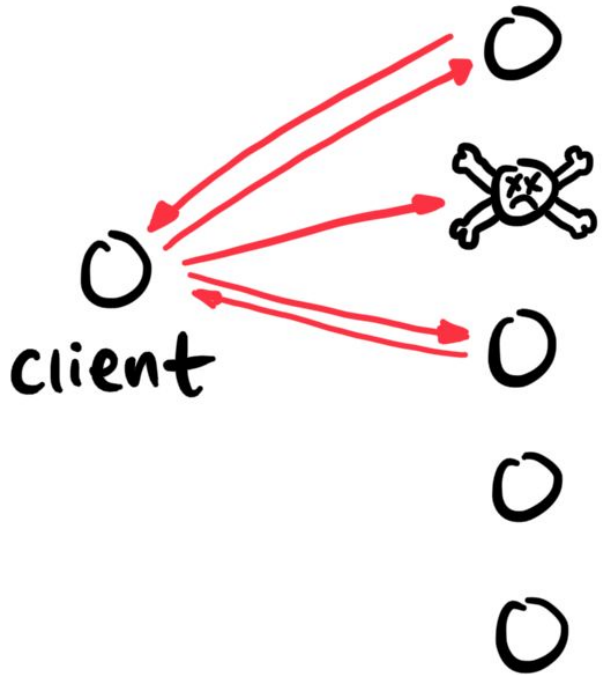
a = Node('a', latency=seconds(1))
b = Node('b', latency=seconds(2))
c = Node('c', latency=seconds(3))
d = Node('d', latency=seconds(4))
e = Node('e', latency=seconds(5))
f = Node('f', latency=seconds(6))
grid = QuorumSystem(reads=a*b*c + d*e*f)

grid.latency(read_fraction=0.5, optimize='latency') # 0:00:03.500000
grid.network_load(read_fraction=0.5, optimize='network') # 2.5
```

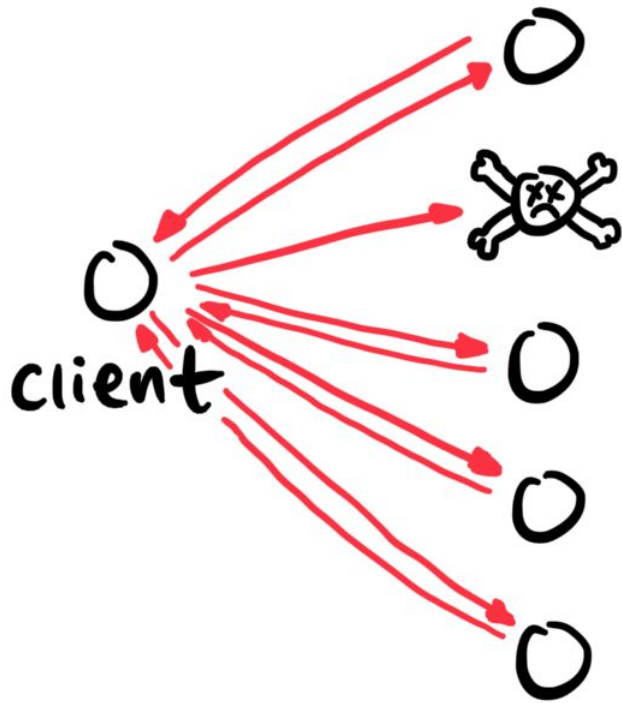
Cheap



Cheap



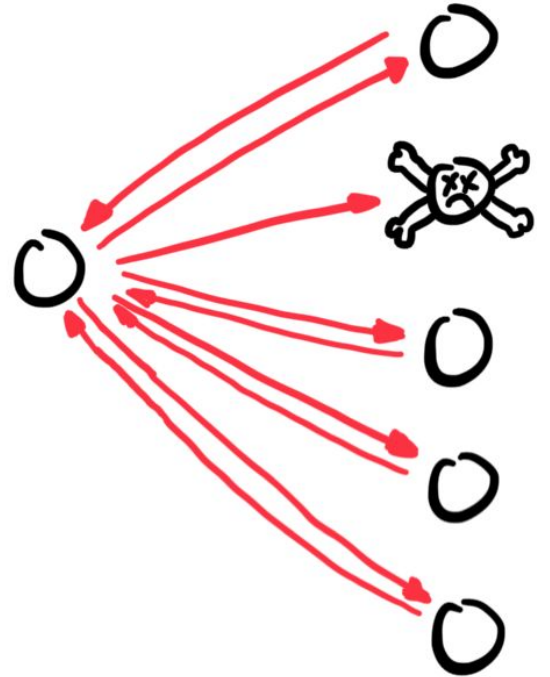
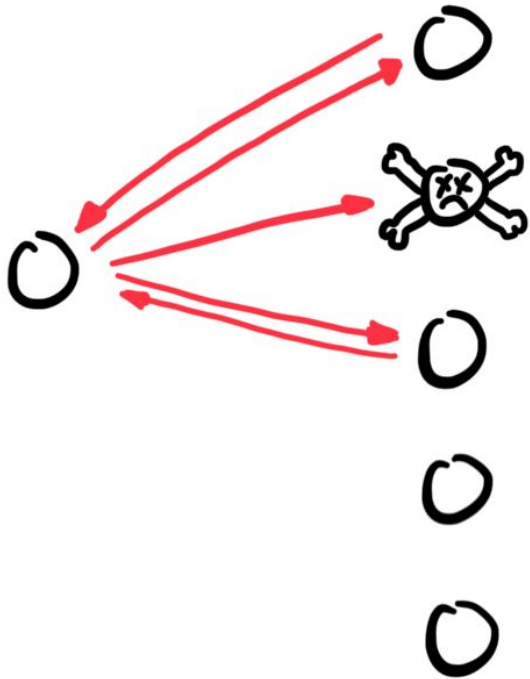
Not Cheap



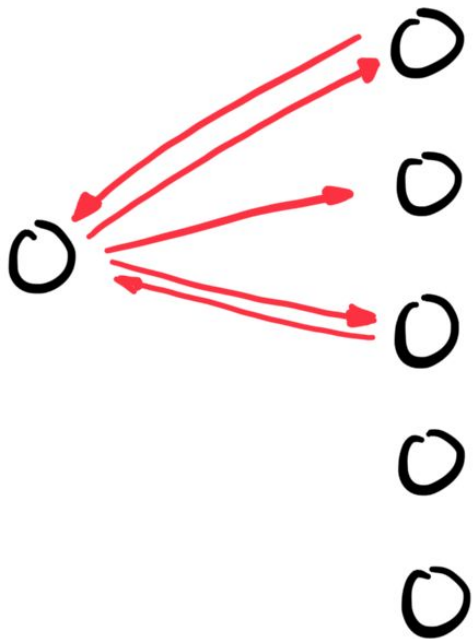
Cheap

???

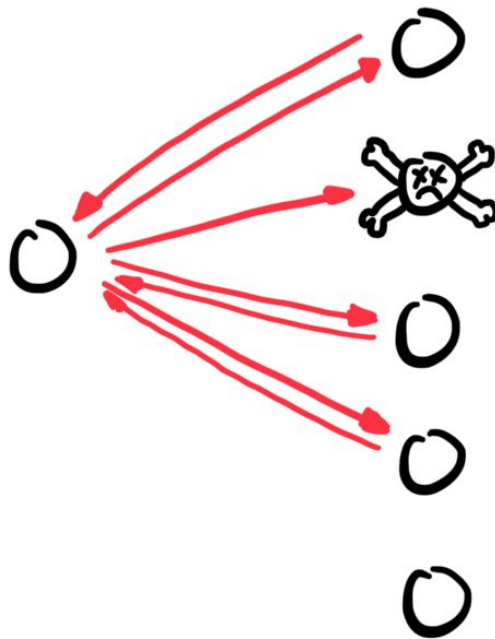
Not Cheap



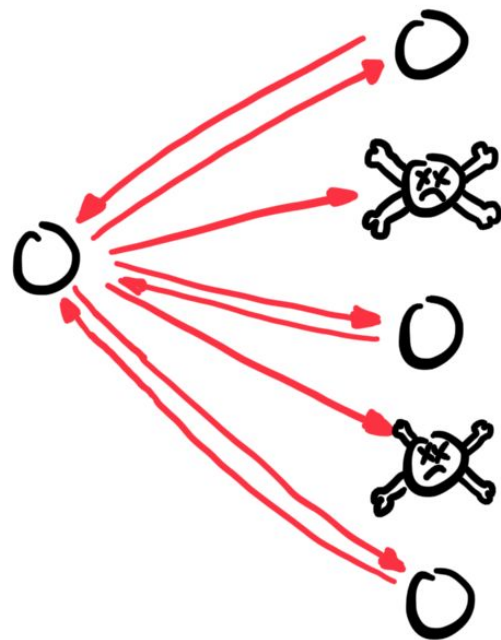
0-resilient



1-resilient



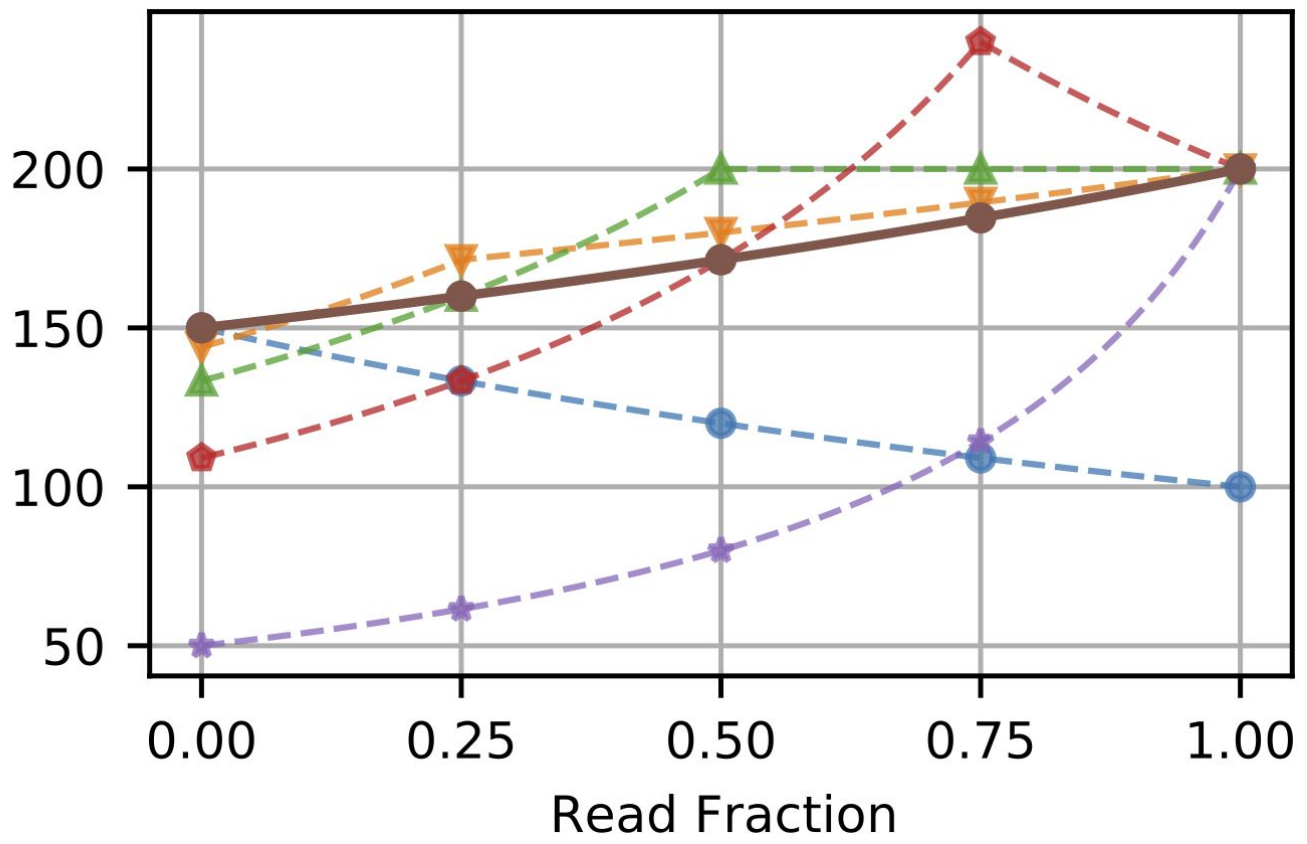
2-resilient

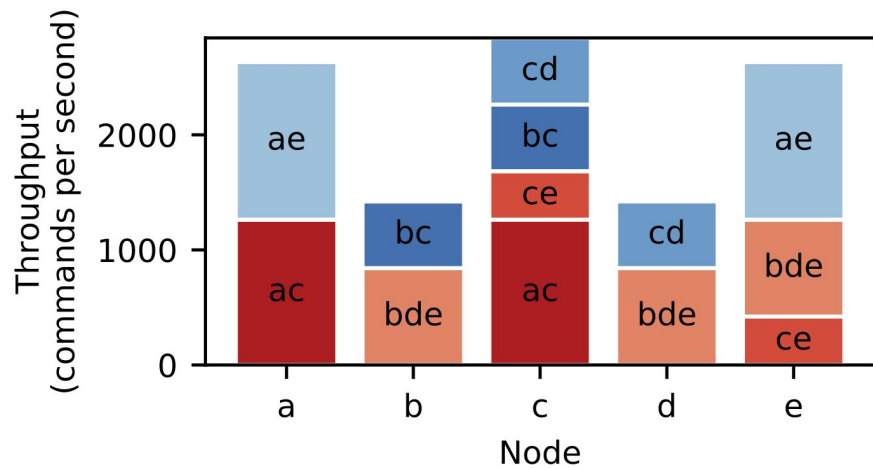
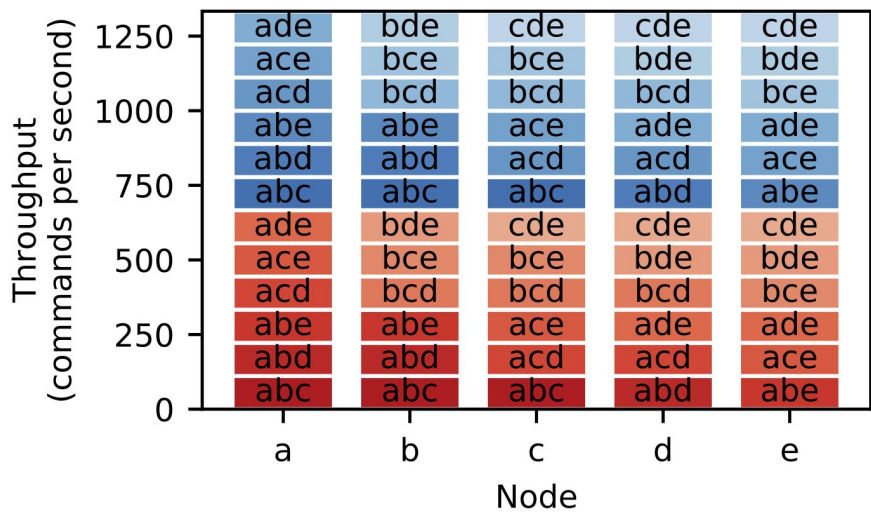


```
quorum_system.capacity(read_fraction=0, f=1)
```

```
quorum_system, strategy = search(  
    nodes=[a, b, c, d, e, f],  
    resilience=1,  
    f=1,  
    read_fraction=0.75,  
    optimize='load',  
    latency_limit=seconds(4),  
    network_limit=4,  
    timeout=seconds(60),  
)
```

Capacity (commands per second)





THANKS

github.com/mwhittaker/quoracle

```
pip install quoracle
```