

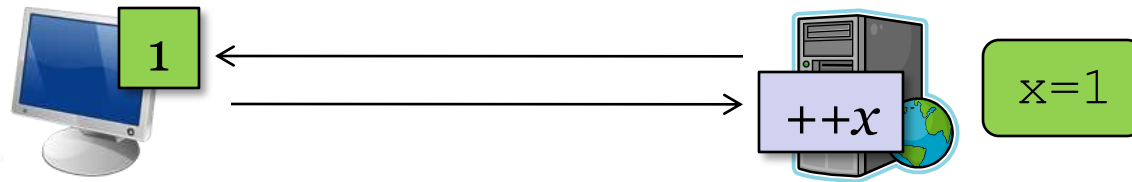
Tolerating Latency in Replicated State Machines through Client Speculation

April 22, 2009

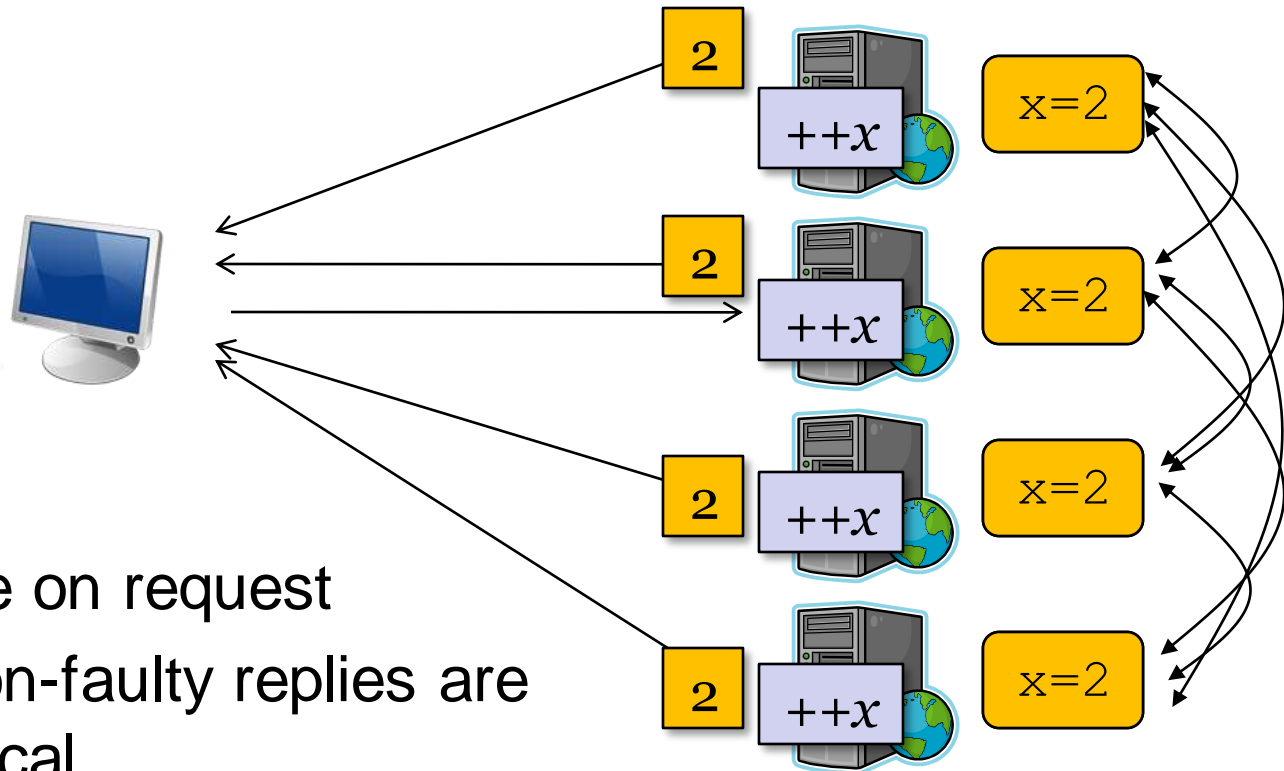
Benjamin Wester¹, James Cowling², Edmund B. Nightingale³,
Peter M. Chen¹, Jason Flinn¹, Barbara Liskov²

University of Michigan¹, MIT CSAIL², Microsoft Research³

Simple Service Configuration

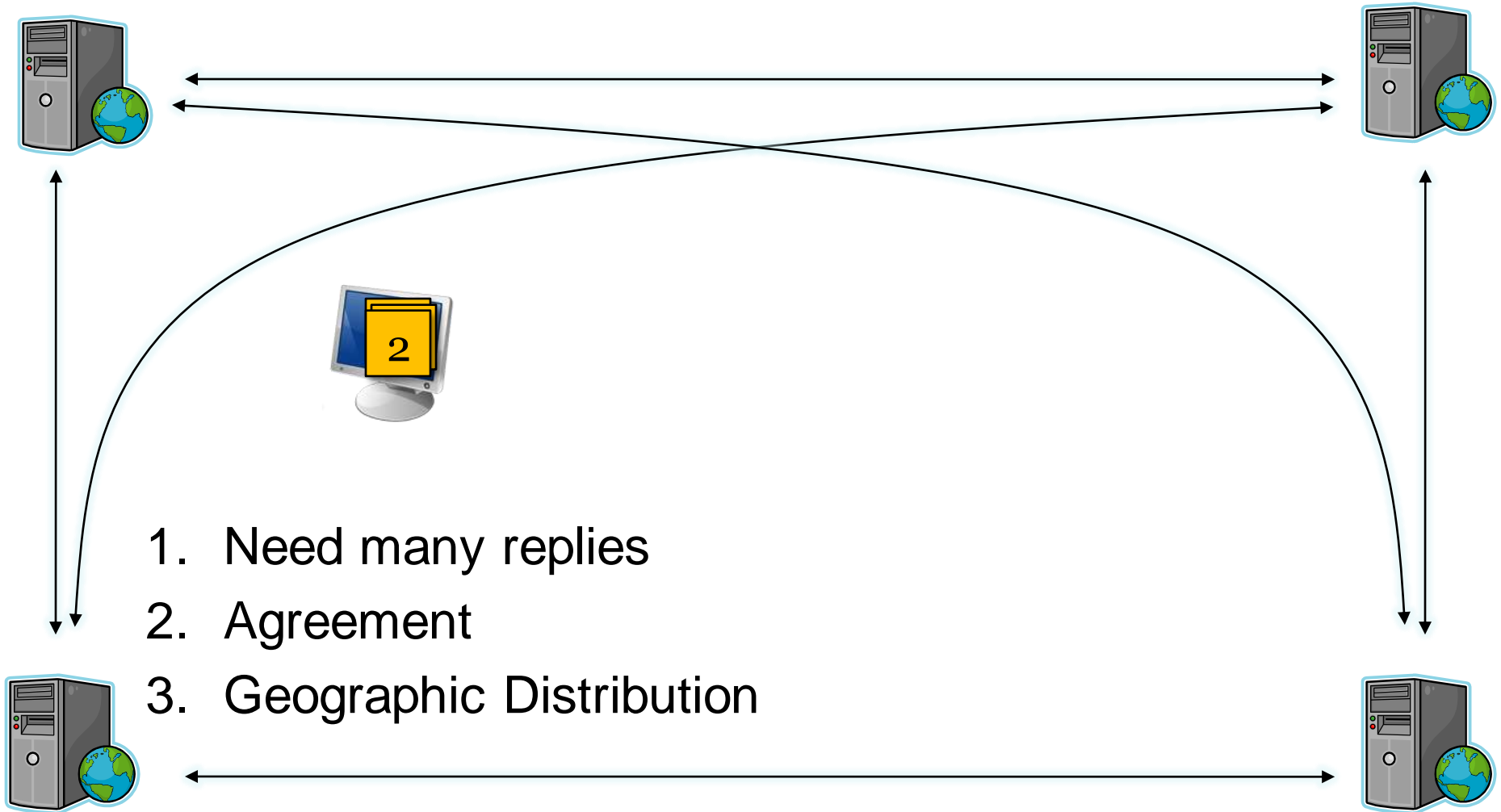


Replicated State Machines (RSM)



- Agree on request
- All non-faulty replies are identical

RSMs have high latency



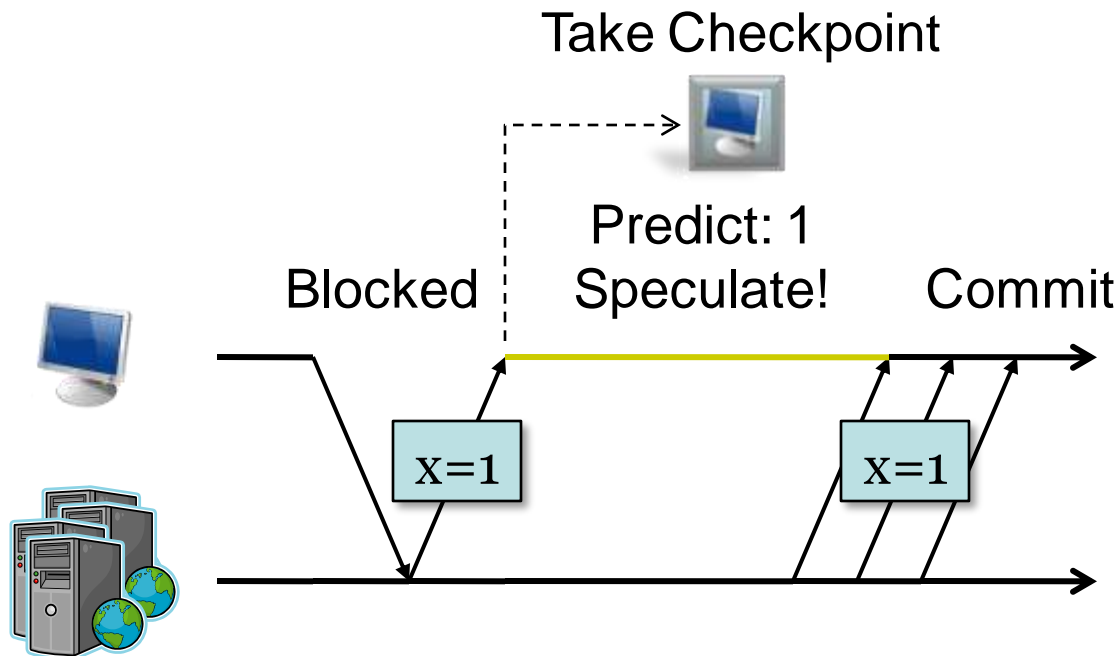
Hide the Latency

- Use speculative execution inside RSM
- Speculate before consensus is reached
 - Without faults, any reply predicts consensus value
 - Let client continue after receiving one reply

Overview

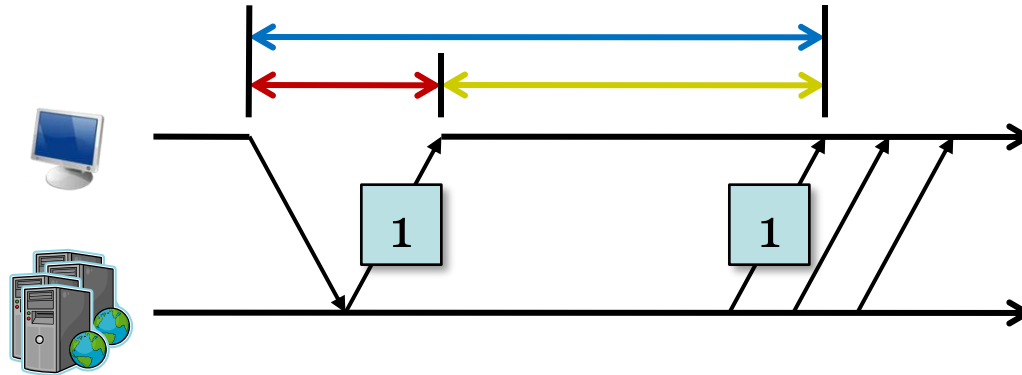
- Introduction
- Improving RSMs with speculation
- Application to PBFT
- Performance
- Conclusion

Speculative Execution in RSM



- Continue processing while waiting

Critical path: first reply

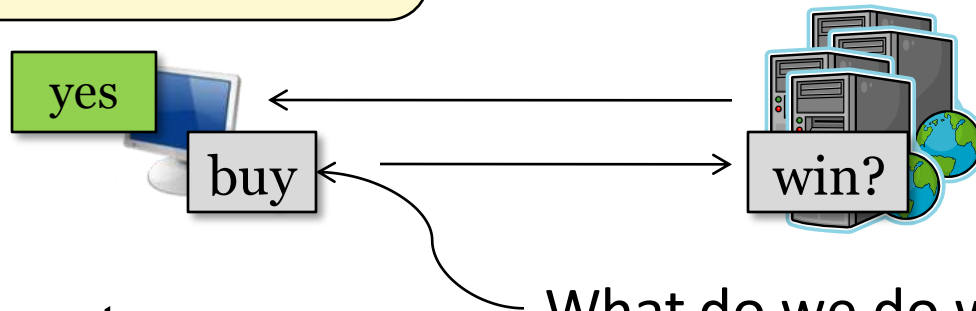


- Completion latency less relevant
- First reply latency sets critical path
 - Speed
 - Accuracy
- Other desirable properties
 - Throughput
 - Stability under contention
 - Smaller number of replicas

Requests while speculative

```
while !check_lottery():  
    submit_tps()  
    buy_corvette()
```

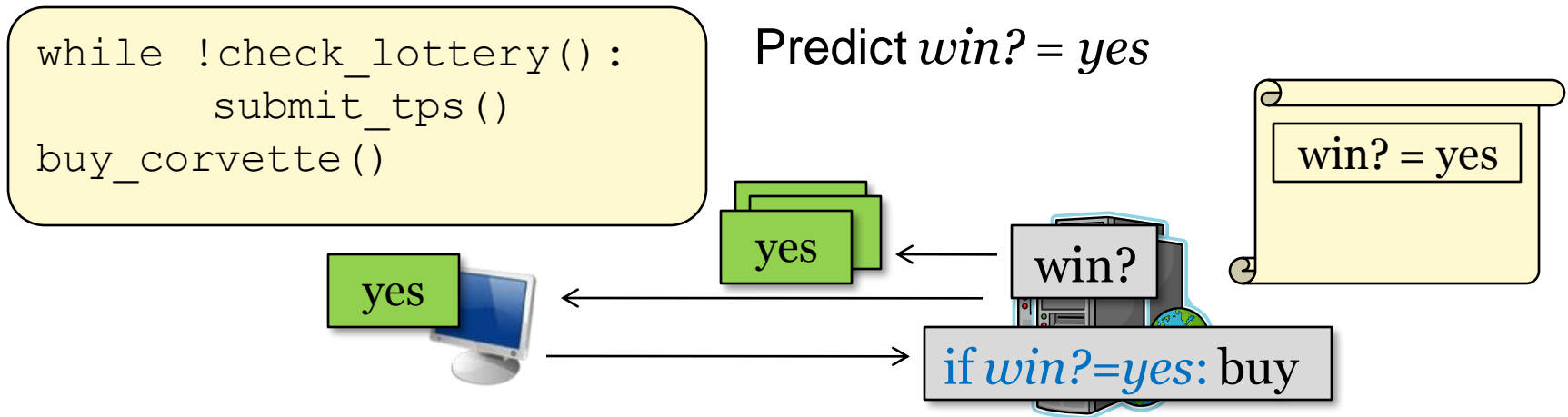
Predict *win?* = *yes*



1. Hold request
 - Bad performance
2. Distributed commit/rollback
 - State tracking complex

What do we do with this?

Resolve speculations on the replicas



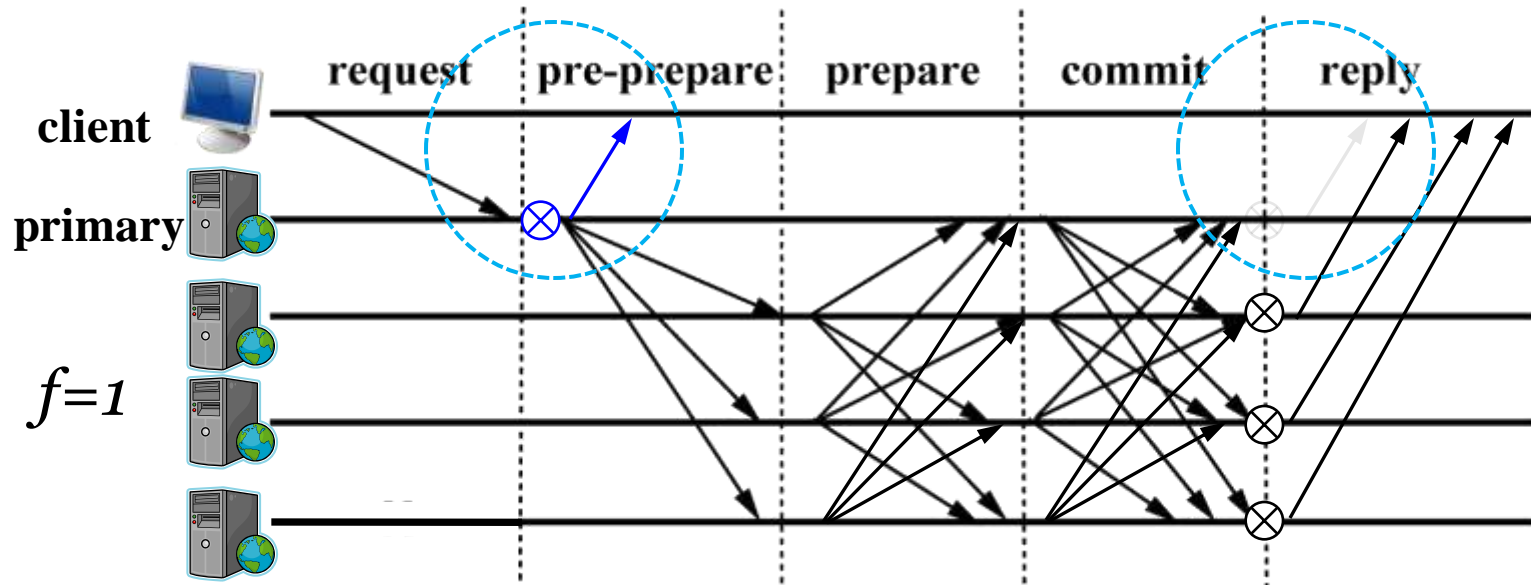
- Explicitly encode dependencies as **predicates**
- No special request handling needed
- Replicas need to log past replies
- Local decision at replicas matches client

Overview

- Introduction
- Improving RSMs with speculation
- **Application to PBFT**
- Performance
- Conclusion

Practical BFT-CS

[Castro and Liskov 1999]



Additional Details

- Tentative execution
 - PBFT/PBFT-CS complete in 4 phases
- Read-only optimization
 - Accurate answer from backup replica
- Failure threshold
 - Bound worst-case failure
- Correctness

Overview

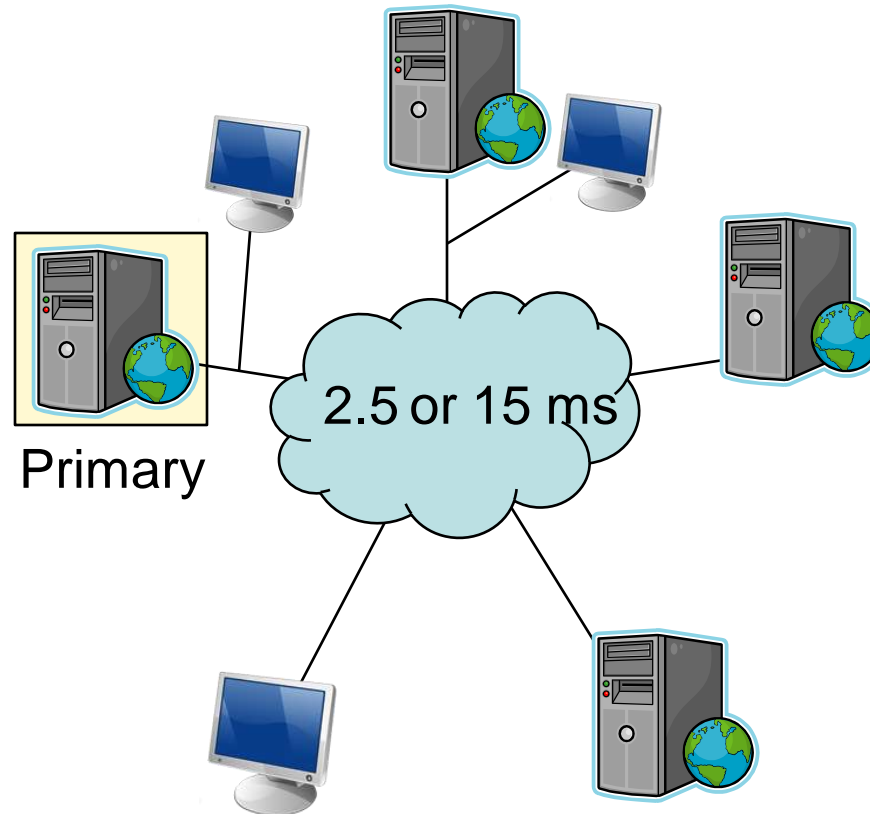
- Introduction
- Improving RSMs with speculation
- Application to PBFT
- **Performance**
- Conclusion

Benchmarks

- Shared counter
 - Simple checkpoint
 - No computation
- NFS: Apache httpd build
 - Complex checkpoint
 - Significant computation

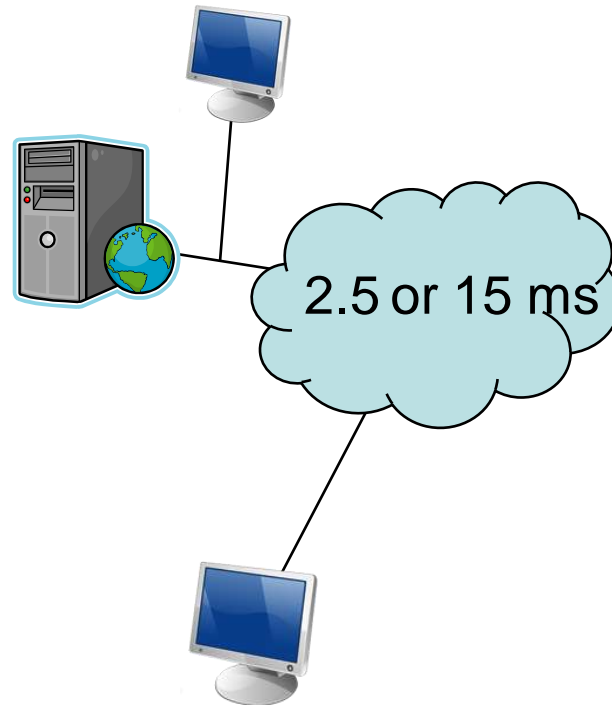
Topology

1. Primary-local
2. Primary-remote
3. Uniform



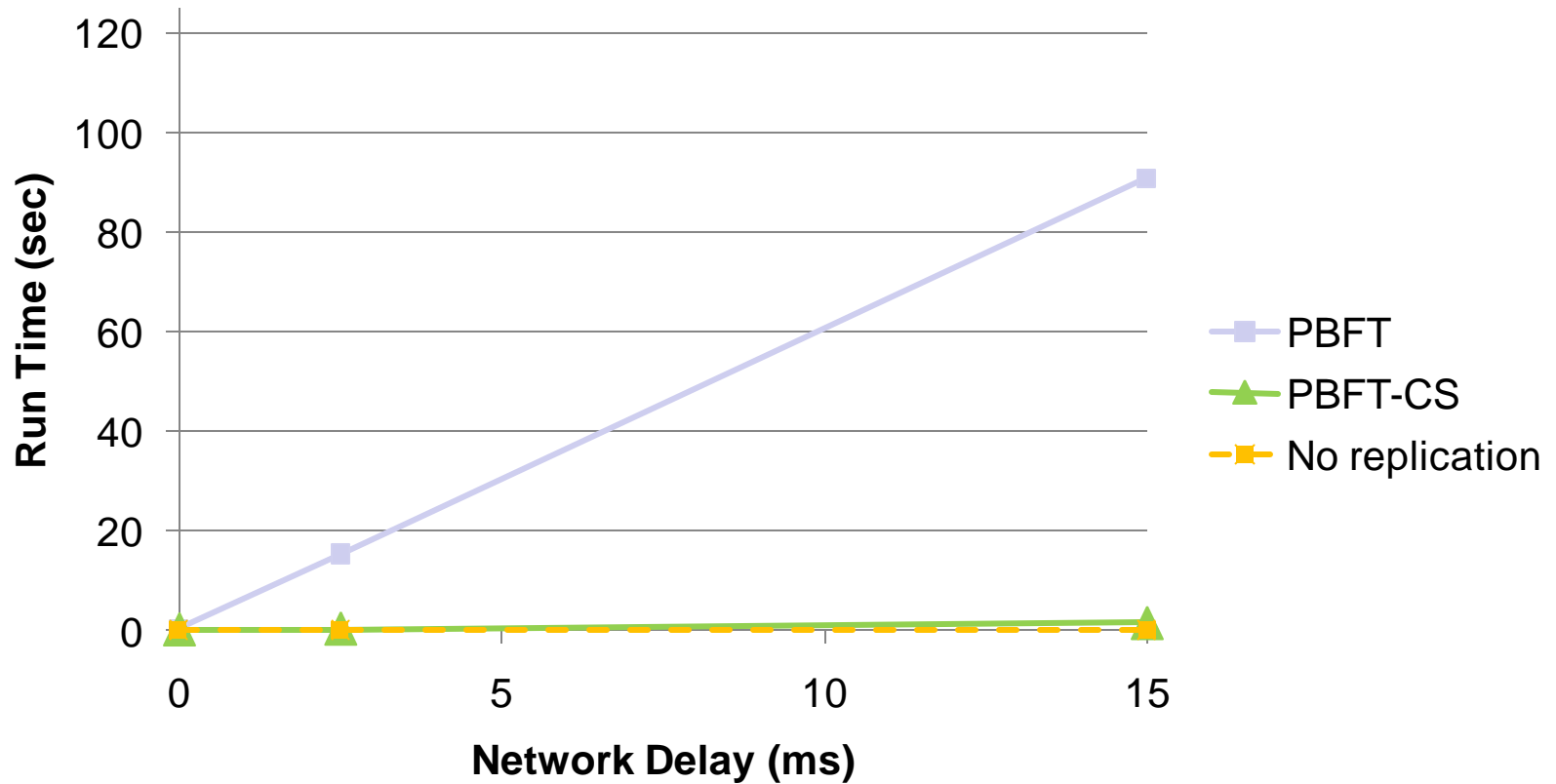
Base case: no replication

1. Primary-local
2. Primary-remote
3. Uniform



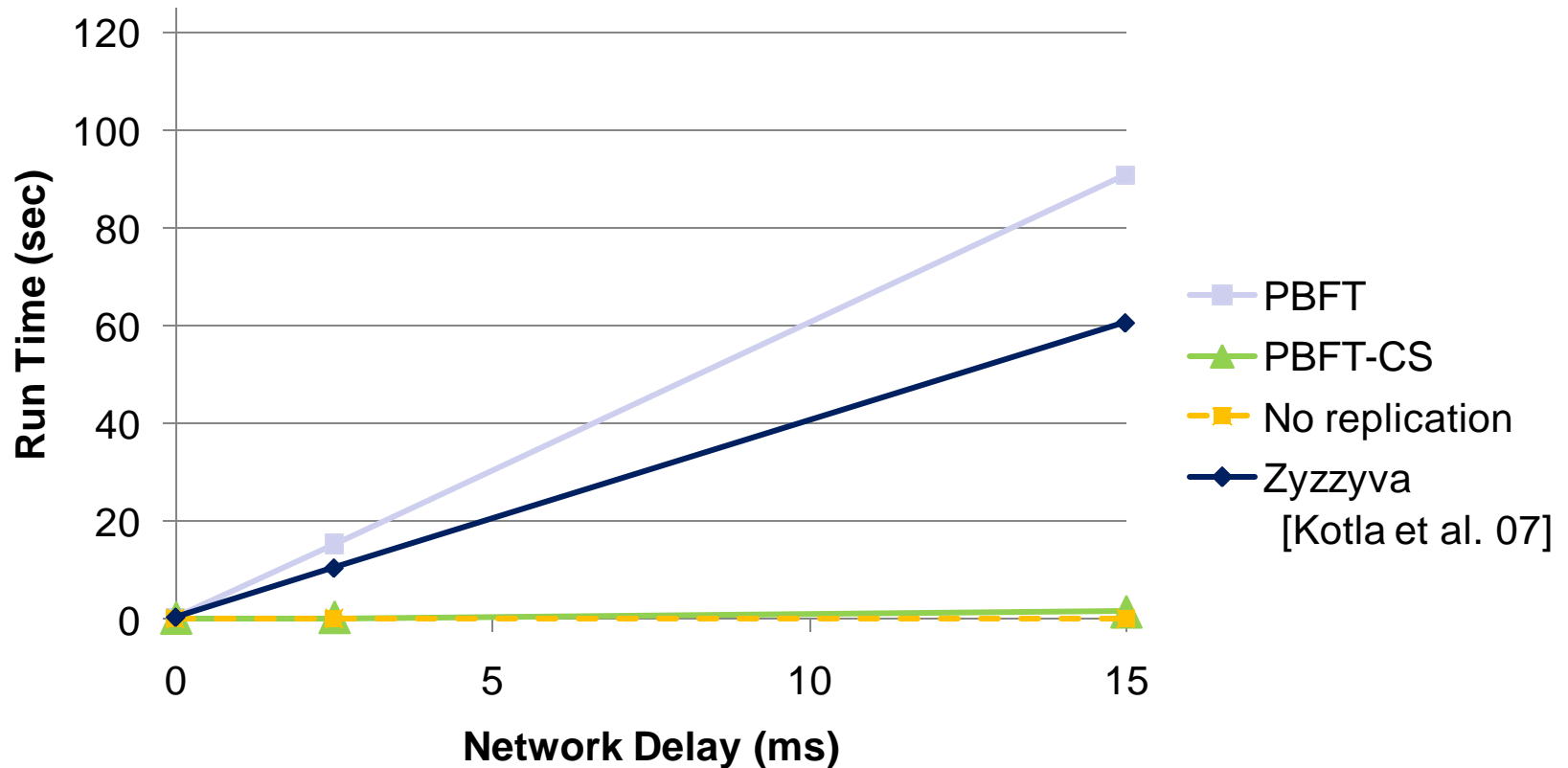
Shared Counter

Primary-local topology



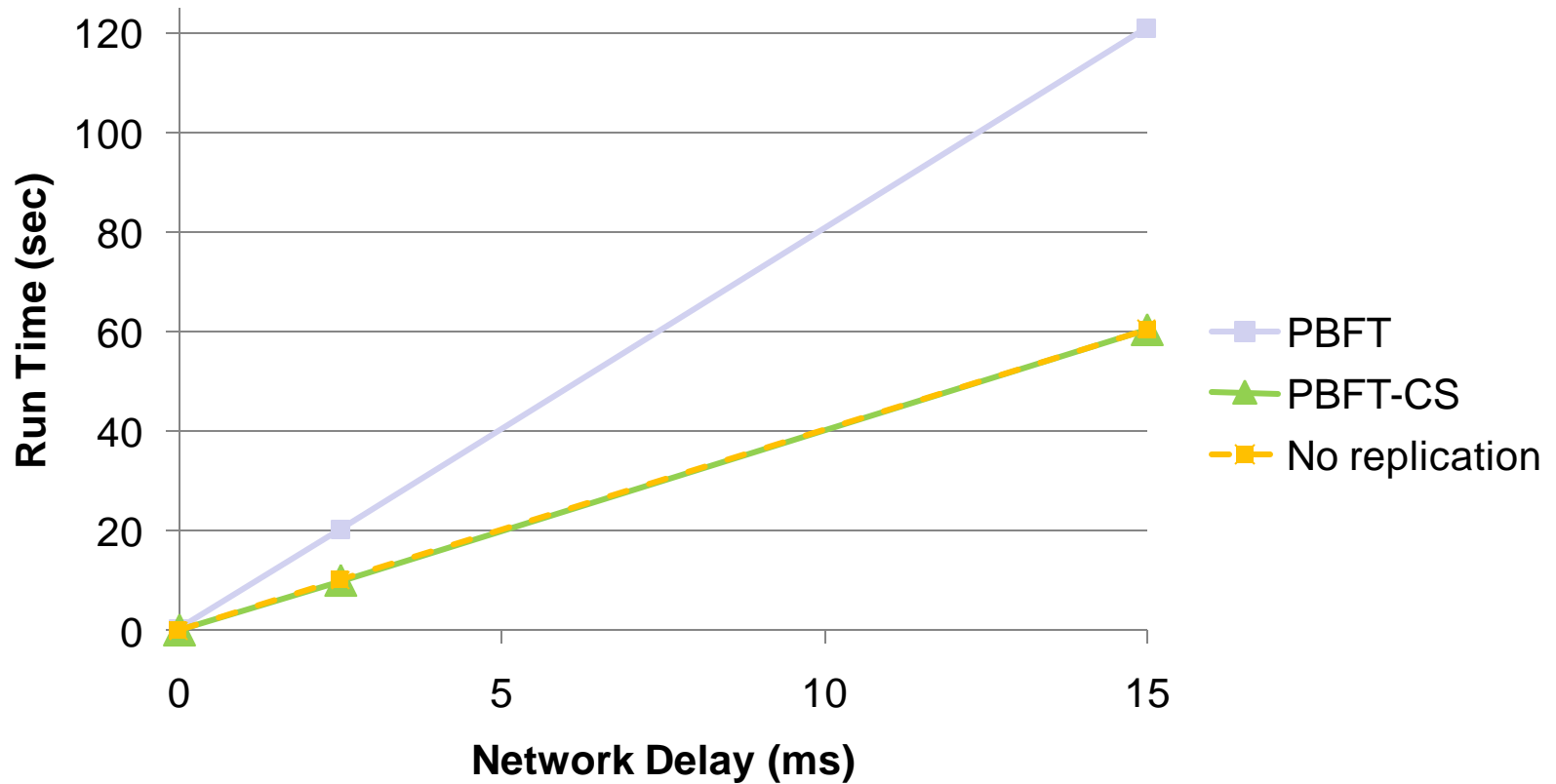
Shared Counter

Primary-local topology



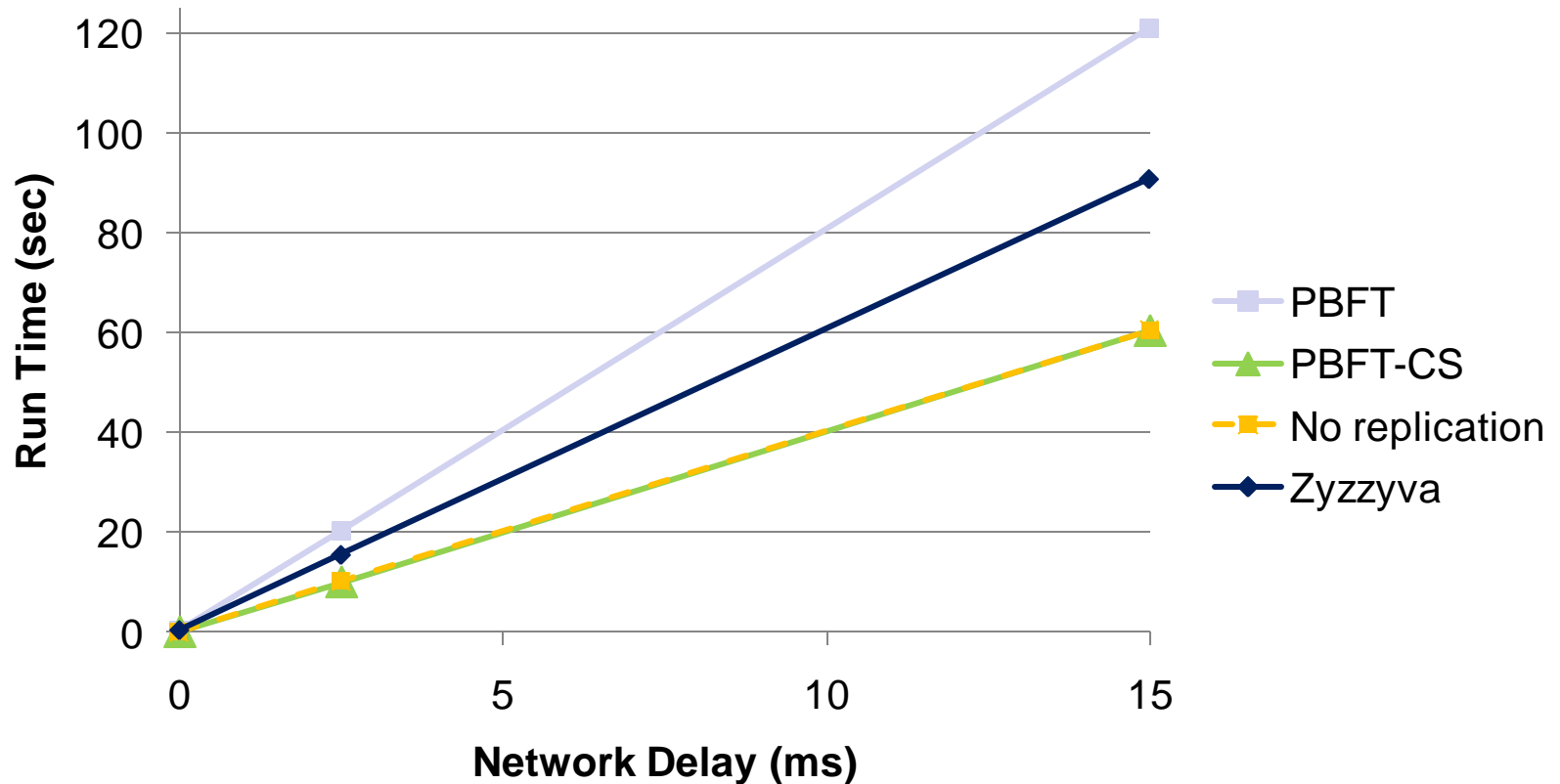
Shared Counter

Uniform & Primary-remote topology



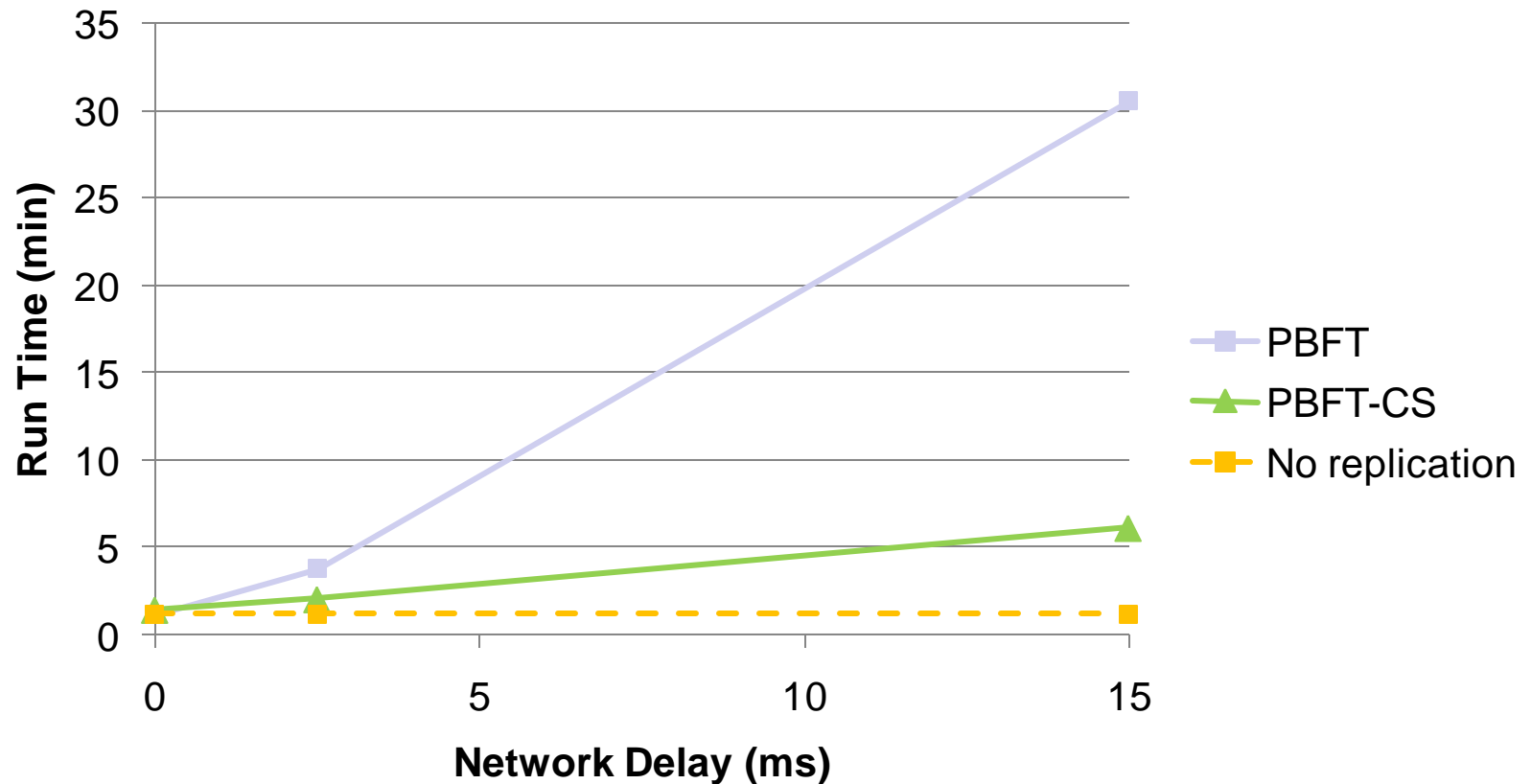
Shared Counter

Uniform & Primary-remote topology



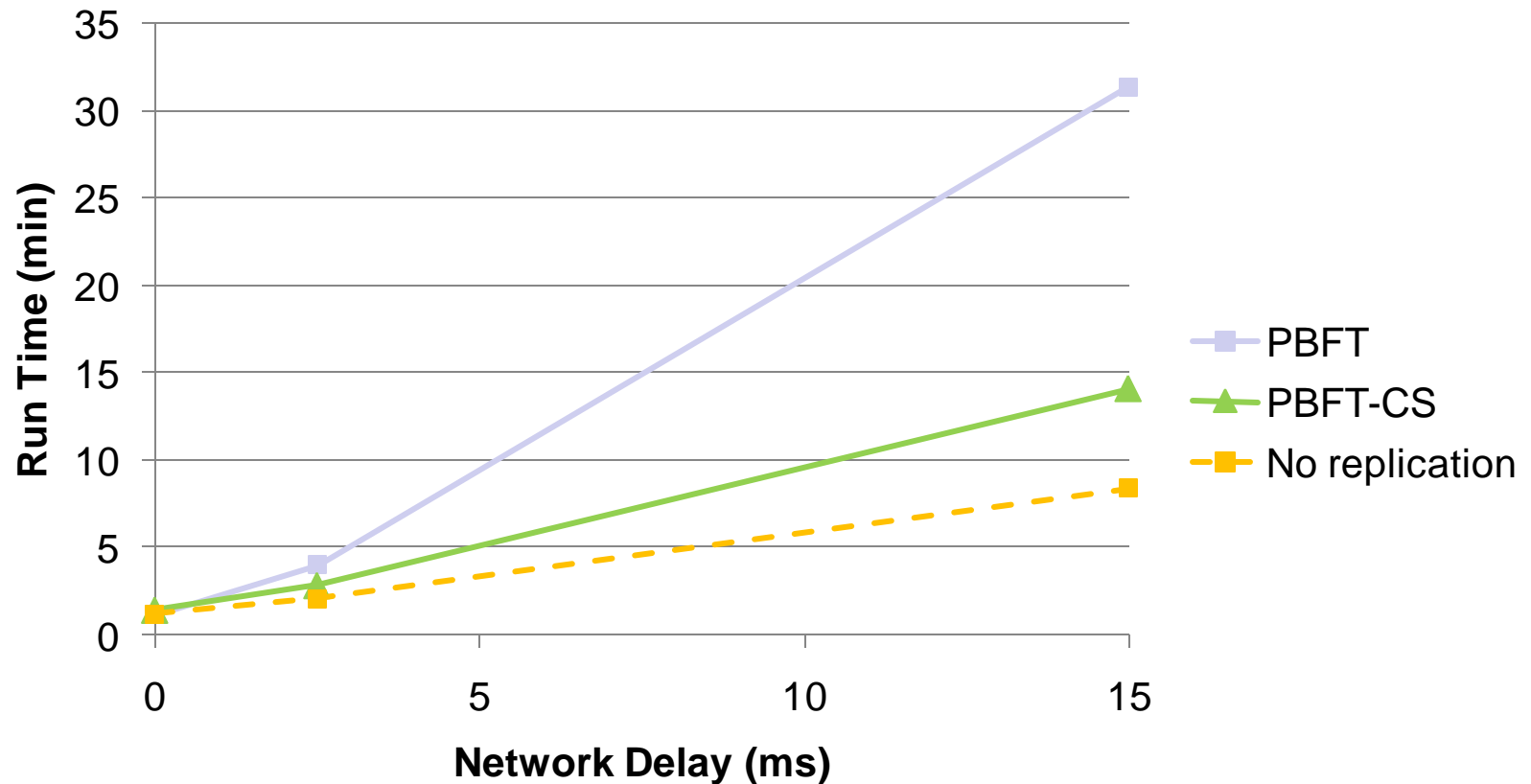
NFS: Apache build

Primary-local topology



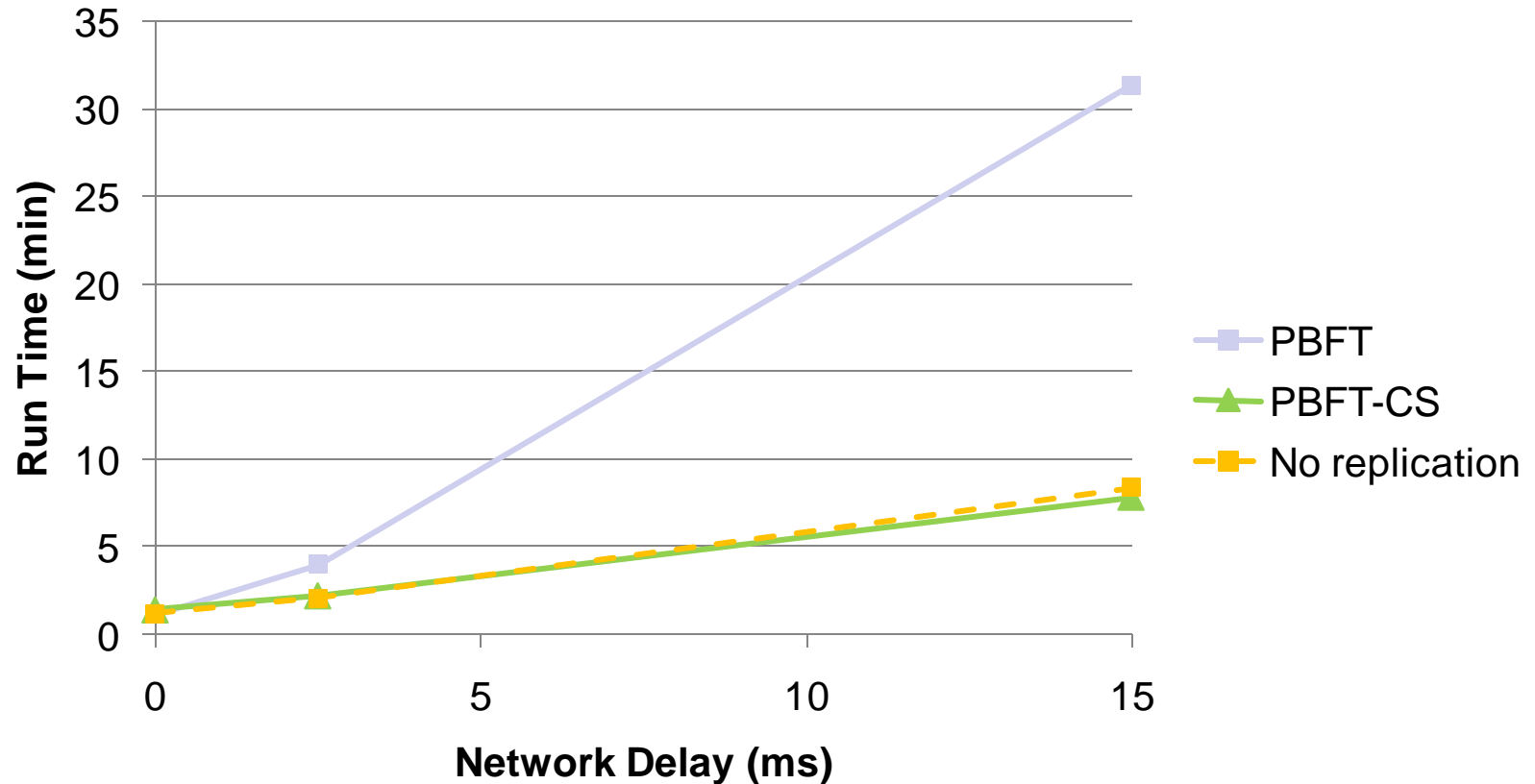
NFS: Apache build

Uniform topology



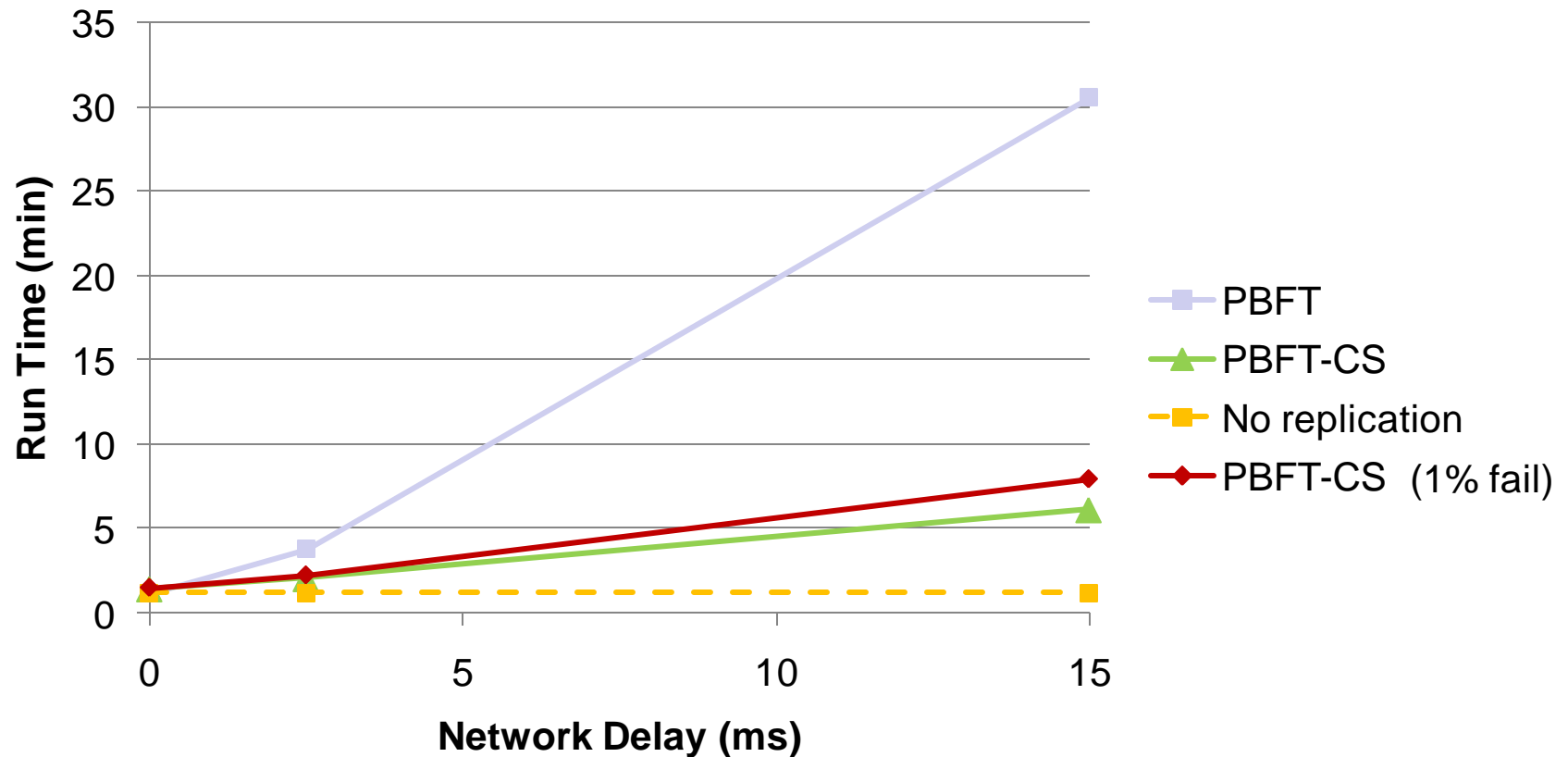
NFS: Apache build

Primary-remote topology



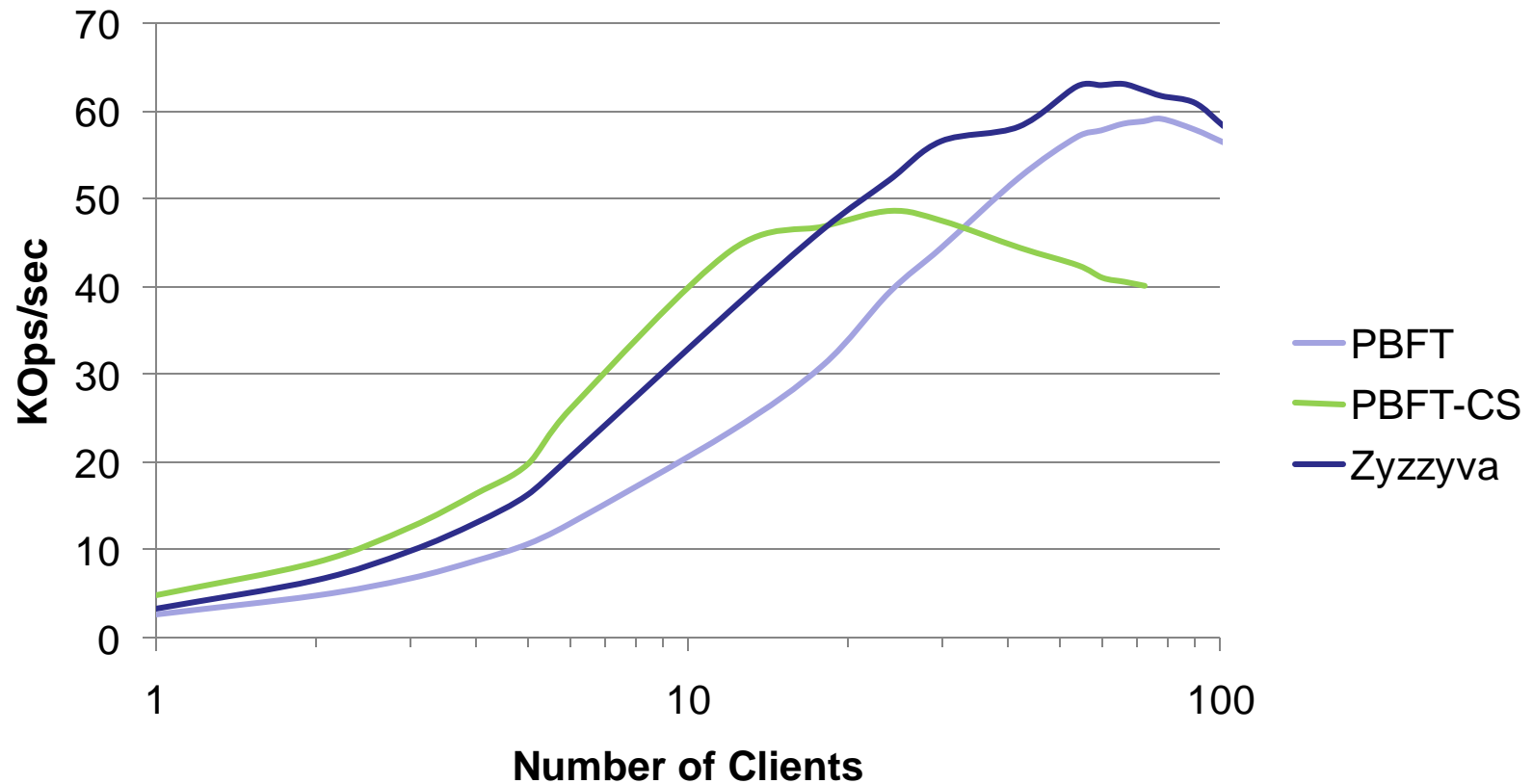
NFS: With Failure

Primary-local topology



Throughput (Shared Counter)

LAN topology



Conclusion

- Integrate client speculation within RSMs
- Predicated requests: performance without complexity
- Clients less sensitive to latency between replicas
- 5x speedup over non-speculative protocol

Makes WAN deployments more practical