

# AWARE: Resilient Wide-Area Byzantine Consensus Using Adaptive Weighted Replication

Christian Berger

University of Passau,  
Germany

Hans P. Reiser

João Sousa

LASIGE, Faculdade de Ciências,  
Universidade de Lisboa, Portugal

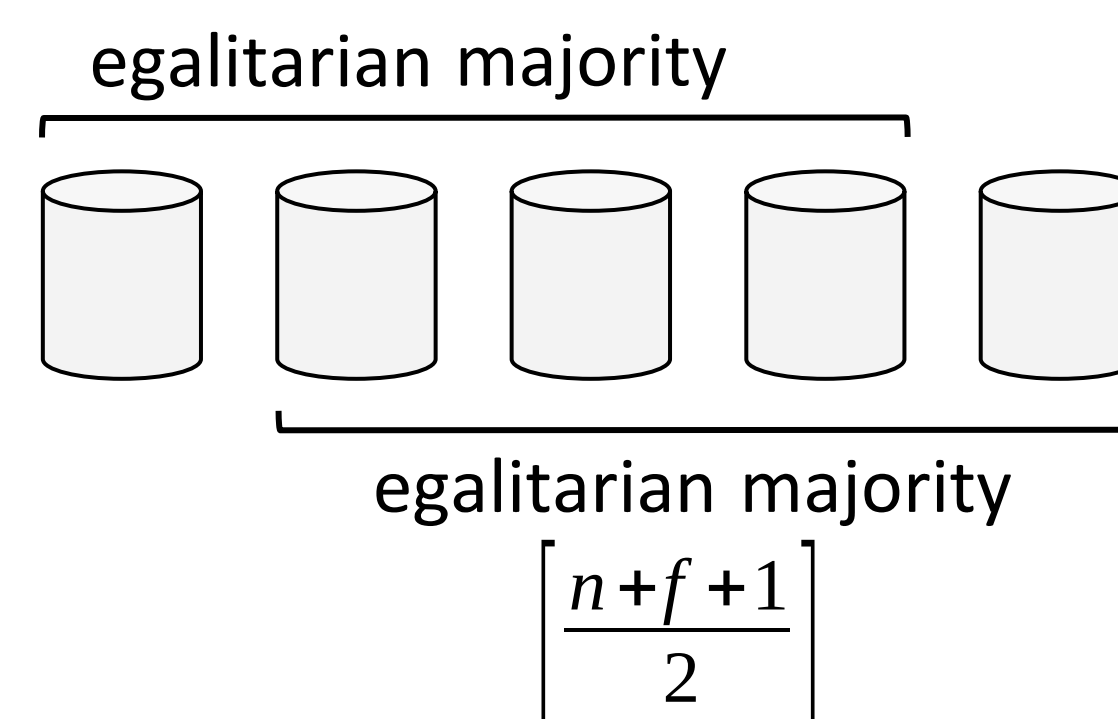
Alysson Bessani

## Problem

- ① In **geo-replicated systems**, the heterogeneous latencies of connections between replicas limit the system's ability to achieve fast consensus
- ➔ WHEAT uses  $\Delta$  additional spare replicas and **weighted replication** to **faster make progress** by accessing a proportionally smaller quorum of replicas
- ② However, the benefit of weighted replication depends on choosing an optimal weight configuration (a non-trivial problem!)
- ③ The **environment of the SMR system** (i.e network characteristics) may change at runtime and thus the optimal configuration may also change
- ➔ AWARE enables geo-replicated systems to **adapt to their environment**

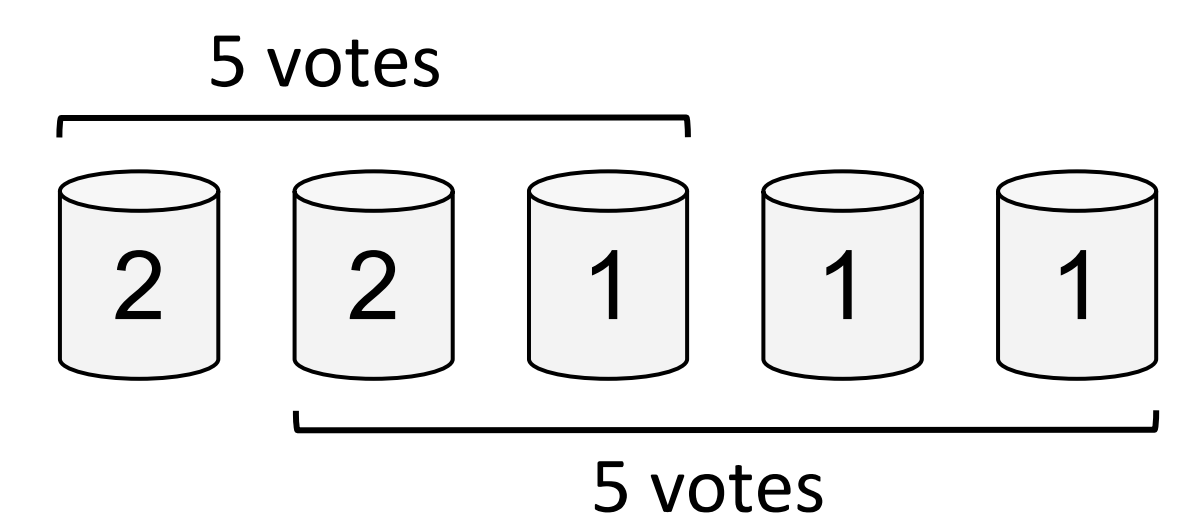
## Weighted Replication

### 1) Egalitarian majority



### 2) Weighted

(from  $2f + 1$  to  $n - f$  replicas)

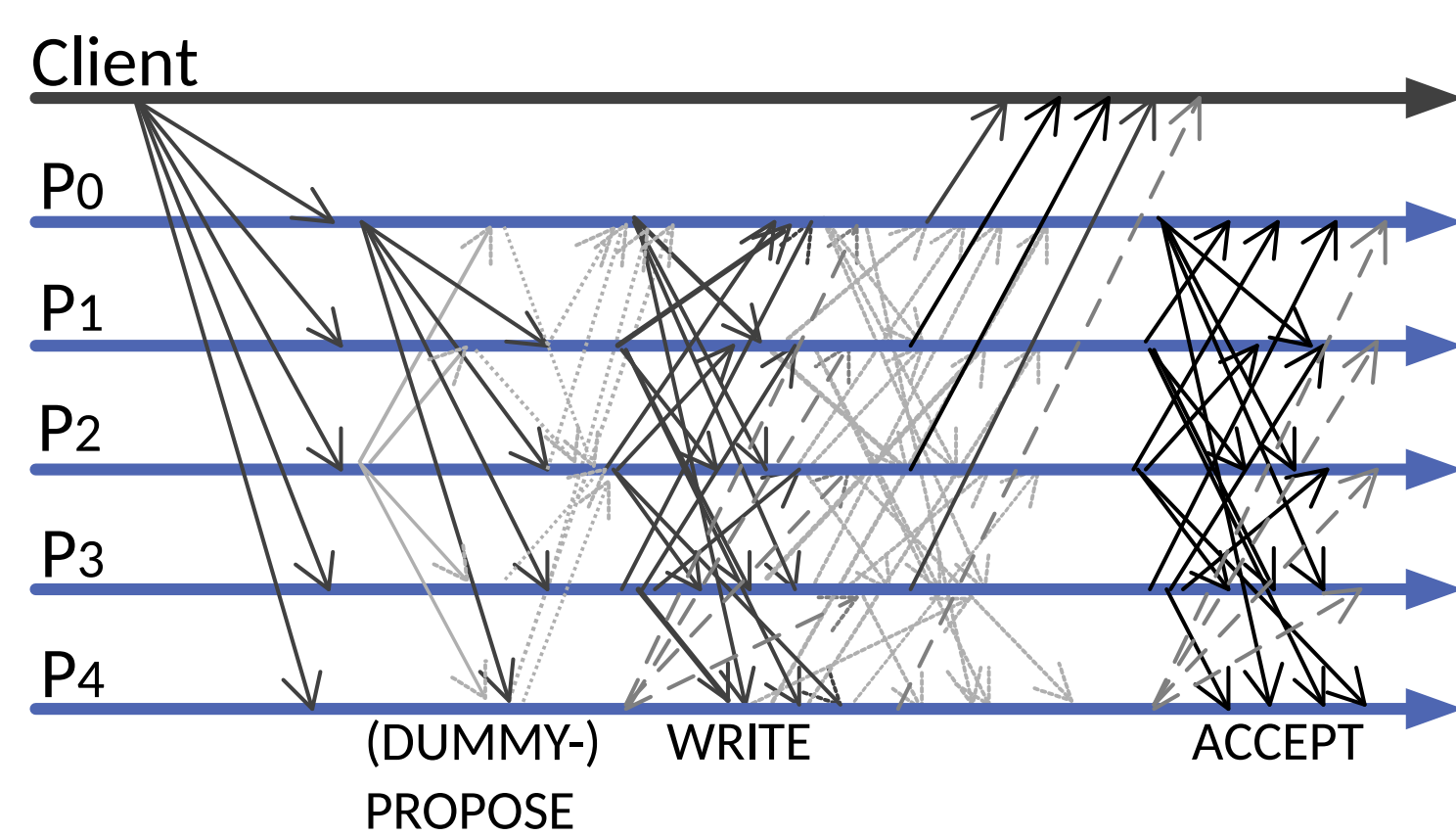


Binary weighting scheme: only  $V_{max}$  or  $V_{min}$  can be assigned

- ★ Weighted replication is **safe** and does **not violate the resilience bound  $f$**

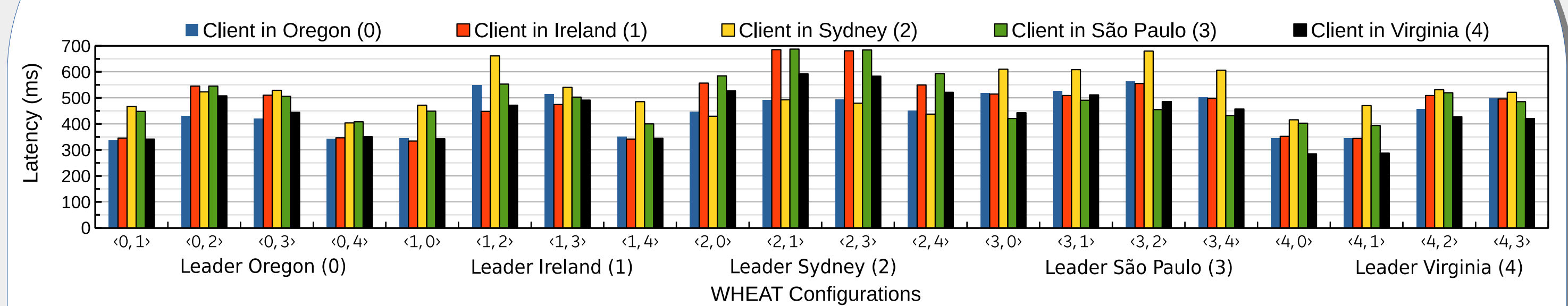
## Monitoring

- ★ AWARE uses **reliable self-monitoring as decision-making basis** for adapting replicas' voting weights or leader position at runtime
- ① **We measure the Propose latency of non-leaders:** periodically, an alternately selected dummy leader broadcasts a dummy proposal
- ② **We measure the Write latency:** replicas immediately respond by sending acknowledgments
- ★ Replicas periodically **disseminate their measurements with total order**, thus maintain the same latency matrix after some specific consensus instance
- ➔ AWARE maintains these **synchronized matrices** for both Propose and Write latencies  $\hat{M}^P$  and  $\hat{M}^W$  used for decisions later



## Evaluation

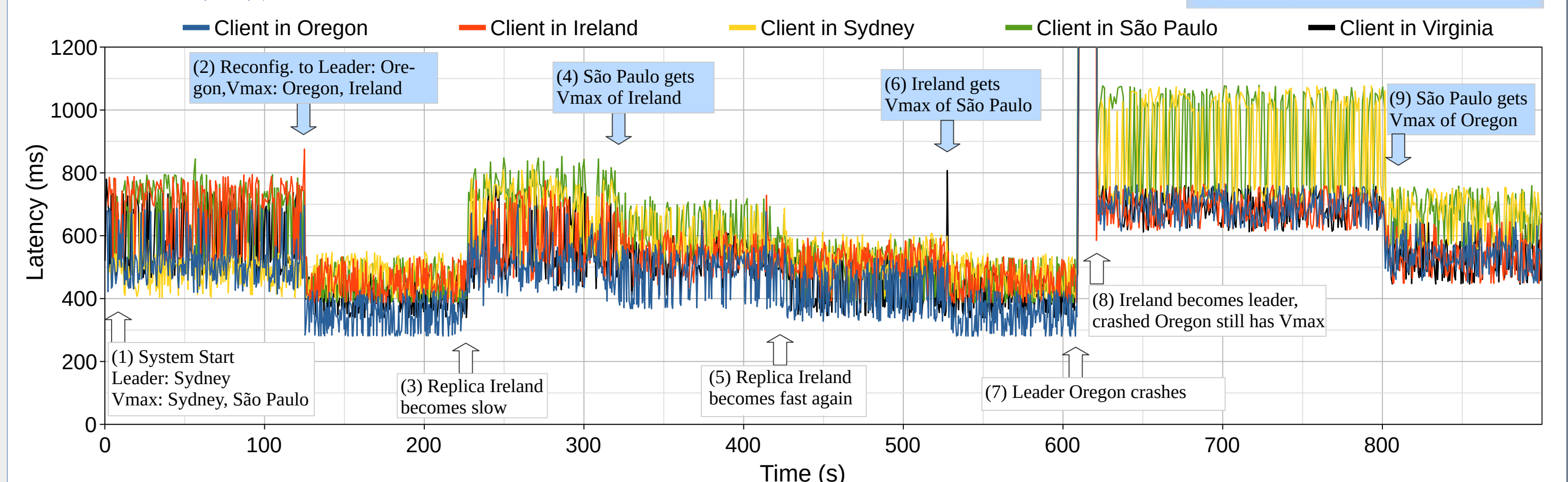
### ★ Measured average request latency of 11th to 90th percentile across clients in different regions



Configuration  $\langle L, R \rangle$  means L is the leader and R is the other replica (besides the leader) with a voting weight of  $V_{max} = 2$

- ① **Tuning voting weights** can reduce latency (compare configs with same leader)
- ② **Leader relocation** may be necessary for achieving optimal consensus latency
- ③ **A global optimum does not exist** but a few pareto-optimal configurations dominate poorer performing configurations

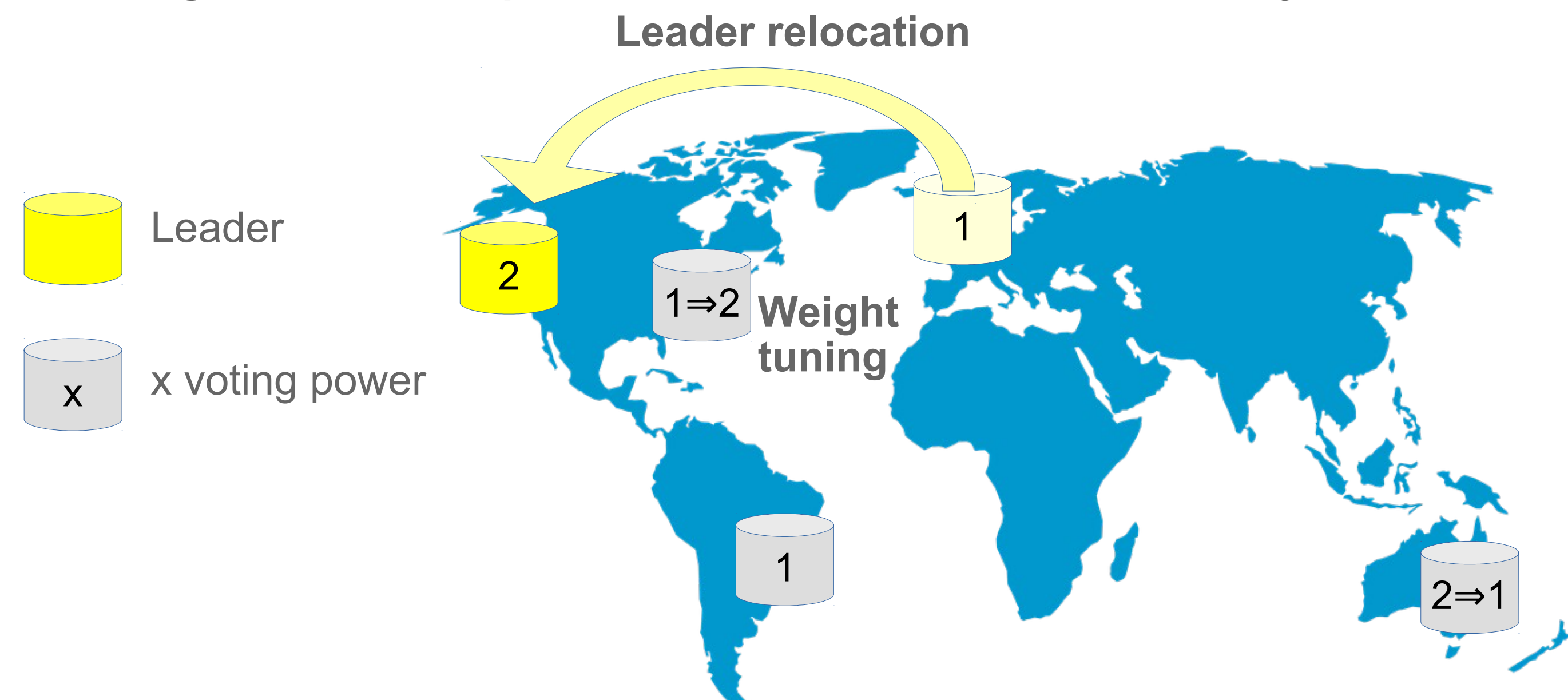
### ★ Runtime behavior of AWARE



- ① **Ease of deployment:** AWARE provides the *needed automation* for finding an optimal configuration by tuning voting weights and/or relocating the leader
- ② **Adjusting to varying network conditions:** if the quality of communication links varies, AWARE *dynamically adjusts to new conditions* by shifting high voting power to replicas that are the fastest in a recent time frame
- ③ **Compensating for faults:** even if  $f$  replicas with high voting power become unavailable and restrict quorum variability, for non-malicious behavior, AWARE detects this and restores the availability of up to  $f$  ( $V_{max} - V_{min}$ ) voting power by redistributing high voting weights

## Self-Optimization

- ★ AWARE continuously strives for latency gains at runtime. We **optimize voting weights and leader position to minimize consensus latency**



- ★ After specific consensus instances, replicas **deterministically solve** an optimization problem: *PredictLatency* is a function to predict the latency of the consensus protocol using the measured latencies in  $\hat{M}^P$ ,  $\hat{M}^W$  and all possible weight distributions  $W \in \mathbf{W}$  and permitted leaders  $l \in \mathbf{L}$ :

$$(\hat{l}, \hat{W}) = \arg \min_{W \in \mathbf{W}, l \in \mathbf{L}} \text{PredictLatency}(l, W, \hat{M}^P, \hat{M}^W)$$

- ➔ Replicas **safely reconfigure** to a new weight or leader configuration if it minimizes the system's consensus latency

## Conclusions

- ★ World-spanning Byzantine consensus systems can benefit from dynamic self-optimizing approaches in combination with weighted replication
- ★ AWARE enriches the idea of weighted replication by providing the needed automation to adapt to changing environmental conditions
- ★ Evaluation results show that the potential for latency and throughput gains is substantial. Specifically, the best configuration performs about 38.7% faster on average in terms of observed latency across clients' sites than the median