

Adaptive Wide Area REplication

geographically as

AWARE

→ automated and dynamic voting-weight tuning and partitioning scheme
↳ deterministic, self-optimisation algorithm

↳ dynamically optimise consensus latency by self-reliantly finding a fast configuration → latency gains

latency optimisation scheme

efficient, adaptive Byzantine consensus

WHEAT + BFT-SMART

underlying weighting scheme

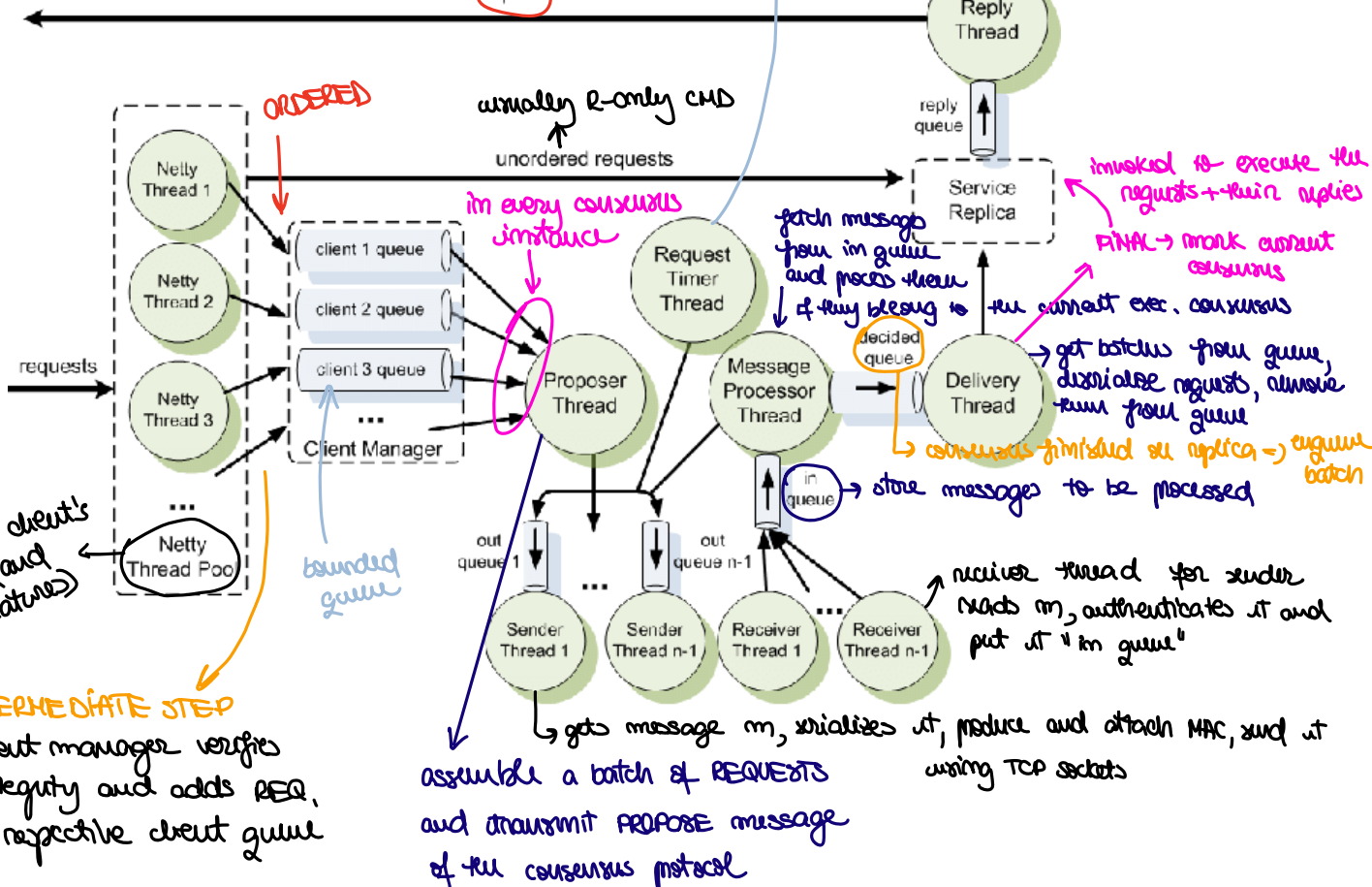
replication protocol

1. BFT-SMART

client waits for specific quorum of matching replies

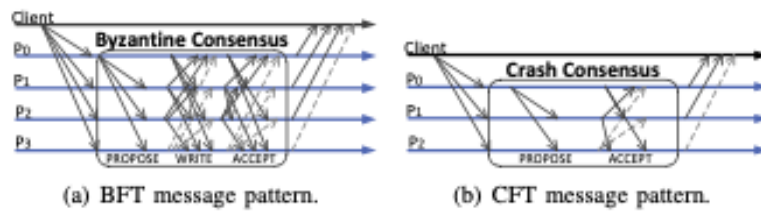
replies

periodically activated to verify if issue request remained more than a pre-defined time on the PENDING REQUESTS QUEUE



⇒ TOTAL-ORDER execution of requests

⇒ employs dynamically scalable replica set + provides modular architecture (use separated and exchangeable components)

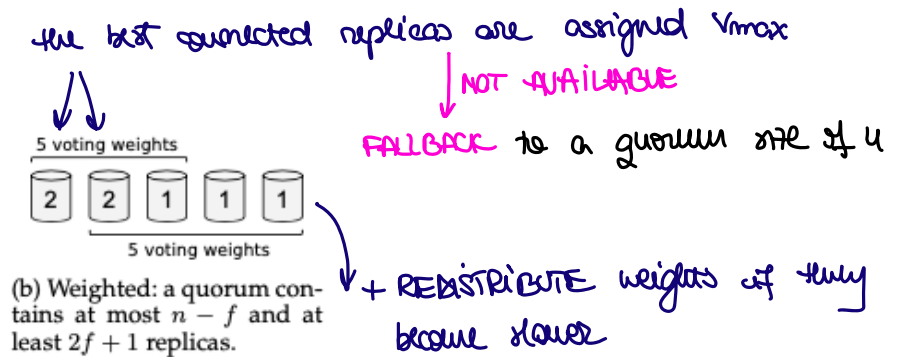
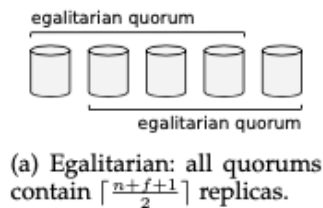


Byzantine consensus algorithm

- PROPOSE → leader broadcasts a message that contains a batch of requests that need to be decided to all other replicas
 - WRITE
 - ACCEPT
- ALL-TO-ALL broadcast used for commitment → each replica i waits for a quorum Q_i ; $\lceil \frac{n+f+1}{2} \rceil$ replicas to proceed
- $i \neq j \Rightarrow |Q_i \cap Q_j| \geq f+1$
- overlap with at least 2 correct replicas

2. WEIGHT

→ add more replicas to a system



SAFE WEIGHT DISTRIBUTION SCHEME

→ m replicas tolerating f Byzantine faults and containing D additional replicas

$$m = 3f + D + 1$$

$Q_v = 2(f+D) + 1 \rightarrow$ wait for quorum formation

$$V_{min} = 1$$

$V_{max} = 1 + \frac{D}{f} \rightarrow$ attributed to $2f$ best-connected replicas

SELF-MONITORING

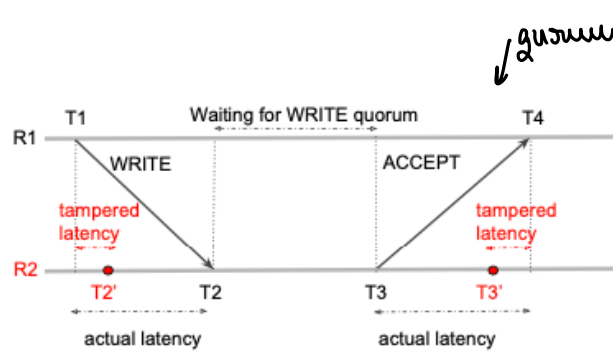


Figure 5: Problem with timestamps and Byzantine replicas.

quorum-based measurements \Rightarrow DON'T ALLOW measuring about cmt latencies

NON-MALICIOUS string \Rightarrow piggyback responses carrying TIMESTAMPS

\Rightarrow FAVOR ONE-SIDED measurements \rightarrow require only the measuring replica to be correct

\Downarrow

WRITE AND ACCEPT \rightarrow replicas immediately respond with WRITE-RESPONSE after receiving a WRITE
 \downarrow
 introduces monitoring overhead

AWARE

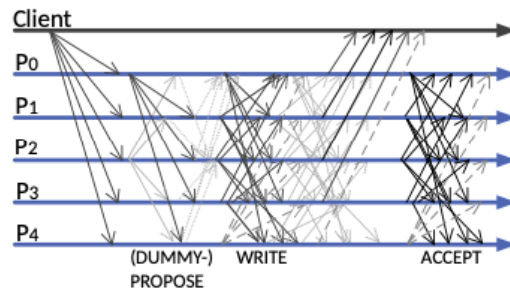


Figure 6: Message flow of AWARE ($f = 1; \Delta = 1$).

highest accuracy in leader election

\rightarrow MONITORING WINDOW \rightarrow # recent monitoring messages to be used for latency comp.

Response to WRITE \rightarrow send WRITE-RESPONSE messages (which include a challenge)

$\hookrightarrow L_i = (e_{i,0}, \dots, e_{i,m-1})$ latency vector of replica $i \rightarrow$ report after SYNCHRONIZATION PERIOD

DUMMY-PROPOSE \rightarrow use ROTATION SCHEME & one additional replica simultaneously sends a proposal for a dummy batch, but without starting a new consensus instance and all replicas disregard the proposal

\downarrow
OPTIONAL

\rightarrow measure the time new-leaders need to PURPOSE

Bounding monitoring overhead \rightarrow specify parameter $w \in [0, 1]$ determines the maximum overhead induced by the monitoring procedure

CALCULATION INTERVAL $c \rightarrow$ # consensus instances after which a calculation and possibly a reconfiguration is triggered

All replicas maintain SYNCHRONISED LATENCY MATRICES M^p, M^w

1. INITIALIZATION

$$M^x[i, j] = \begin{cases} \infty, & i=j \\ 0, & \text{otherwise} \end{cases}$$

← latency of replica i to j
measured by i for message type x

2. UPDATE using BFT-SMART invoke

$$\text{invokeOrdered}(\text{MEASURE}, L_i^p, L_i^w)$$

⇒ global total-order on all client requests and measurement messages

3. SANITIZE matrices

$$\hat{M}^x[i, j] = \max(M^x[i, j], M^x[j, i])$$

SELF-OPTIMIZATION

↳ replicas deterministically reconfigure to a new weight configuration and/or leader position

1. VOTING WEIGHTS TUNING

↳ search for weight distribution that optimises the system's consensus latency

2. LEADER RELOCATION

→ allow protocol to provide a suitable set of 2 candidates → AWARE chooses one

→ optimization goal α → threshold by which a predicted consensus must be faster to trigger reconfiguration

↓
evaluated at most once every
calculation interval

MONITORING

1. Each replica i collects its latency measurements (a moving median) in a vector $L_i = \langle l_{i,0}, \dots, l_{i,n-1} \rangle$;
2. Periodically, each replica i disseminates its vectors for $PROPOSE$ (L_i^P) and $WRITE$ (L_i^W) with total order by calling $invokeOrdered(MEASURE, L_i^P, L_i^W)$;
3. Once a replica i decides a batch, that batch may contain messages $\langle MEASURE, L_j^P, L_j^W \rangle$ from some replica $j \in I$. It uses these vectors to update its synchronized matrices M_i^P and M_i^W , i.e., for replicas $k = 0, \dots, n-1$ assigning $M_i^W[j, k] = L_j^W[k]$ that is the information that i has about the latency between replica j and all other replicas measured by j . This applies to the maintained latency information for both $PROPOSE$ (M_i^P) and $WRITE$ (M_i^W).
4. When a defined number (specified by the calculation interval c) of consensus instances is reached, all replicas have the same matrices M^P and M^W , e.g., with stopping condition
4. The calculation interval determines the frequency of optimizations and can be configured in our implementation. Our default value is 500.

Algorithm 1: *formQV* computes the times replicas form weighted quorums of $Q_v = 2fV_{max} + 1$ voting weights.

Data: replica set I , latency matrix \hat{M}^W , times $(T_i^{current})_{i \in I}$ and voting weights $(V_i)_{i \in I}$
Result: times $(T_i^{next})_{i \in I}$ replicas can advance to the next protocol stage

```

1 for  $i \in I$  do
2    $received_i \leftarrow$  new PriorityQueue()
3   for  $j \in I$  do
4      $received_i.add(\langle T_j^{current} + \hat{M}^W[j, i], V_j \rangle)$ 
5   for  $i \in I$  do
6      $weight \leftarrow 0$ 
7     while  $weight < Q_v$  do
8        $\langle T_{next}, V_{next} \rangle \leftarrow received_i.dequeue()$ 
9        $weight \leftarrow weight + V_{next}$ 
10       $T_i^{next} \leftarrow T_{next}$ 
11 return  $(T_i^{next})_{i \in I}$ 

```

Handwritten notes:
 → compute the time the replica receives the write from each others
 → for every replica
 → compute the time that gathered enough voting weights

Algorithm 2: *PredictLatency* computes the consensus latency (amortized over multiple rounds)

Data: replica set I , leader p , system sizes n, f, Δ , weight config. $W = \langle R_{max}, R_{min} \rangle$, latency matrices for $PROPOSE$ \hat{M}^P and $WRITE$ \hat{M}^W , consensus rounds r
Result: consensus latency of the AWARE leader

```

1  $V_{max} \leftarrow 1 + \frac{\Delta}{f}$   $V_{min} \leftarrow 1$   $V_i \leftarrow \begin{cases} V_{max}, & \text{if } i \in R_{max} \\ V_{min}, & \text{otherwise} \end{cases}$ 
2  $\forall i \in I: offset_i \leftarrow 0$ 
3 while  $r > 0$  do
4   for  $i \in I$  do
5      $T_i^{PROPOSED} \leftarrow \max(\hat{M}^P[p, i], offset_i)$ 
6    $(T^{WRITTEN})_{i \in I} \leftarrow formQV(I, \hat{M}^W, (T_i^{PROPOSED})_{i \in I}, (V_i)_{i \in I})$ 
7    $(T^{ACCEPTED})_{i \in I} \leftarrow formQV(I, \hat{M}^W, (T^{WRITTEN})_{i \in I}, (V_i)_{i \in I})$ 
8   for  $i \in I$  do
9      $offset_i \leftarrow T_i^{ACCEPTED} - T_p^{ACCEPTED}$ 
10   $consensusLatencies_r \leftarrow T_p^{ACCEPTED}$ 
11   $r \leftarrow r - 1$ 
12 return average of  $consensusLatencies$ 

```

Handwritten notes:
 → the time each replica receives leader's proposal
 → leader
 → the time each replica receives leader's proposal

$M^W[i, j] = L_i^W[j]$ if replica i sent its $WRITE$ measurements within the last c consensus instances, or $M^W[i, j] = +\infty$ if i did not send its measurements. The same applies to M^P .

5. The next step is to deterministically sanitize the matrices to avoid the influence of malicious replicas (see §4.3), generating \hat{M}^P and \hat{M}^W .
6. Now, every replica solves the following optimization problem, where *PredictLatency* (Algorithm 2) is a function for predicting the latency of the consensus protocol using the latencies in \hat{M}^P, \hat{M}^W , and a set of weight distributions $W \in \mathcal{W}$ and permitted leaders $l \in \mathcal{L}$:

$$\langle \hat{l}, \hat{W} \rangle = \arg \min_{W \in \mathcal{W}, l \in \mathcal{L}} PredictLatency(l, W, \hat{M}^P, \hat{M}^W) \quad (7)$$

In the end, the configuration $\langle \hat{l}, \hat{W} \rangle$ that provides optimal leader consensus latency is the one selected for the next reconfiguration if the predicted latency is better than the current configuration by the factor α (optimization goal). Note that since this procedure is deterministic, the $\langle l, W \rangle$ is the same in all replicas.

7. In case the replicas find a faster weight configuration, they update their view to respect the new voting weights for the following consensus instances. Optionally, if the system uses leader relocation, the replicas might also trigger a view change to elect a faster leader.

DETERMINISTIC LATENCY PREDICTION

→ PAST CONFIGURATION → the one that yields low consensus latency from the leader's perspective

OPTIMIZATION