

## Chapter 19

# Quorum Systems

What happens if a single server is no longer powerful enough to service all your customers? The obvious choice is to add more servers and to use the majority approach (e.g. Paxos, Chapter 15) to guarantee consistency. However, even if you buy one million servers, a client still has to access more than half of them per request! While you gain fault-tolerance, your efficiency can at most be doubled. Do we have to give up on consistency?

Let us take a step back: We used majorities because majority sets always overlap. But are majority sets the only sets that guarantee overlap? In this chapter we study the theory behind overlapping sets, known as quorum systems.

**Definition 19.1** (quorum, quorum system). *Let  $V = \{v_1, \dots, v_n\}$  be a set of nodes. A **quorum**  $Q \subseteq V$  is a subset of these nodes. A **quorum system**  $\mathcal{S} \subset 2^V$  is a set of quorums s.t. every two quorums intersect, i.e.,  $Q_1 \cap Q_2 \neq \emptyset$  for all  $Q_1, Q_2 \in \mathcal{S}$ .*

### Remarks:

- When a quorum system is being used, a client selects a quorum, acquires a lock (or ticket) on all nodes of the quorum, and when done releases all locks again. The idea is that no matter which quorum is chosen, its nodes will intersect with the nodes of every other quorum.
- What can happen if two quorums try to lock their nodes at the same time?
- A quorum system  $\mathcal{S}$  is called **minimal** if  $\forall Q_1, Q_2 \in \mathcal{S} : Q_1 \not\subseteq Q_2$ .
- The simplest quorum system imaginable consists of just one quorum, which in turn just consists of one server. It is known as **Singleton** (or primary copy).
- In the **Majority** quorum system, every quorum has  $\lfloor \frac{n}{2} \rfloor + 1$  nodes.
- Can you think of other simple quorum systems?

## 19.1 Load and Work

**Definition 19.2** (access strategy). An **access strategy**  $Z$  defines the probability  $P_Z(Q)$  of accessing a quorum  $Q \in \mathcal{S}$  s.t.  $\sum_{Q \in \mathcal{S}} P_Z(Q) = 1$ .

**Definition 19.3** (load).

- The **load** of access strategy  $Z$  on a node  $v_i$  is  $L_Z(v_i) = \sum_{Q \in \mathcal{S}; v_i \in Q} P_Z(Q)$ . The load is the probability that  $v_i \in Q$  if  $Q$  is sampled from  $\mathcal{S}$ .
- The **load** induced by access strategy  $Z$  on a quorum system  $\mathcal{S}$  is the maximal load induced by  $Z$  on any node in  $\mathcal{S}$ , i.e.,  $L_Z(\mathcal{S}) = \max_{v_i \in \mathcal{S}} L_Z(v_i)$ .
- The **load** of a quorum system  $\mathcal{S}$  is  $L(\mathcal{S}) = \min_Z L_Z(\mathcal{S})$ .

**Definition 19.4** (work).

- The **work** of a quorum  $Q \in \mathcal{S}$  is the number of nodes in  $Q$ ,  $W(Q) = |Q|$ .
- The **work** induced by access strategy  $Z$  on a quorum system  $\mathcal{S}$  is the expected number of nodes accessed, i.e.,  $W_Z(\mathcal{S}) = \sum_{Q \in \mathcal{S}} P_Z(Q) \cdot W(Q)$ .
- The **work** of a quorum system  $\mathcal{S}$  is  $W(\mathcal{S}) = \min_Z W_Z(\mathcal{S})$ .

**Remarks:**

- Note that you cannot choose different access strategies  $Z$  for work and load, you have to pick a single  $Z$  for both.
- Observe that for some access strategy  $Z$  the work  $W_Z(\mathcal{S})$  can be rewritten as  $\sum_{Q \in \mathcal{S}} \sum_{v \in Q} P_Z(Q)$ ; we can then swap summation order to get  $\sum_{v \in V} \sum_{Q \in \mathcal{S}; v \in Q} P_Z(Q) = \sum_{v \in V} L_Z(v)$ . In other words, the work induced by  $Z$  is the sum of the individual loads induced on the nodes.
- We illustrate the above concepts with a small example. Let  $V = \{v_1, v_2, v_3, v_4, v_5\}$  and  $\mathcal{S} = \{Q_1, Q_2, Q_3, Q_4\}$ , with  $Q_1 = \{v_1, v_2\}$ ,  $Q_2 = \{v_1, v_3, v_4\}$ ,  $Q_3 = \{v_2, v_3, v_5\}$ ,  $Q_4 = \{v_2, v_4, v_5\}$ . If we choose the access strategy  $Z$  s.t.  $P_Z(Q_1) = 1/2$  and  $P_Z(Q_2) = P_Z(Q_3) = P_Z(Q_4) = 1/6$ , then the node with the highest load is  $v_2$  with  $L_Z(v_2) = 1/2 + 1/6 + 1/6 = 5/6$ , i.e.,  $L_Z(\mathcal{S}) = 5/6$ . Regarding work, we have  $W_Z(\mathcal{S}) = 1/2 \cdot 2 + 1/6 \cdot 3 + 1/6 \cdot 3 + 1/6 \cdot 3 = 15/6$ .
- Can you come up with a better access strategy for  $\mathcal{S}$ ?
- If every quorum  $Q$  in a quorum system  $\mathcal{S}$  has the same number of elements,  $\mathcal{S}$  is called *uniform*.
- What is the minimum load a quorum system can have?

Primary Copy vs. Majority		Singleton	Majority
How many nodes need to be accessed?	(Work)	1	$\approx n/2$
What is the load of the busiest node?	(Load)	1	$\approx 1/2$

Table 19.5: First comparison of the Singleton and Majority quorum systems. Note that the Singleton quorum system can be a good choice when the failure probability of every single node is  $> 1/2$ .

**Theorem 19.6.** *Let  $\mathcal{S}$  be a quorum system. Then  $L(\mathcal{S}) \geq 1/\sqrt{n}$  holds.*

*Proof.* Let  $Q = \{v_1, \dots, v_q\}$  be a quorum of minimal size in  $\mathcal{S}$ , with  $|Q| = q$ . Let  $Z$  be an access strategy for  $\mathcal{S}$ . Every other quorum in  $\mathcal{S}$  intersects in at least one element with this quorum  $Q$ . Each time a quorum is accessed, at least one node in  $Q$  is accessed as well, yielding a lower bound of  $L_Z(v_i) \geq 1/q$  for some  $v_i \in Q$ .

Furthermore, as  $Q$  is minimal, at least  $q$  nodes need to be accessed, yielding  $W(\mathcal{S}) \geq q$ . Thus,  $W_Z(\mathcal{S}) = \sum_{v \in V} L_Z(v) \geq q$ , so  $L_Z(v_i) \geq q/n$  for some  $v_i \in V$ . In other words, as each time  $q$  nodes are accessed, the load of the most accessed node is at least  $q/n$ .

Combining both ideas leads to  $L_Z(\mathcal{S}) \geq \max(1/q, q/n) \Rightarrow L_Z(\mathcal{S}) \geq 1/\sqrt{n}$ . Thus,  $L(\mathcal{S}) \geq 1/\sqrt{n}$ , as  $Z$  can be *any* access strategy.  $\square$

**Remarks:**

- Can we achieve this load?

## 19.2 Grid Quorum Systems

**Definition 19.7** (Basic Grid quorum system). *Assume  $\sqrt{n} \in \mathbb{N}$ , and arrange the  $n$  nodes in a square matrix with side length of  $\sqrt{n}$ , i.e., in a grid. The basic **Grid** quorum system consists of  $\sqrt{n}$  quorums, with each containing the full row  $i$  and the full column  $i$ , for  $1 \leq i \leq \sqrt{n}$ .*

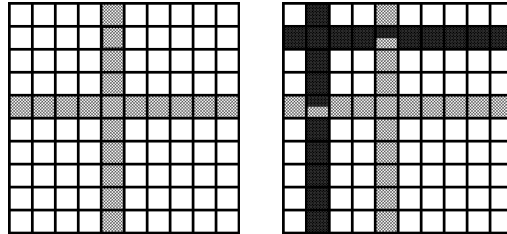


Figure 19.8: The basic version of the Grid quorum system, where each quorum  $Q_i$  with  $1 \leq i \leq \sqrt{n}$  uses row  $i$  and column  $i$ . The size of each quorum is  $2\sqrt{n} - 1$  and two quorums overlap in exactly two nodes. Thus, when the access strategy  $Z$  is uniform (i.e., the probability of each quorum is  $1/\sqrt{n}$ ), the work is  $2\sqrt{n} - 1$ , and the load of every node is in  $\Theta(1/\sqrt{n})$ .

**Remarks:**

- Consider the right picture in Figure 19.8: The two quorums intersect in two nodes. If both quorums were to be accessed at the same time, it is not guaranteed that at least one quorum will lock all of its nodes, as they could enter a deadlock!
- In the case of just two quorums, one could solve this by letting the quorums just intersect in one node, see Figure 19.9. However, already with three quorums the same situation could occur again, progress is not guaranteed!

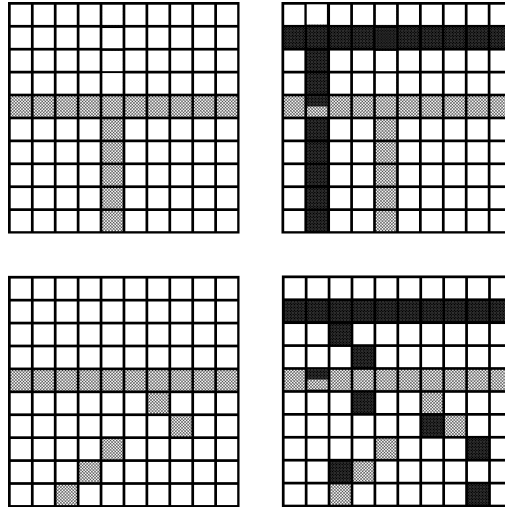


Figure 19.9: There are other ways to choose quorums in the grid s.t. pairwise different quorums only intersect in one node. The size of each quorum is between  $\sqrt{n}$  and  $2\sqrt{n} - 1$ , i.e., the work is in  $\Theta(\sqrt{n})$ . When the access strategy  $Z$  is uniform, the load of every node is in  $\Theta(1/\sqrt{n})$ .

- However, by deviating from the “access all at once” strategy, we can guarantee progress if the nodes are totally ordered!

---

**Algorithm 19.10** Sequential Locking Strategy for a Quorum  $Q$

---

- 1: Attempt to lock the nodes one by one, ordered by their identifiers
  - 2: Should a node be already locked, release all locks and start over
- 

**Theorem 19.11.** *If each quorum is accessed by Algorithm 19.10, at least one quorum will obtain a lock for all of its nodes.*

*Proof.* We prove the theorem by contradiction. Assume no quorum can make progress, i.e., for every quorum we have: At least one of its nodes is locked by another quorum. Let  $v$  be the node with the highest identifier that is locked by some quorum  $Q$ . Observe that  $Q$  already locked all of its nodes with a smaller identifier than  $v$ , otherwise  $Q$  would have restarted. As all nodes with a higher identifier than  $v$  are not locked,  $Q$  either has locked all of its nodes or can make progress – a contradiction. As the set of nodes is finite, one quorum will eventually be able to lock all of its nodes.  $\square$

**Remarks:**

- But now we are back to sequential accesses in a distributed system? Let’s do it concurrently with the same idea, i.e., resolving conflicts by the ordering of the nodes. Then, a quorum that locked the highest identifier so far can always make progress!

**Theorem 19.13.** *If the nodes and quorums use Algorithm 19.12, at least one quorum will obtain a lock for all of its nodes.*

**Algorithm 19.12** Concurrent Locking Strategy for a Quorum  $Q$ 

**Invariant:** Let  $v_Q \in Q$  be the highest identifier of a node locked by  $Q$  s.t. all nodes  $v_i \in Q$  with  $v_i < v_Q$  are locked by  $Q$  as well. Should  $Q$  not have any lock, then  $v_Q$  is set to 0.

---

```

1: repeat
2:   Attempt to lock all nodes of the quorum  $Q$ 
3:   for each node  $v \in Q$  that was not able to be locked by  $Q$  do
4:     exchange  $v_Q$  and  $v_{Q'}$  with the quorum  $Q'$  that locked  $v$ 
5:     if  $v_Q > v_{Q'}$  then
6:        $Q'$  releases lock on  $v$  and  $Q$  acquires lock on  $v$ 
7:     end if
8:   end for
9: until all nodes of the quorum  $Q$  are locked

```

---

*Proof.* The proof is analogous to the proof of Theorem 19.11: Assume for contradiction that no quorum can make progress. However, at least the quorum with the highest  $v_Q$  can always make progress – a contradiction! As the set of nodes is finite, at least one quorum will eventually be able to acquire a lock on all of its nodes.  $\square$

**Remarks:**

- What if a quorum locks all of its nodes and then crashes? Is the quorum system dead now? This issue can be prevented by, e.g., using leases instead of locks: leases have a timeout, i.e., a lock is released eventually. But what happens if a quorum is slow and its acquired leases expire before it can acquire all leases?

## 19.3 Fault Tolerance

**Definition 19.14** (resilience). *If any  $f$  nodes from a quorum system  $\mathcal{S}$  can fail s.t. there is still a quorum  $Q \in \mathcal{S}$  without failed nodes, then  $\mathcal{S}$  is  $f$ -resilient. The largest such  $f$  is the **resilience**  $R(\mathcal{S})$ .*

**Theorem 19.15.** *Let  $\mathcal{S}$  be a Grid quorum system where each of the  $n$  quorums consists of a full row and a full column.  $\mathcal{S}$  has a resilience of  $\sqrt{n} - 1$ .*

*Proof.* If all  $\sqrt{n}$  nodes on the diagonal of the grid fail, then every quorum will have at least one failed node. Should less than  $\sqrt{n}$  nodes fail, then there is a row and a column without failed nodes.  $\square$

**Remarks:**

- The Grid quorum system in Theorem 19.15 is different from the Basic Grid quorum system described in Definition 19.7. In each quorum in the Basic Grid quorum system the row and column index are identical, while in the Grid quorum system of Theorem 19.15 this is not the case.

**Definition 19.16** (failure probability). *Assume that every node works with a fixed probability  $p$  (in the following we assume concrete values, e.g.  $p > 1/2$ )*

or  $p \geq 2/3$ ). The **failure probability**  $F_p(\mathcal{S})$  of a quorum system  $\mathcal{S}$  is the probability that at least one node of every quorum fails.

**Remarks:**

- The **asymptotic failure probability** is  $F_p(\mathcal{S})$  for  $n \rightarrow \infty$ .

**Facts 19.17.** A version of a **Chernoff bound** states the following:

Let  $x_1, \dots, x_n$  be independent Bernoulli-distributed random variables with  $Pr[x_i = 1] = p_i$  and  $Pr[x_i = 0] = 1 - p_i = q_i$ , then for  $X := \sum_{i=1}^n x_i$  and  $\mu := \mathbb{E}[X] = \sum_{i=1}^n p_i$  the following holds:

$$\text{for all } 0 < \delta < 1: Pr[X \leq (1 - \delta)\mu] \leq e^{-\mu\delta^2/2}.$$

**Theorem 19.18.** The asymptotic failure probability of the Majority quorum system is 0, for  $p > 1/2$ .

*Proof.* In a Majority quorum system each quorum contains exactly  $\lfloor \frac{n}{2} \rfloor + 1$  nodes and each subset of nodes with cardinality  $\lfloor \frac{n}{2} \rfloor + 1$  forms a quorum. If only  $\lfloor \frac{n}{2} \rfloor$  nodes work, then the Majority quorum system fails. Otherwise there is at least one quorum available. In order to calculate the failure probability we define the following random variables:

$$x_i = \begin{cases} 1, & \text{if node } i \text{ works, happens with probability } p \\ 0, & \text{if node } i \text{ fails, happens with probability } q = 1 - p \end{cases}$$

and  $X := \sum_{i=1}^n x_i$ , with  $\mu = np$ , whereas  $X$  corresponds to the number of working nodes. To estimate the probability that the number of working nodes is less than  $\lfloor \frac{n}{2} \rfloor + 1$  we will make use of the Chernoff inequality from above. By setting  $\delta = 1 - \frac{1}{2p}$  we obtain  $F_P(\mathcal{S}) = Pr[X \leq \lfloor \frac{n}{2} \rfloor] \leq Pr[X \leq \frac{n}{2}] = Pr[X \leq (1 - \delta)\mu]$ .

With  $\delta = 1 - \frac{1}{2p}$  we have  $0 < \delta \leq 1/2$  due to  $1/2 < p \leq 1$ . Thus, we can use the Chernoff bound and get  $F_P(\mathcal{S}) \leq e^{-\mu\delta^2/2} \in e^{-\Omega(n)}$ .  $\square$

**Theorem 19.19.** The asymptotic failure probability of the Grid quorum system is 1 for  $p > 0$ .

*Proof.* Consider the  $n = d \cdot d$  nodes to be arranged in a  $d \times d$  grid. A quorum always contains one full row. In this estimation we will make use of the Bernoulli inequality which states that for all  $n \in \mathbb{N}, x \geq -1 : (1 + x)^n \geq 1 + nx$ .

The system fails, if in each row at least one node fails (which happens with probability  $1 - p^d$  for a particular row, as all nodes work with probability  $p^d$ ). Therefore we can bound the failure probability from below with:

$$F_p(\mathcal{S}) \geq Pr[\text{at least one failure per row}] = (1 - p^d)^d \geq 1 - dp^d \xrightarrow[n \rightarrow \infty]{} 1. \quad \square$$

**Remarks:**

- Now we have a quorum system with optimal load (the Grid) and one with fault-tolerance (Majority), but what if we want both?

**Definition 19.20** (B-Grid quorum system). Consider  $n = dhr$  nodes, arranged in a rectangular grid with  $h \cdot r$  rows and  $d$  columns. Each group of  $r$  rows is a band, and  $r$  elements in a column restricted to a band are called a mini-column. A quorum consists of one mini-column in every band and one element from each mini-column of one band; thus every quorum has  $d + hr - 1$  elements. The **B-Grid** quorum system consists of all such quorums.

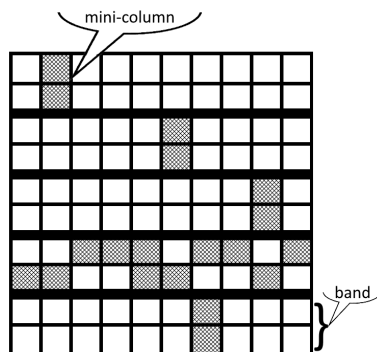


Figure 19.21: A B-Grid quorum system with  $n = 100$  nodes,  $d = 10$  columns,  $h \cdot r = 10$  rows,  $h = 5$  bands, and  $r = 2$ . The depicted quorum has a  $d + hr - 1 = 10 + 5 \cdot 2 - 1 = 19$  nodes. If the access strategy  $Z$  is chosen uniformly, then we have a work of  $d + hr - 1$  and a load of  $\frac{d+hr-1}{n}$ . By setting  $d = \sqrt{n}$  and  $r = \ln d$ , we obtain a work of  $\Theta(\sqrt{n})$  and a load of  $\Theta(1/\sqrt{n})$ .

**Theorem 19.22.** *The asymptotic failure probability of the B-Grid quorum system is 0, for  $p \geq \frac{2}{3}$ .*

*Proof.* Suppose  $n = dhr$  and the elements are arranged in a grid with  $d$  columns and  $h \cdot r$  rows. The B-Grid quorum system does fail if in each band a complete mini-column fails, because then it is not possible to choose a band where in each mini-column an element is still working. It also fails if in a band an element in each mini-column fails. If none of those cases holds, then the B-Grid system does not fail. Those events may not be independent of each other, but with the help of the union bound, we can upper bound the failure probability with the following equation:

$$\begin{aligned} F_p(\mathcal{S}) &\leq Pr[\text{in every band a complete mini-column fails}] \\ &\quad + Pr[\text{in a band at least one element of every m.-col. fails}] \\ &\leq (d(1-p)^r)^h + h(1-p^r)^d \end{aligned}$$

We use  $d = \sqrt{n}$ ,  $r = \ln d$ , and  $0 \leq 1 - p \leq 1/3$ . Using  $n^{\ln x} = x^{\ln n}$ , we have  $d(1-p)^r \leq d \cdot d^{\ln 1/3} \approx d^{-0.1}$ , and hence for large enough  $d$  the whole first term is bounded from above by  $d^{-0.1h} \ll 1/d^2 = 1/n$ .

Regarding the second term, we have  $p \geq 2/3$ , and  $h = d/\ln d < d$ . Hence we can bound the term from above by  $d(1 - d^{\ln 2/3})^d \approx d(1 - d^{-0.4})^d$ . Using  $(1 + t/n)^n \leq e^t$ , we get (again, for large enough  $d$ ) an upper bound of  $d(1 - d^{-0.4})^d = d(1 - d^{0.6}/d)^d \leq d \cdot e^{-d^{0.6}} = d^{(-d^{0.6}/\ln d)+1} \ll d^{-2} = 1/n$ . In total, we have  $F_p(\mathcal{S}) \in O(1/n)$ .  $\square$

## 19.4 Byzantine Quorum Systems

While failed nodes are bad, they are still easy to deal with: just access another quorum where all nodes can respond! Byzantine nodes make life more difficult however, as they can pretend to be a regular node, i.e., one needs more sophisticated methods to deal with them. We need to ensure that the intersection

	Singleton	Majority	Grid	B-Grid*
Work	<b>1</b>	$\approx n/2$	$\Theta(\sqrt{n})$	$\Theta(\sqrt{n})$
Load	1	$\approx 1/2$	$\Theta(\mathbf{1}/\sqrt{\mathbf{n}})$	$\Theta(\mathbf{1}/\sqrt{\mathbf{n}})$
Resilience	0	$\approx \mathbf{n}/\mathbf{2}$	$\Theta(\sqrt{n})$	$\Theta(\sqrt{n})$
F. Prob.**	$1 - p$	$\rightarrow \mathbf{0}$	$\rightarrow 1$	$\rightarrow \mathbf{0}$

Table 19.23: Overview of the different quorum systems regarding resilience, work, load, and their asymptotic failure probability. The best entries in each row are set in bold.

\* Setting  $d = \sqrt{n}$  and  $r = \ln d$ .

\*\* Assuming prob.  $q = 1 - p$  is constant but significantly less than  $1/2$ .

of two quorums always contains a non-byzantine (correct) node and furthermore, the byzantine nodes should not be allowed to infiltrate every quorum. In this section we study three counter-measures of increasing strength, and their implications on the load of quorum systems.

**Definition 19.24** (*f*-disseminating). A quorum system  $\mathcal{S}$  is *f-disseminating* if (1) the intersection of two different quorums always contains  $f + 1$  nodes, and (2) for any set of  $f$  byzantine nodes, there is at least one quorum without byzantine nodes.

**Remarks:**

- Thanks to (2), even with  $f$  byzantine nodes, the byzantine nodes cannot stop all quorums by just pretending to have crashed. At least one quorum will survive. We will also keep this assumption for the upcoming more advanced byzantine quorum systems.
- Byzantine nodes can also do something worse than crashing - they could falsify data! Nonetheless, due to (1), there is at least one non-byzantine node in every quorum intersection. If the data is self-verifying by, e.g., authentication, then this one node is enough.
- If the data is not self-verifying, then we need another mechanism.

**Definition 19.25** (*f*-masking). A quorum system  $\mathcal{S}$  is *f-masking* if (1) the intersection of two different quorums always contains  $2f + 1$  nodes, and (2) for any set of  $f$  byzantine nodes, there is at least one quorum without byzantine nodes.

**Remarks:**

- Note that except for the second condition, an *f*-masking quorum system is the same as a  $2f$ -disseminating system. The idea is that the non-byzantine nodes (at least  $f + 1$ ) can outvote the byzantine ones (at most  $f$ ), but only if all non-byzantine nodes are up-to-date!
- This raises an issue not covered yet in this chapter. If we access some quorum and update its values, this change still has to be disseminated to the other nodes in the byzantine quorum system. Opaque quorum systems deal with this issue, which are discussed at the end of this section.



- One can show that  $f$ -disseminating quorum systems need more than  $3f$  nodes and  $f$ -masking quorum systems need more than  $4f$  nodes. In other words,  $f < n/3$ , or  $f < n/4$ . Essentially, the quorums may not contain too many nodes, and the different intersection properties lead to the different bounds.

**Theorem 19.26.** *Let  $\mathcal{S}$  be an  $f$ -disseminating quorum system. Then  $L(\mathcal{S}) \geq \sqrt{(f+1)/n}$  holds.*

**Theorem 19.27.** *Let  $\mathcal{S}$  be an  $f$ -masking quorum system. Then  $L(\mathcal{S}) \geq \sqrt{(2f+1)/n}$  holds.*

*Proofs of Theorems 19.26 and 19.27.* The proofs follow the proof of Theorem 19.6, by observing that now not just one element is accessed from a minimal quorum, but  $f+1$  or  $2f+1$ , respectively.  $\square$

**Definition 19.28** ( $f$ -masking Grid quorum system). *A  $f$ -masking Grid quorum system is constructed as the grid quorum system, but each quorum contains one full column and  $f+1$  rows of nodes, with  $2f+1 \leq \sqrt{n}$ .*

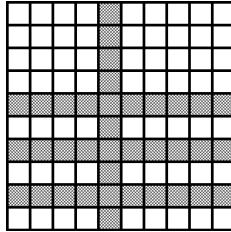


Figure 19.29: An example how to choose a quorum in the  $f$ -masking Grid with  $f = 2$ , i.e.,  $2 + 1 = 3$  rows. The load is in  $\Theta(f/\sqrt{n})$  when the access strategy is chosen to be uniform. Two quorums overlap by their columns intersecting each other's rows, i.e., they overlap in at least  $2f + 2$  nodes.

**Remarks:**

- The  $f$ -masking Grid nearly hits the lower bound for the load of  $f$ -masking quorum systems, but not quite. A small change and we will be optimal asymptotically.

**Definition 19.30** (*M-Grid quorum system*). *The M-Grid quorum system is constructed as the grid quorum as well, but each quorum contains  $\sqrt{f+1}$  rows and  $\sqrt{f+1}$  columns of nodes, with  $2f+1 \leq \sqrt{n}$ .*

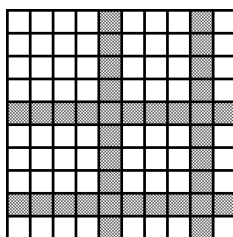


Figure 19.31: An example how to choose a quorum in the  $M$ -Grid with  $f = 3$ , i.e., 2 rows and 2 columns. The load is in  $\Theta(\sqrt{f/n})$  when the access strategy is chosen to be uniform. Two quorums overlap with each row intersecting each other's column, i.e.,  $2\sqrt{f+1}^2 = 2f+2$  nodes.

**Corollary 19.32.** *The  $f$ -masking Grid quorum system and the  $M$ -Grid quorum system are  $f$ -masking quorum systems.*

**Remarks:**

- We achieved nearly the same load as without byzantine nodes! However, as mentioned earlier, what happens if we access a quorum that is not up-to-date, except for the intersection with an up-to-date quorum? Surely we can fix that as well without too much loss?
- This property will be handled in the last part of this chapter by *opaque* quorum systems. It will ensure that the number of correct up-to-date nodes accessed will be larger than the number of out-of-date nodes combined with the byzantine nodes in the quorum (cf. (19.33.1)).

**Definition 19.33** ( *$f$ -opaque quorum system*). *A quorum system  $\mathcal{S}$  is  $f$ -opaque if the following two properties hold for any set of  $f$  byzantine nodes  $F$  and any two different quorums  $Q_1, Q_2$ :*

$$|(Q_1 \cap Q_2) \setminus F| > |(Q_2 \cap F) \cup (Q_2 \setminus Q_1)| \quad (19.33.1)$$

$$F \cap Q = \emptyset \text{ for some } Q \in \mathcal{S} \quad (19.33.2)$$

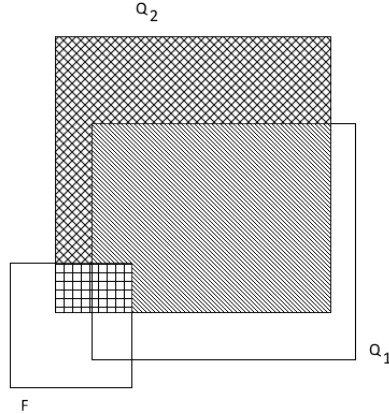


Figure 19.34: Intersection properties of an opaque quorum system. Equation (19.33.1) ensures that the set of non-byzantine nodes in the intersection of  $Q_1, Q_2$  is larger than the set of out of date nodes, even if the byzantine nodes “team up” with those nodes. Thus, the correct up to date value can always be recognized by a majority voting.

**Remarks:**

- For any  $f$ -opaque quorum system, inequality (19.33.1) also holds for  $|F| < f$ . In particular, substituting  $F = \emptyset$  in (19.33.1) gives  $|Q_1 \cap Q_2| > |Q_2 \setminus Q_1|$ ; similarly, one can also deduce that  $|Q_1 \cap Q_2| > |Q_1 \setminus Q_2|$ . Therefore,  $|Q_1| = |Q_1 \setminus Q_2| + |Q_1 \cap Q_2| < 2|Q_1 \cap Q_2|$ , so  $|Q_1 \cap Q_2| > \frac{|Q_1|}{2}$ .

**Theorem 19.35.** *Let  $\mathcal{S}$  be an  $f$ -opaque quorum system. Then,  $f < n/5$ .*

*Proof.* Due to (19.33.2), there exists a quorum  $Q_1$  with size at most  $n - f$ . With (19.33.1),  $|Q_1| > f$  holds. Let  $F_1$  be a set of  $f$  (byzantine) nodes  $F_1 \subset Q_1$ , and with (19.33.2), there exists a  $Q_2 \subseteq V \setminus F_1$ . Thus,  $|Q_1 \cap Q_2| \leq n - 2f$ . With (19.33.1),  $|Q_1 \cap Q_2| > f$  holds. Thus, one could choose  $f$  (byzantine) nodes  $F_2$  with  $F_2 \subset Q_1 \cap Q_2$ . Using (19.33.1) one can bound  $n - 3f$  from below:  $n - 3f \geq |Q_2 \cap Q_1| - |F_2| = |(Q_2 \cap Q_1) \setminus F_2| > |(Q_1 \cap F_2) \cup (Q_1 \setminus Q_2)| = |F_2 \cup (Q_1 \setminus Q_2)| = |F_2| + |Q_1 \setminus Q_2| \geq |F_2| + |F_1| = 2f$ .  $\square$

**Remarks:**

- One can extend the Majority quorum system to be  $f$ -opaque by setting the size of each quorum to contain  $\lceil (2n + 2f)/3 \rceil$  nodes. Then its load is  $1/n \lceil (2n + 2f)/3 \rceil \geq 2/3 + 2f/3n \geq 2/3$ .
- Can we do much better? Sadly, no...

**Theorem 19.36.** *Let  $\mathcal{S}$  be an  $f$ -opaque quorum system. Then  $L(\mathcal{S}) > 1/2$  holds.*

*Proof.* Equation (19.33.1) implies that for  $Q_1, Q_2 \in \mathcal{S}$ , the intersection of  $Q_1$  and  $Q_2$  is more than half their size, i.e.,  $|Q_1 \cap Q_2| > |Q_1|/2$ . Assuming  $\mathcal{S} =$

$\{Q_1, Q_2, \dots\}$ , the total load induced by an access strategy  $Z$  on nodes in  $Q_1$  is:

$$\sum_{v \in Q_1} \sum_{i: v \in Q_i} P_Z(Q_i) = \sum_i \sum_{v \in Q_1 \cap Q_i} P_Z(Q_i) = \sum_i P_Z(Q_i) |Q_1 \cap Q_i| > \frac{|Q_1|}{2}.$$

Using the pigeonhole principle, there must be at least one node in  $Q_1$  with load greater than  $1/2$ .  $\square$

## Chapter Notes

Historically, a quorum is the minimum number of members of a deliberative body necessary to conduct the business of that group. Their use has inspired the introduction of quorum systems in computer science since the late 1970s/early 1980s. Early work focused on Majority quorum systems [Lam78, Gif79, Tho79], with the notion of minimality introduced shortly after [GB85]. The Grid quorum system was first considered in [Mae85], with the B-Grid being introduced in [NW94]. The latter article and [PW95] also initiated the study of load and resilience.

The  $f$ -masking Grid quorum system and opaque quorum systems are from [MR98], and the  $M$ -Grid quorum system was introduced in [MRW97]. Both papers also mark the start of the formal study of Byzantine quorum systems. The  $f$ -masking and the  $M$ -Grid have asymptotic failure probabilities of 1, more complex systems with better values can be found in these papers as well.

Quorum systems have also been extended to cope with nodes dynamically leaving and joining, see, e.g., the dynamic paths quorum system in [NW05].

For a further overview on quorum systems, we refer to the book by Vukolić [Vuk12] and the article by Merideth and Reiter [MR10].

This chapter was written in collaboration with Klaus-Tycho Förster.

## Bibliography

- [GB85] Hector Garcia-Molina and Daniel Barbará. How to assign votes in a distributed system. *J. ACM*, 32(4):841–860, 1985.
- [Gif79] David K. Gifford. Weighted voting for replicated data. In Michael D. Schroeder and Anita K. Jones, editors, *Proceedings of the Seventh Symposium on Operating System Principles, SOSP 1979, Asilomar Conference Grounds, Pacific Grove, California, USA, 10-12, December 1979*, pages 150–162. ACM, 1979.
- [Lam78] Leslie Lamport. The implementation of reliable distributed multiprocess systems. *Computer Networks*, 2:95–114, 1978.
- [Mae85] Mamoru Maekawa. A square root  $N$  algorithm for mutual exclusion in decentralized systems. *ACM Trans. Comput. Syst.*, 3(2):145–159, 1985.
- [MR98] Dahlia Malkhi and Michael K. Reiter. Byzantine quorum systems. *Distributed Computing*, 11(4):203–213, 1998.

- [MR10] Michael G. Merideth and Michael K. Reiter. Selected results from the latest decade of quorum systems research. In Bernadette Charron-Bost, Fernando Pedone, and André Schiper, editors, *Replication: Theory and Practice*, volume 5959 of *Lecture Notes in Computer Science*, pages 185–206. Springer, 2010.
- [MRW97] Dahlia Malkhi, Michael K. Reiter, and Avishai Wool. The load and availability of byzantine quorum systems. In James E. Burns and Hagit Attiya, editors, *Proceedings of the Sixteenth Annual ACM Symposium on Principles of Distributed Computing, Santa Barbara, California, USA, August 21-24, 1997*, pages 249–257. ACM, 1997.
- [NW94] Moni Naor and Avishai Wool. The load, capacity and availability of quorum systems. In *35th Annual Symposium on Foundations of Computer Science, Santa Fe, New Mexico, USA, 20-22 November 1994*, pages 214–225. IEEE Computer Society, 1994.
- [NW05] Moni Naor and Udi Wieder. Scalable and dynamic quorum systems. *Distributed Computing*, 17(4):311–322, 2005.
- [PW95] David Peleg and Avishai Wool. The availability of quorum systems. *Inf. Comput.*, 123(2):210–223, 1995.
- [Tho79] Robert H. Thomas. A majority consensus approach to concurrency control for multiple copy databases. *ACM Trans. Database Syst.*, 4(2):180–209, 1979.
- [Vuk12] Marko Vukolic. *Quorum Systems: With Applications to Storage and Consensus*. Synthesis Lectures on Distributed Computing Theory. Morgan & Claypool Publishers, 2012.