

# Fault-Tolerant Distributed Transactions on Blockchain

## *Beyond the Design of PBFT*



Suyash Gupta



Jelle Hellings



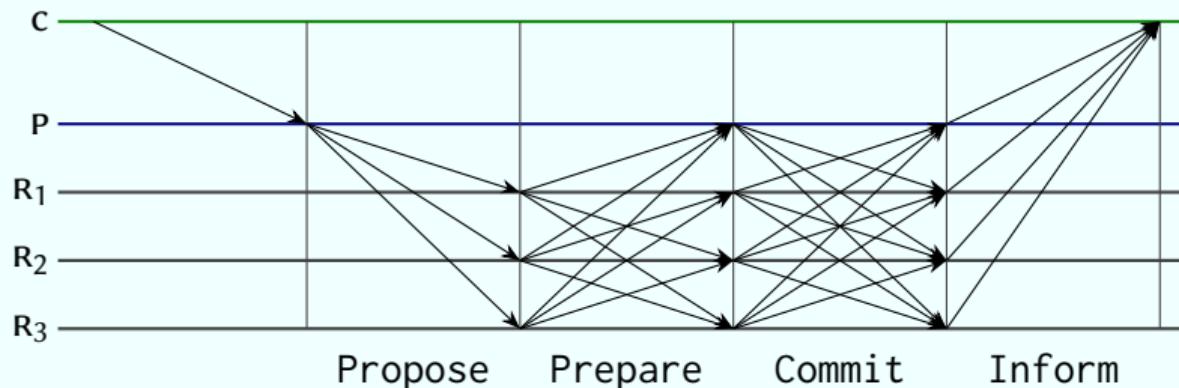
Mohammad Sadoghi

**UCDAVIS**  
UNIVERSITY OF CALIFORNIA

 **Expolab**  
Creativity Unfolded

**McMaster**  
University 

## Previously: PBFT



## Central Question

What is the *expected performance* of PBFT? Motivate!

# On the Performance of Consensus

**Consensus throughput** Decisions per second made by consensus.

**Consensus latency** Duration of a single round of consensus.

**Resource utilization** The cost of consensus (e.g., computational, network bandwidth).  
*Imbalance* in resource utilized by replicas (e.g., primary).

# On the Performance of Consensus

**Consensus throughput** Decisions per second made by consensus.

**Consensus latency** Duration of a single round of consensus.

**Resource utilization** The cost of consensus (e.g., computational, network bandwidth).  
*Imbalance* in resource utilized by replicas (e.g., primary).

**Complexity** Complexity of normal-case and of recovery (e.g., view-change).

**Failure Model** The types of failures consensus can deal with.

# On the Performance of Consensus

**Consensus throughput** Decisions per second made by consensus.

**Consensus latency** Duration of a single round of consensus.

**Resource utilization** The cost of consensus (e.g., computational, network bandwidth).  
*Imbalance* in resource utilized by replicas (e.g., primary).

**Complexity** Complexity of normal-case and of recovery (e.g., view-change).

**Failure Model** The types of failures consensus can deal with.

**Client latency** Duration between a client request and the outcome.

# On the Performance of Consensus

**Consensus throughput** Decisions per second made by consensus.

**Consensus latency** Duration of a single round of consensus.

**Resource utilization** The cost of consensus (e.g., computational, network bandwidth).  
*Imbalance* in resource utilized by replicas (e.g., primary).

**Complexity** Complexity of normal-case and of recovery (e.g., view-change).

**Failure Model** The types of failures consensus can deal with.

**Client latency** Duration between a client request and the outcome.

- ▶ *Low loads*: Function of the consensus latency.
- ▶ *High loads*: Function of the consensus throughput.

## Determining the Performance Variables

Number of replicas determines the amount of messages exchanged.

Network bandwidth determines how long it takes to exchange these messages.

Message delay determines how long it takes for sent messages to arrive.

Computational speed determines the speed by which messages are processed.

## Determining the Performance Variables

Number of replicas determines the amount of messages exchanged.

Network bandwidth determines how long it takes to exchange these messages.

Message delay determines how long it takes for sent messages to arrive.

Computational speed determines the speed by which messages are processed.

System processing client transactions

## Determining the Performance Variables

Number of replicas determines the amount of messages exchanged.

Network bandwidth determines how long it takes to exchange these messages.

Message delay determines how long it takes for sent messages to arrive.

Computational speed determines the speed by which messages are processed.

### System processing client transactions

- ▶ Bottlenecks *outside consensus*: speed by which replicas *execute transactions*.

## Determining the Performance Variables

Number of replicas determines the amount of messages exchanged.

Network bandwidth determines how long it takes to exchange these messages.

Message delay determines how long it takes for sent messages to arrive.

Computational speed determines the speed by which messages are processed.

### System processing client transactions

- ▶ Bottlenecks *outside consensus*: speed by which replicas *execute transactions*.
- ▶ Computational speed typically sufficient when *parallelization* is used.

# Determining the Performance Variables

Number of replicas determines the amount of messages exchanged.

Network bandwidth determines how long it takes to exchange these messages.

Message delay determines how long it takes for sent messages to arrive.

Computational speed determines the speed by which messages are processed.

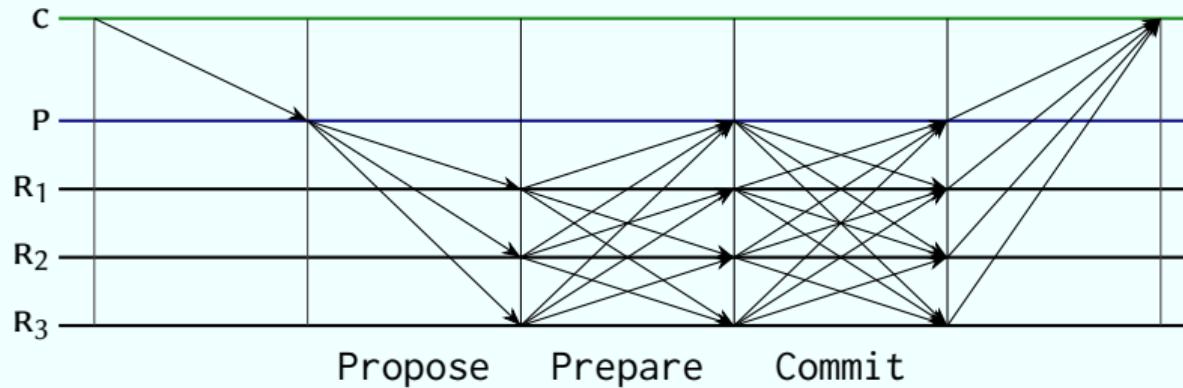
## System processing client transactions

- ▶ Bottlenecks *outside consensus*: speed by which replicas *execute transactions*.
- ▶ Computational speed typically sufficient when *parallelization* is used.

**Bottleneck in practice:** consensus performance in terms of throughput and latency  
(as a function of *network bandwidth* and *message delay*).

# The Single-Round Cost of PBFT (Sketch)

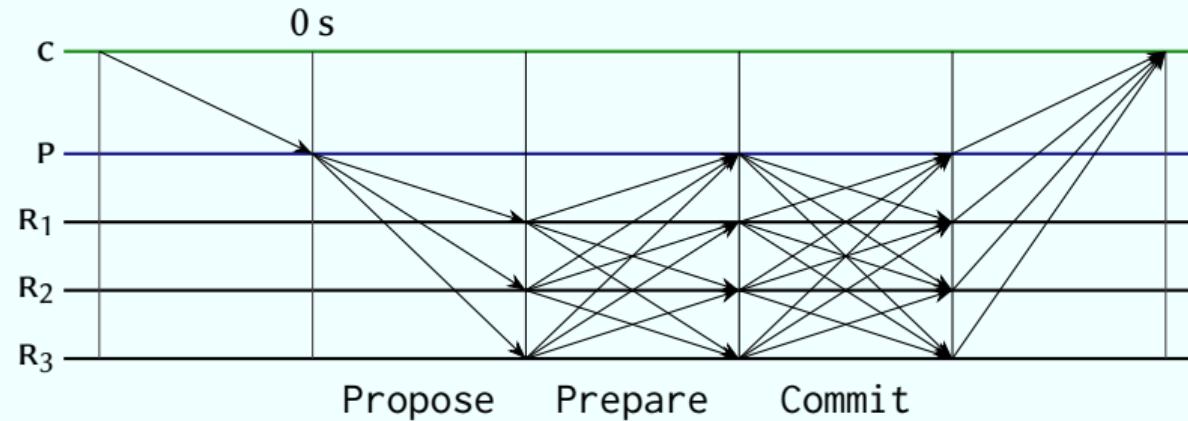
Assumption: Network bandwidth  $B = 100 \text{ MiB/s}$  and delay  $\delta = 15 \text{ ms}$



# The Single-Round Cost of PBFT (Sketch)

Assumption: Network bandwidth  $B = 100 \text{ MiB/s}$  and delay  $\delta = 15 \text{ ms}$

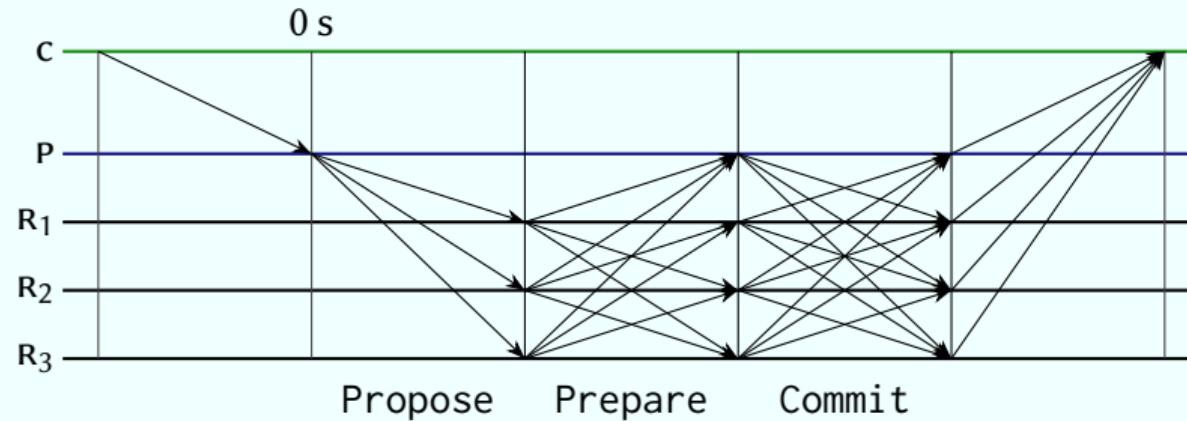
Propose:  $s_t = 4048 \text{ B}$  each.



# The Single-Round Cost of PBFT (Sketch)

Assumption: Network bandwidth  $B = 100 \text{ MiB/s}$  and delay  $\delta = 15 \text{ ms}$

Propose:  $s_t = 4048 \text{ B}$  each.

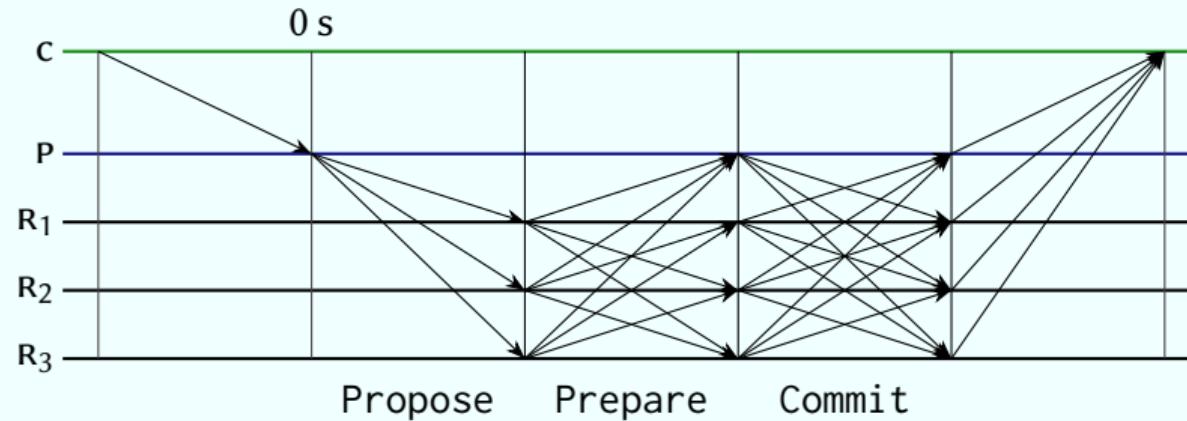


- ▶  $n - 1$  messages
- ▶  $s_t \text{ B}$  each

# The Single-Round Cost of PBFT (Sketch)

Assumption: Network bandwidth  $B = 100 \text{ MiB/s}$  and delay  $\delta = 15 \text{ ms}$

Propose:  $s_t = 4048 \text{ B}$  each.

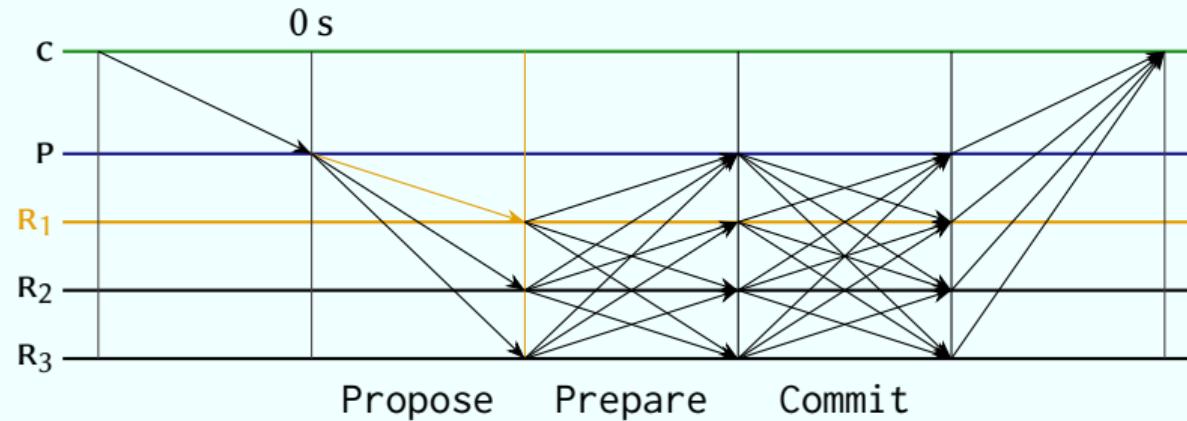


- ▶  $n - 1$  messages
  - ▶  $s_t \text{ B}$  each
- $\left. \right\} (n - 1)s_t = 3 \cdot 4048 = 12\,144 \text{ B}$

# The Single-Round Cost of PBFT (Sketch)

Assumption: Network bandwidth  $B = 100 \text{ MiB/s}$  and delay  $\delta = 15 \text{ ms}$

Propose:  $s_t = 4048 \text{ B}$  each.

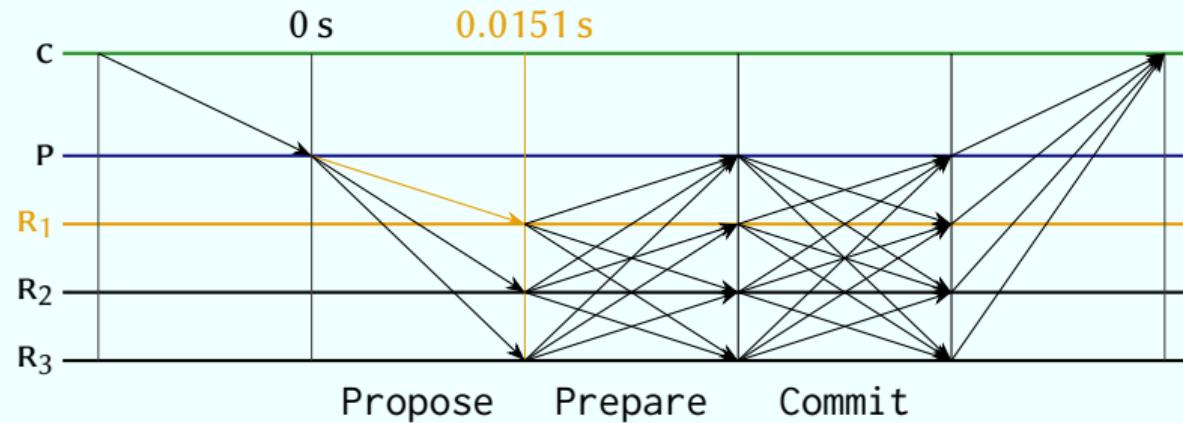


- ▶  $n - 1$  messages
  - ▶  $s_t \text{ B}$  each
  - ▶  $B \text{ MiB/s}$
  - ▶ Last byte arrives after  $\delta$
- $\left. \begin{array}{l} \\ \\ \end{array} \right\} (n - 1)s_t = 3 \cdot 4048 = 12\,144 \text{ B}$

# The Single-Round Cost of PBFT (Sketch)

Assumption: Network bandwidth  $B = 100 \text{ MiB/s}$  and delay  $\delta = 15 \text{ ms}$

Propose:  $s_t = 4048 \text{ B}$  each.

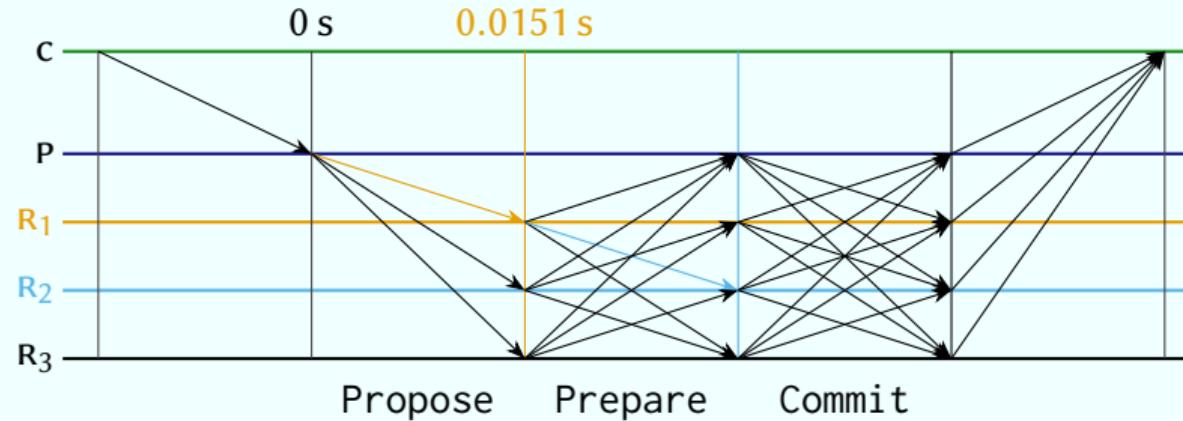


- ▶  $n - 1$  messages
  - ▶  $s_t \text{ B}$  each
  - ▶  $B \text{ MiB/s}$
  - ▶ Last byte arrives after  $\delta$
- $\left. \begin{array}{l} (\mathbf{n}-1)s_t = 3 \cdot 4048 = 12\,144 \text{ B} \\ \frac{(\mathbf{n}-1)s_t}{B} + \delta = \frac{12\,144}{100 \cdot 2^{20}} + 0.015 \approx 0.0151 \text{ s.} \end{array} \right\}$

# The Single-Round Cost of PBFT (Sketch)

Assumption: Network bandwidth  $B = 100 \text{ MiB/s}$  and delay  $\delta = 15 \text{ ms}$

Propose:  $s_t = 4048 \text{ B}$  each. Prepare and Commit:  $s_m = 256 \text{ B}$  each.

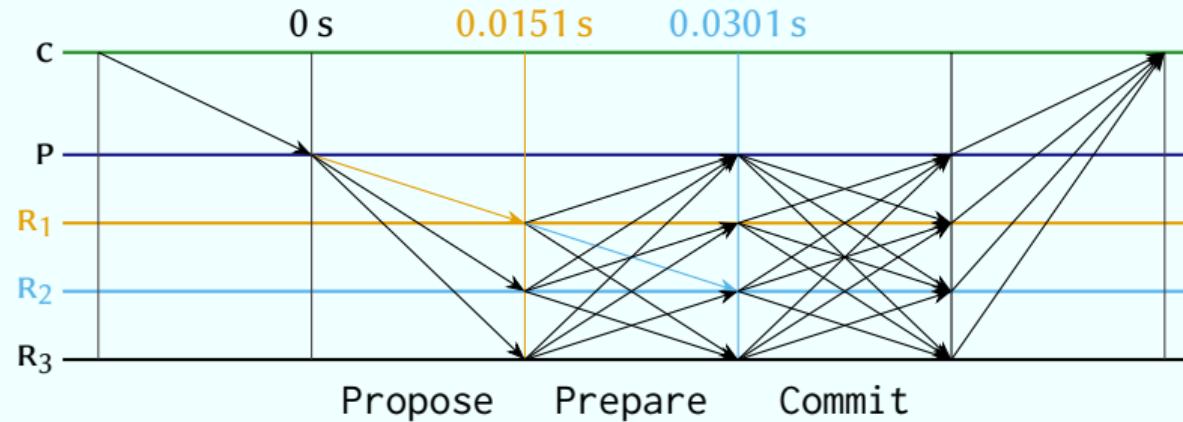


- ▶  $n - 1$  messages
- ▶  $s_m \text{ B}$  each
- ▶  $B \text{ MiB/s}$
- ▶ Last byte arrives after  $\delta$

# The Single-Round Cost of PBFT (Sketch)

Assumption: Network bandwidth  $B = 100 \text{ MiB/s}$  and delay  $\delta = 15 \text{ ms}$

Propose:  $s_t = 4048 \text{ B}$  each. Prepare and Commit:  $s_m = 256 \text{ B}$  each.

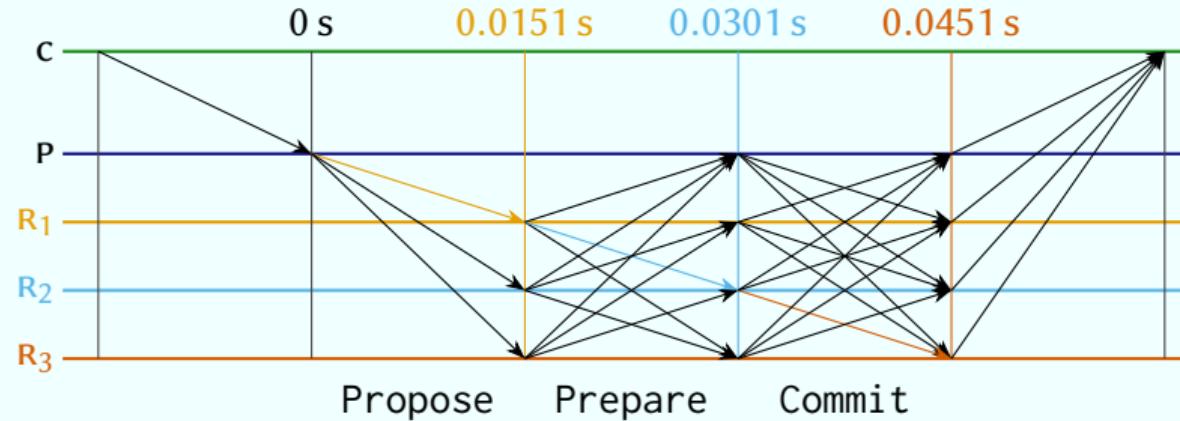


- ▶  $n - 1$  messages
  - ▶  $s_m \text{ B}$  each
  - ▶  $B \text{ MiB/s}$
  - ▶ Last byte arrives after  $\delta$
- $$\left. \begin{array}{l} \frac{(n-1)s_m}{B} + \delta = \frac{768}{100 \cdot 2^{20}} + 0.015 \\ \approx 0.0150 \text{ s.} \end{array} \right\}$$

# The Single-Round Cost of PBFT (Sketch)

Assumption: Network bandwidth  $B = 100 \text{ MiB/s}$  and delay  $\delta = 15 \text{ ms}$

Propose:  $s_t = 4048 \text{ B}$  each. Prepare and Commit:  $s_m = 256 \text{ B}$  each.

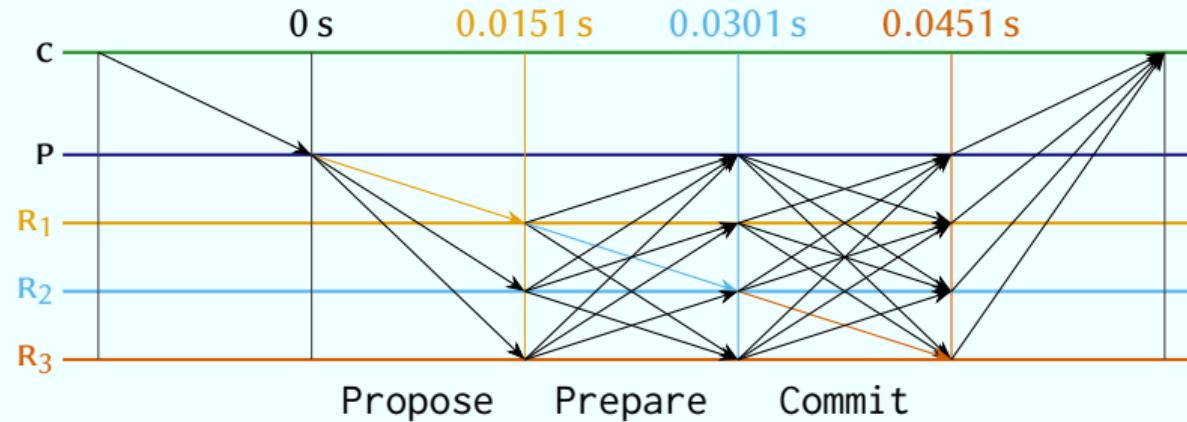


- ▶  $n - 1$  messages
  - ▶  $s_m \text{ B}$  each
  - ▶  $B \text{ MiB/s}$
  - ▶ Last byte arrives after  $\delta$
- $\} \approx 0.0150 \text{ s.}$

# The Single-Round Cost of PBFT (Sketch)

Assumption: Network bandwidth  $B = 100 \text{ MiB/s}$  and delay  $\delta = 15 \text{ ms}$

Propose:  $s_t = 4048 \text{ B}$  each. Prepare and Commit:  $s_m = 256 \text{ B}$  each.

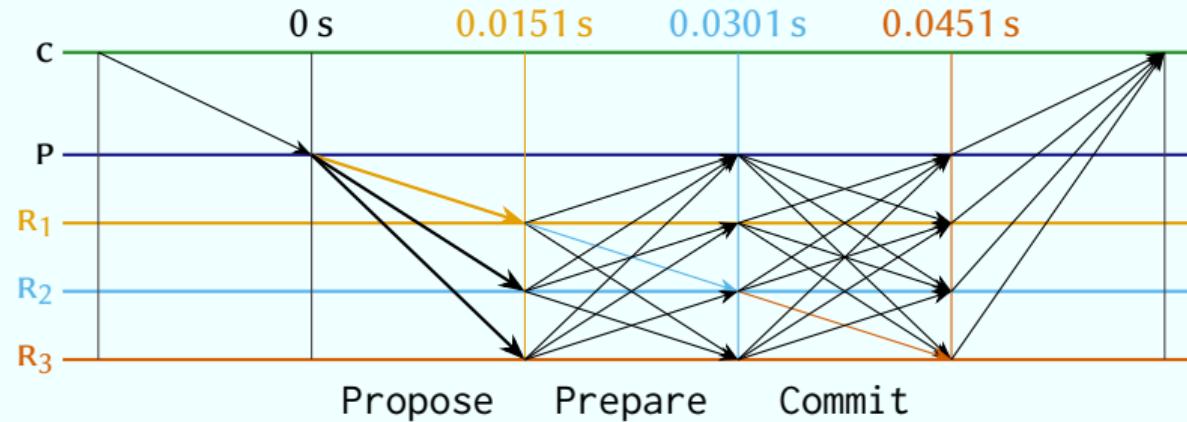


$$\Delta_{\text{PBFT}} =$$

# The Single-Round Cost of PBFT (Sketch)

Assumption: Network bandwidth  $B = 100 \text{ MiB/s}$  and delay  $\delta = 15 \text{ ms}$

Propose:  $s_t = 4048 \text{ B}$  each. Prepare and Commit:  $s_m = 256 \text{ B}$  each.

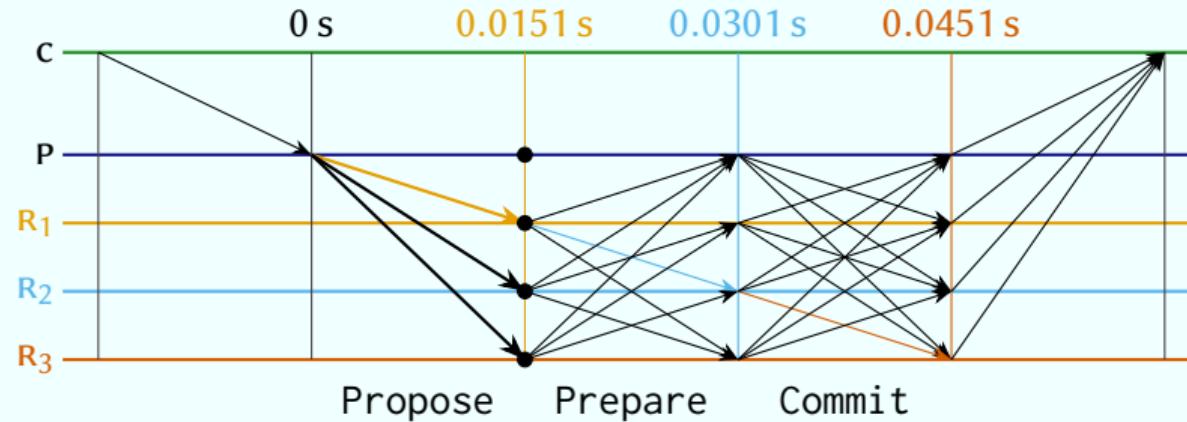


$$\Delta_{\text{PBFT}} = \frac{(n - 1)s_t}{B}$$

# The Single-Round Cost of PBFT (Sketch)

Assumption: Network bandwidth  $B = 100 \text{ MiB/s}$  and delay  $\delta = 15 \text{ ms}$

Propose:  $s_t = 4048 \text{ B}$  each. Prepare and Commit:  $s_m = 256 \text{ B}$  each.

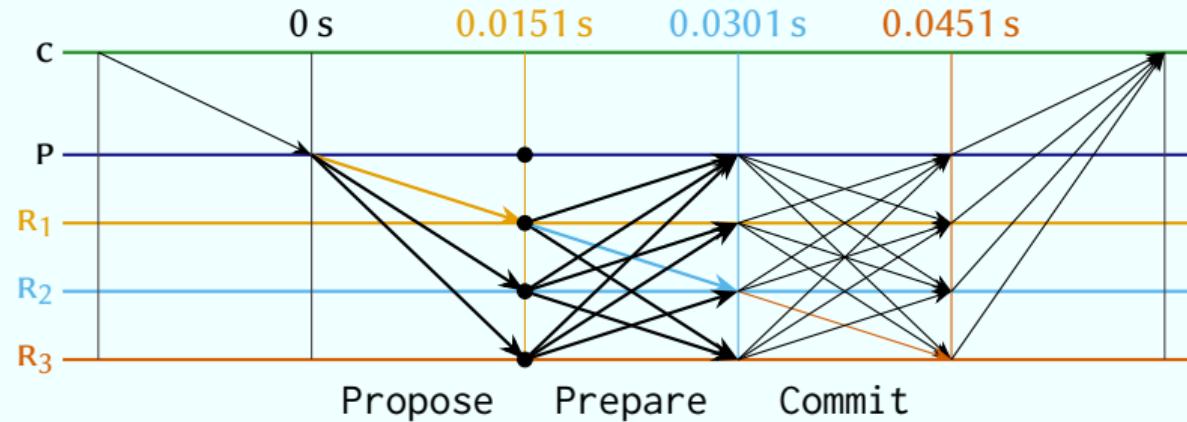


$$\Delta_{\text{PBFT}} = \frac{(n - 1)s_t}{B} + \delta$$

# The Single-Round Cost of PBFT (Sketch)

Assumption: Network bandwidth  $B = 100 \text{ MiB/s}$  and delay  $\delta = 15 \text{ ms}$

Propose:  $s_t = 4048 \text{ B}$  each. Prepare and Commit:  $s_m = 256 \text{ B}$  each.

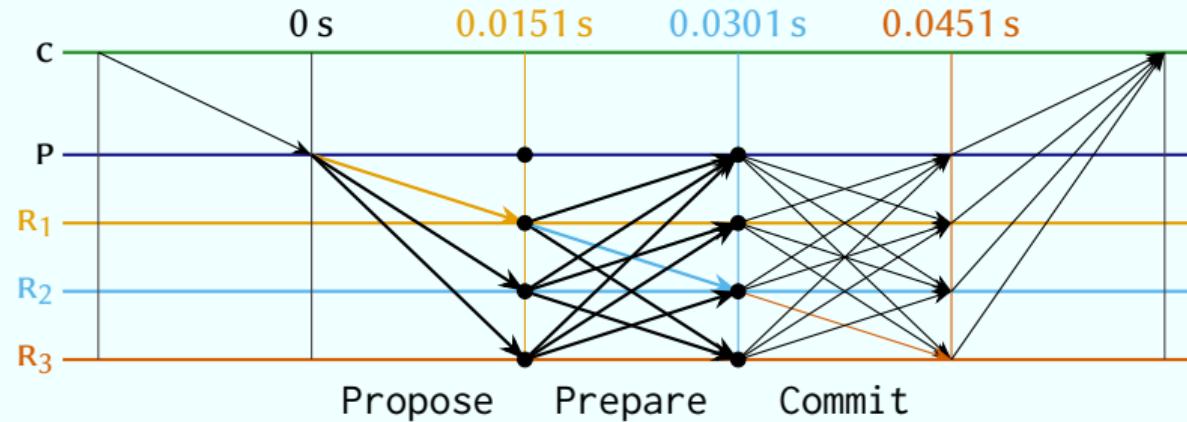


$$\Delta_{\text{PBFT}} = \frac{(\mathbf{n} - 1)s_t}{B} + \delta + \frac{(\mathbf{n} - 1)s_m}{B}$$

# The Single-Round Cost of PBFT (Sketch)

Assumption: Network bandwidth  $B = 100 \text{ MiB/s}$  and delay  $\delta = 15 \text{ ms}$

Propose:  $s_t = 4048 \text{ B}$  each. Prepare and Commit:  $s_m = 256 \text{ B}$  each.

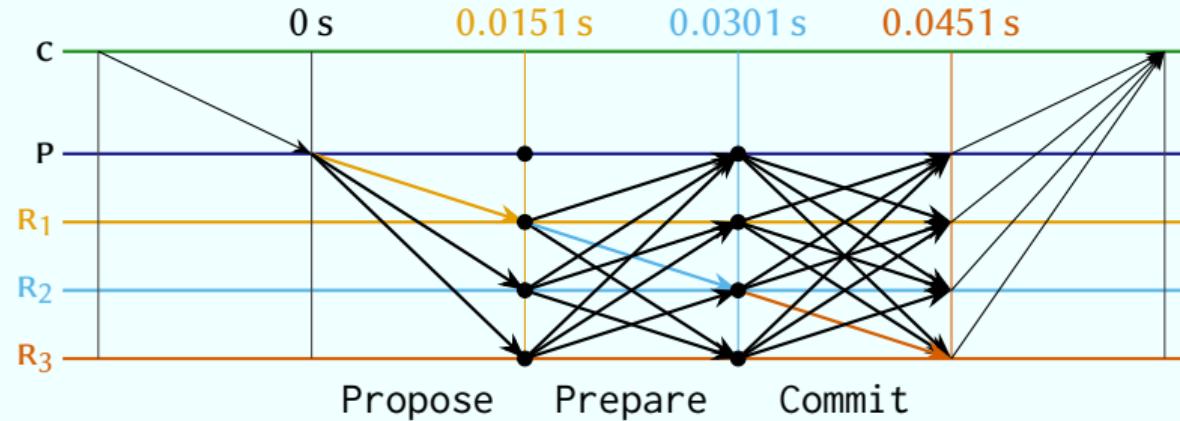


$$\Delta_{\text{PBFT}} = \frac{(\mathbf{n} - 1)s_t}{B} + \delta + \frac{(\mathbf{n} - 1)s_m}{B} + \delta$$

# The Single-Round Cost of PBFT (Sketch)

Assumption: Network bandwidth  $B = 100 \text{ MiB/s}$  and delay  $\delta = 15 \text{ ms}$

Propose:  $s_t = 4048 \text{ B}$  each. Prepare and Commit:  $s_m = 256 \text{ B}$  each.

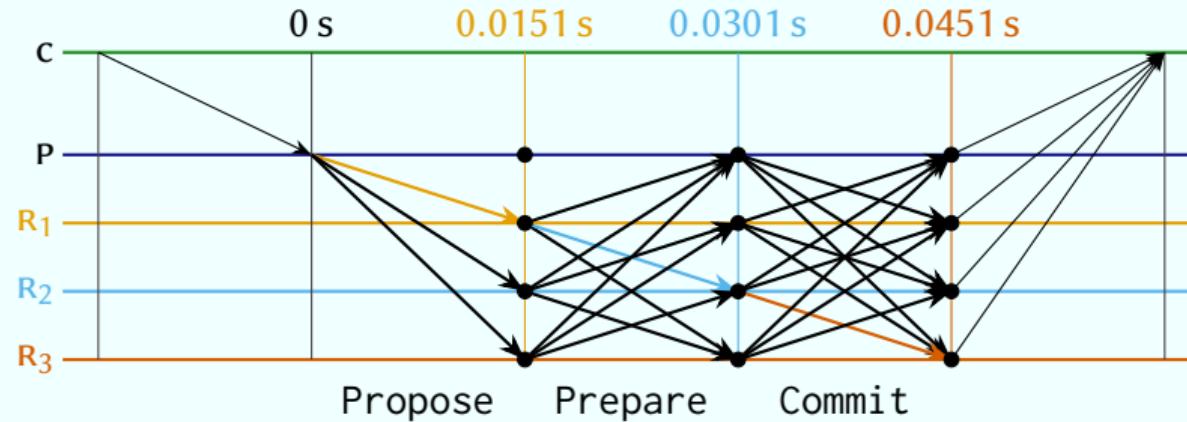


$$\Delta_{\text{PBFT}} = \frac{(\mathbf{n} - 1)s_t}{B} + \delta + \frac{(\mathbf{n} - 1)s_m}{B} + \delta + \frac{(\mathbf{n} - 1)s_m}{B}$$

# The Single-Round Cost of PBFT (Sketch)

Assumption: Network bandwidth  $B = 100 \text{ MiB/s}$  and delay  $\delta = 15 \text{ ms}$

Propose:  $s_t = 4048 \text{ B}$  each. Prepare and Commit:  $s_m = 256 \text{ B}$  each.

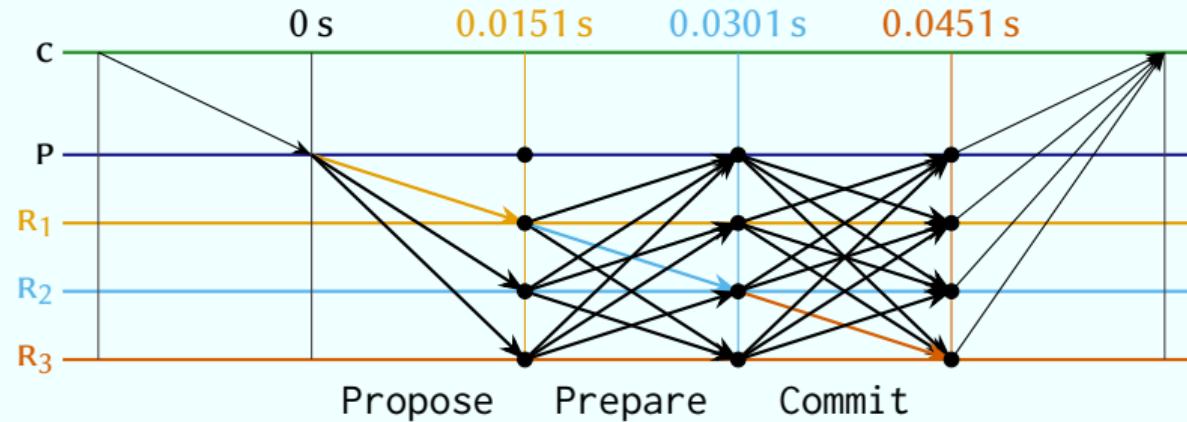


$$\Delta_{\text{PBFT}} = \frac{(\mathbf{n} - 1)s_t}{B} + \delta + \frac{(\mathbf{n} - 1)s_m}{B} + \delta + \frac{(\mathbf{n} - 1)s_m}{B} + \delta$$

# The Single-Round Cost of PBFT (Sketch)

Assumption: Network bandwidth  $B = 100 \text{ MiB/s}$  and delay  $\delta = 15 \text{ ms}$

Propose:  $s_t = 4048 \text{ B}$  each. Prepare and Commit:  $s_m = 256 \text{ B}$  each.

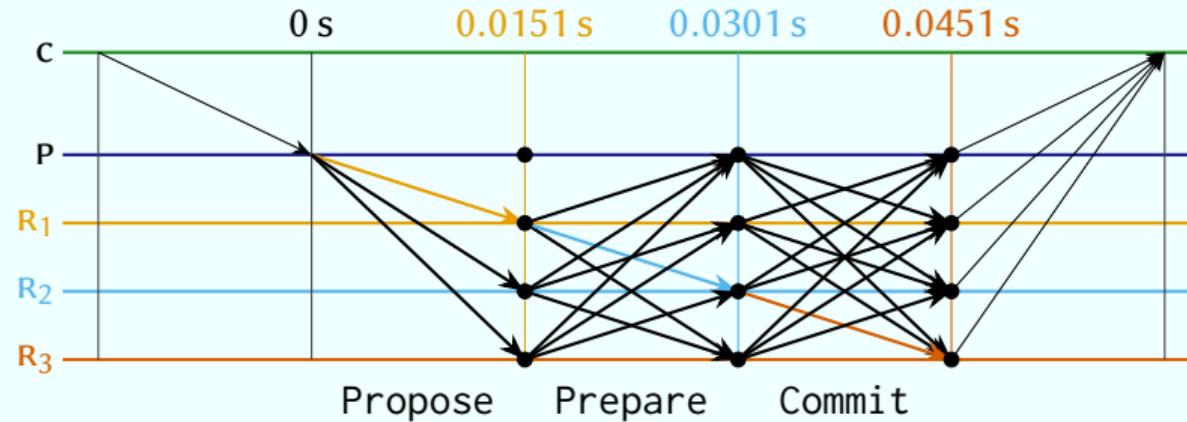


$$\begin{aligned}\Delta_{\text{PBFT}} &= \frac{(\mathbf{n} - 1)s_t}{B} + \delta + \frac{(\mathbf{n} - 1)s_m}{B} + \delta + \frac{(\mathbf{n} - 1)s_m}{B} + \delta \\ &= \frac{(\mathbf{n} - 1)s_t + 2(\mathbf{n} - 1)s_m}{B} + 3\delta\end{aligned}$$

# The Single-Round Cost of PBFT (Sketch)

Assumption: Network bandwidth  $B = 100 \text{ MiB/s}$  and delay  $\delta = 15 \text{ ms}$

Propose:  $s_t = 4048 \text{ B}$  each. Prepare and Commit:  $s_m = 256 \text{ B}$  each.



$$\begin{aligned}\Delta_{\text{PBFT}} &= \frac{(\mathbf{n} - 1)s_t}{B} + \delta + \frac{(\mathbf{n} - 1)s_m}{B} + \delta + \frac{(\mathbf{n} - 1)s_m}{B} + \delta \\ &= \frac{(\mathbf{n} - 1)s_t + 2(\mathbf{n} - 1)s_m}{B} + 3\delta \\ &\approx 3\delta \text{ (assuming } \textit{high delay relative to bandwidth}).\end{aligned}$$

## The Throughput of PBFT

Sequential: Next consensus round starts after finishing the current round

$$T_{\text{PBFT}} = \frac{1}{\Delta_{\text{PBFT}}} = \frac{B}{(\mathbf{n} - 1)s_t + 2(\mathbf{n} - 1)s_m + 3B\delta}.$$

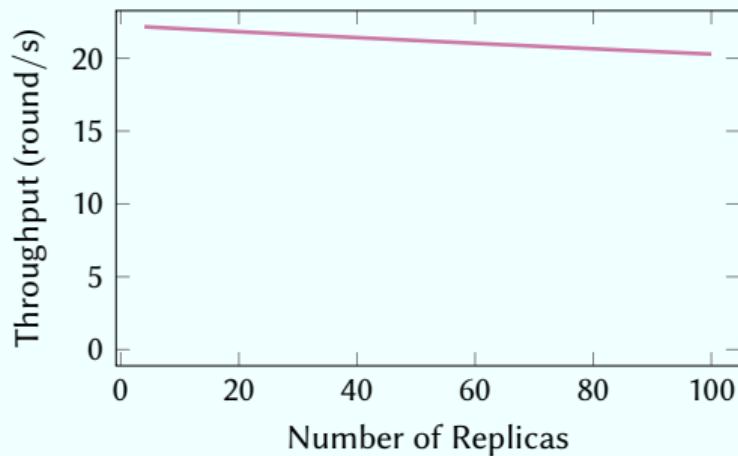
# The Throughput of PBFT

Sequential: Next consensus round starts after finishing the current round

$$T_{\text{PBFT}} = \frac{1}{\Delta_{\text{PBFT}}} = \frac{B}{(\mathbf{n} - 1)s_t + 2(\mathbf{n} - 1)s_m + 3B\delta}.$$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$

$(\delta = 15 \text{ ms})$

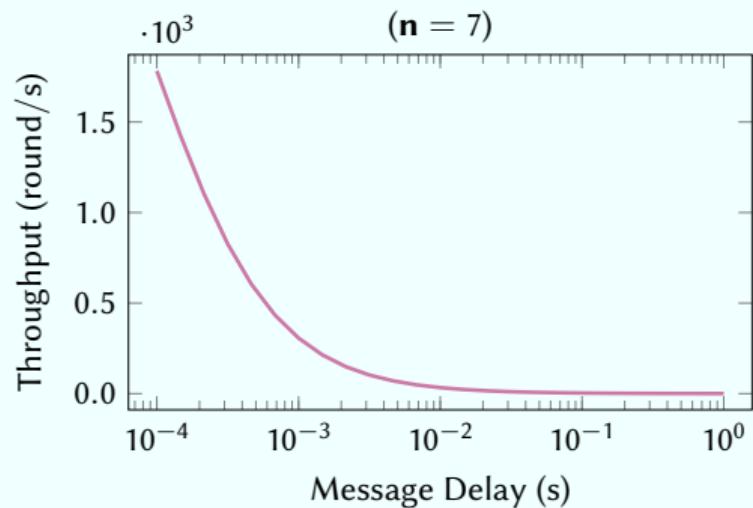
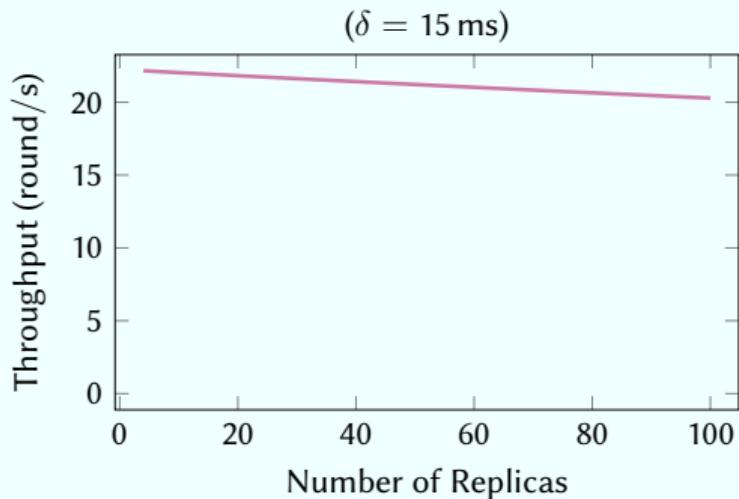


# The Throughput of PBFT

Sequential: Next consensus round starts after finishing the current round

$$T_{\text{PBFT}} = \frac{1}{\Delta_{\text{PBFT}}} = \frac{B}{(\mathbf{n} - 1)s_t + 2(\mathbf{n} - 1)s_m + 3B\delta}.$$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$



## Implementation techniques for PBFT

Realistic wide-area message delays: 10 ms–300 ms

The throughput  $T_{\text{PBFT}}$  of sequential PBFT is *impractically low*.

# Implementation techniques for PBFT

Realistic wide-area message delays: 10 ms–300 ms

The throughput  $T_{\text{PBFT}}$  of sequential PBFT is *impractically low*.

## Fine-tuning PBFT implementations

Batching many transactions per consensus decision.

Out-of-order processing many consensus decisions at the same time.

Overlapping phases of consecutive rounds.

# Batching Client Requests

The cost of a single round of PBFT

Message	Sent by	Size
Propose	Primary	$s_t$
Prepare	Backups	$s_m$
Commit	All	$s_m$

# Batching Client Requests

The cost of a single round of PBFT

Batching: each decision is on **m** transactions.

Message	Sent by	Size	(batch)
Propose	Primary	$s_t$	$ms_t$
Prepare	Backups	$s_m$	$s_m$
Commit	All	$s_m$	$s_m$

# Batching Client Requests

The cost of a single round of PBFT

Batching: each decision is on  $m$  transactions.

Message	Sent by	Size	(batch)
Propose	Primary	$s_t$	$ms_t$
Prepare	Backups	$s_m$	$s_m$
Commit	All	$s_m$	$s_m$
Total:	$2n(n - 1)$	$\mathcal{O}(s_t n + s_m n^2)$	$\mathcal{O}(ms_t n + s_m n^2)$

# The Throughput of PBFT with Batching

Sequential, batching  $m$  transactions per consensus round

$$\Delta_{\text{PBFT-}m} = \frac{m(\mathbf{n} - 1)s_t + 2(\mathbf{n} - 1)s_m}{B} + 3\delta;$$

## The Throughput of PBFT with Batching

Sequential, batching  $m$  transactions per consensus round

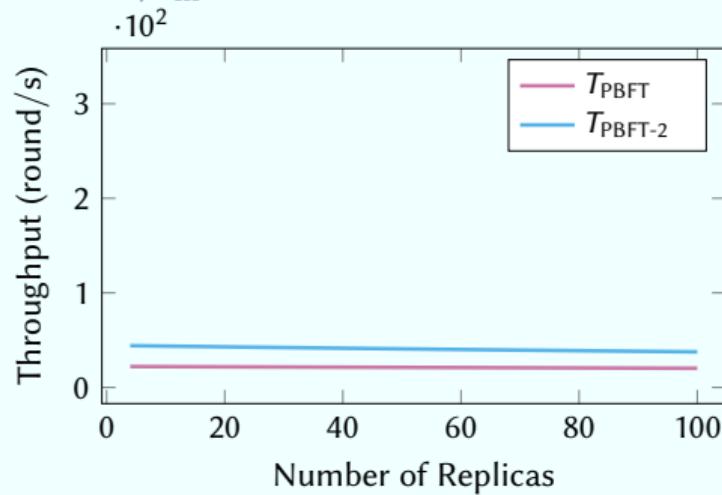
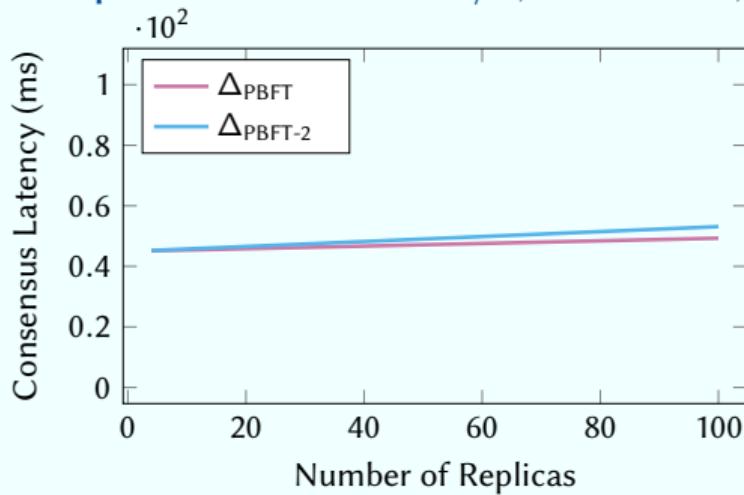
$$\Delta_{\text{PBFT}-m} = \frac{m(\mathbf{n} - 1)s_t + 2(\mathbf{n} - 1)s_m}{B} + 3\delta; \quad T_{\text{PBFT}-m} = \frac{m}{\Delta_{\text{PBFT}-m}}.$$

# The Throughput of PBFT with Batching

Sequential, batching  $m$  transactions per consensus round

$$\Delta_{\text{PBFT}-m} = \frac{m(\mathbf{n} - 1)s_t + 2(\mathbf{n} - 1)s_m}{B} + 3\delta; \quad T_{\text{PBFT}-m} = \frac{m}{\Delta_{\text{PBFT}-m}}.$$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $\delta = 15 \text{ ms}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$

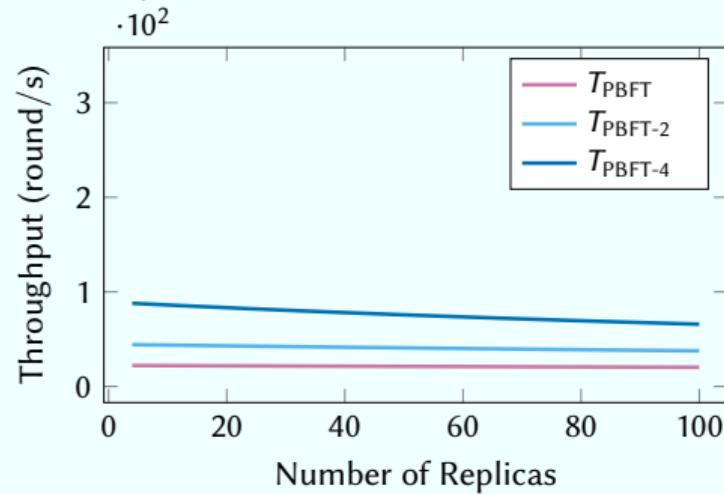
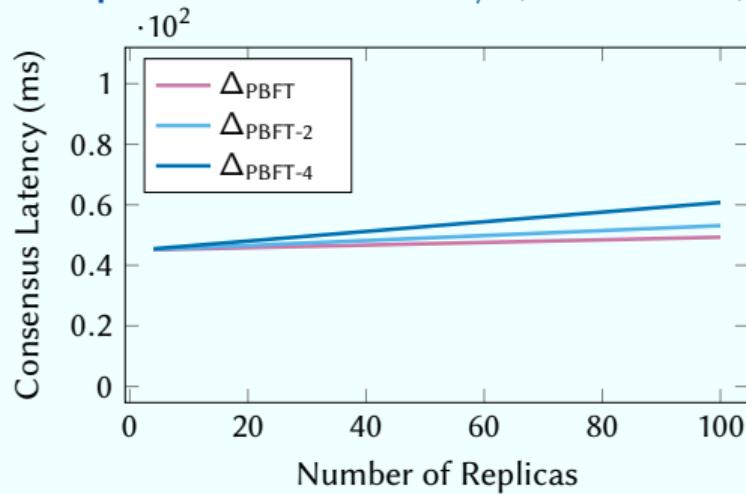


# The Throughput of PBFT with Batching

Sequential, batching  $m$  transactions per consensus round

$$\Delta_{\text{PBFT}-m} = \frac{m(\mathbf{n} - 1)s_t + 2(\mathbf{n} - 1)s_m}{B} + 3\delta; \quad T_{\text{PBFT}-m} = \frac{m}{\Delta_{\text{PBFT}-m}}.$$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $\delta = 15 \text{ ms}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$

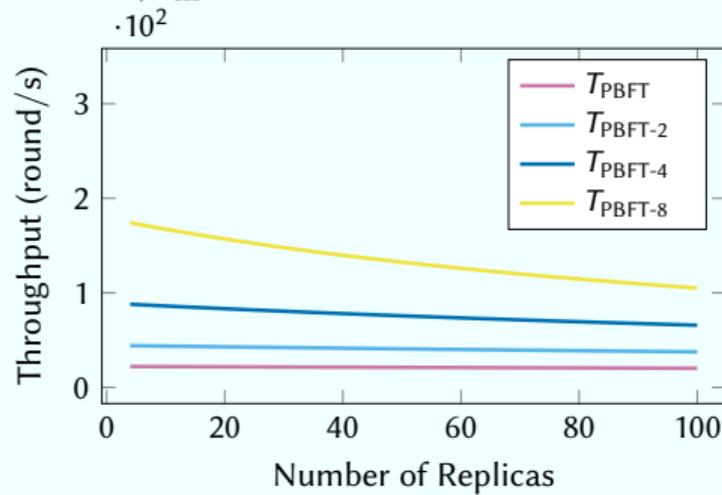
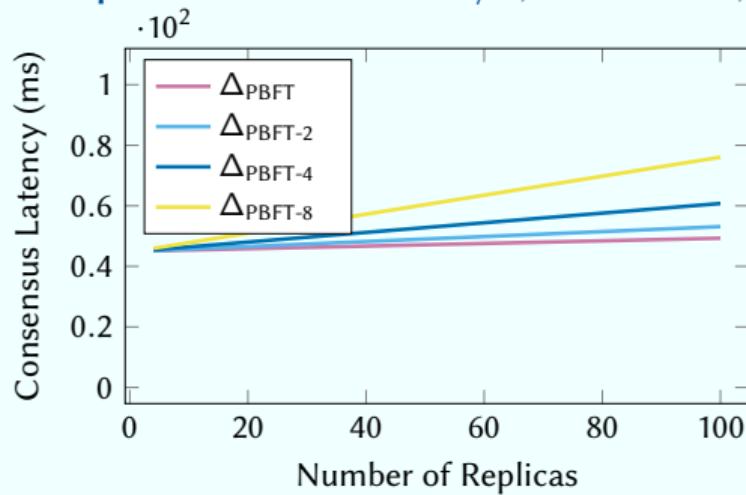


# The Throughput of PBFT with Batching

Sequential, batching  $m$  transactions per consensus round

$$\Delta_{\text{PBFT}-m} = \frac{m(\mathbf{n} - 1)s_t + 2(\mathbf{n} - 1)s_m}{B} + 3\delta; \quad T_{\text{PBFT}-m} = \frac{m}{\Delta_{\text{PBFT}-m}}.$$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $\delta = 15 \text{ ms}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$

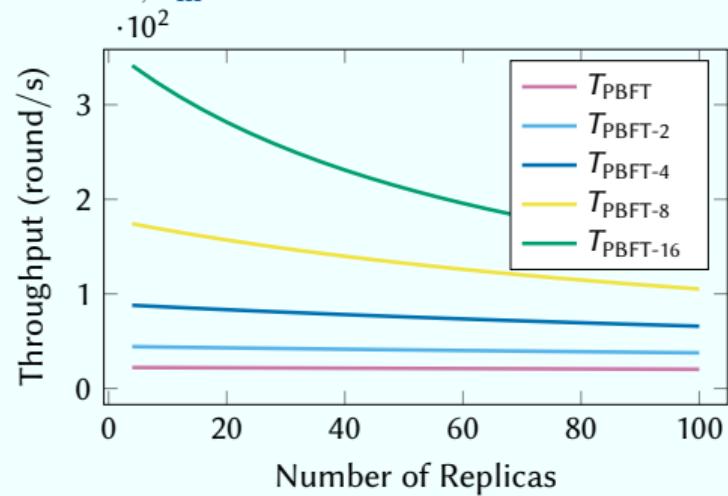
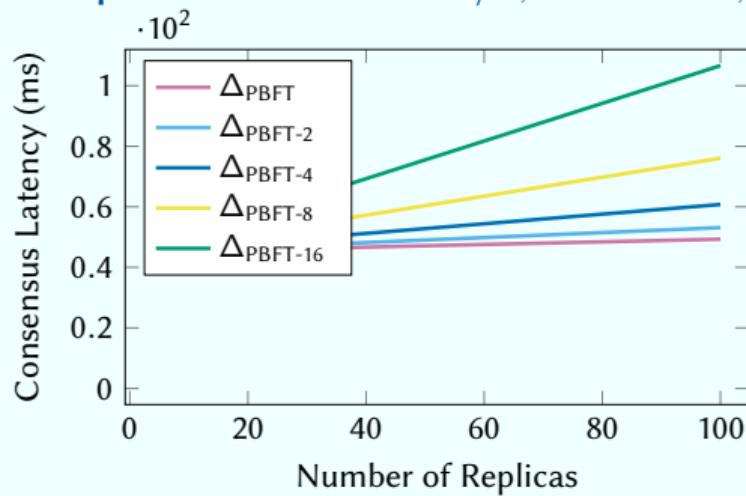


# The Throughput of PBFT with Batching

Sequential, batching  $m$  transactions per consensus round

$$\Delta_{\text{PBFT}-m} = \frac{m(\mathbf{n} - 1)s_t + 2(\mathbf{n} - 1)s_m}{B} + 3\delta; \quad T_{\text{PBFT}-m} = \frac{m}{\Delta_{\text{PBFT}-m}}.$$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $\delta = 15 \text{ ms}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$

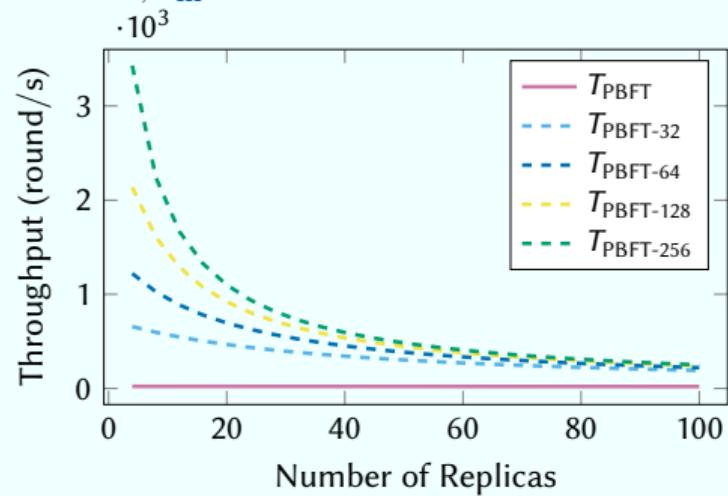
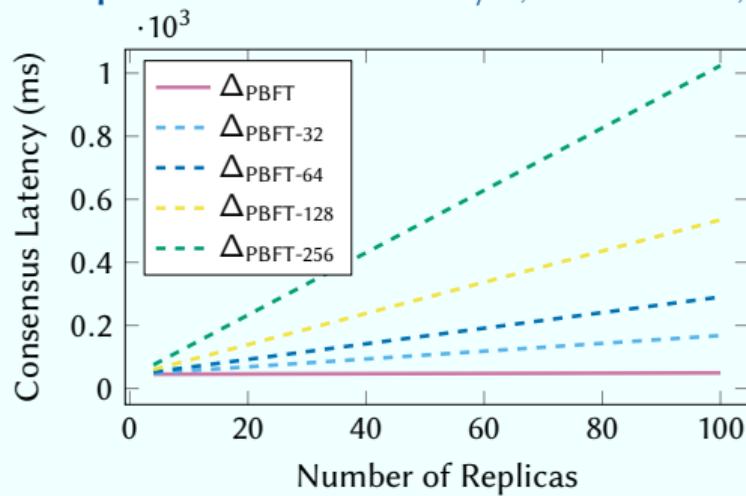


# The Throughput of PBFT with Batching

Sequential, batching  $m$  transactions per consensus round

$$\Delta_{\text{PBFT}-m} = \frac{m(\mathbf{n} - 1)s_t + 2(\mathbf{n} - 1)s_m}{B} + 3\delta; \quad T_{\text{PBFT}-m} = \frac{m}{\Delta_{\text{PBFT}-m}}.$$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $\delta = 15 \text{ ms}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$



## Using Batching to Improve Throughput Scalability

	Messages <i>(per trans.)</i>	Size <i>(per trans.)</i>
PBFT	$2n(n - 1)$	$O(s_t n + s_m n^2)$

## Using Batching to Improve Throughput Scalability

	Messages <i>(per trans.)</i>	Size <i>(per trans.)</i>
PBFT	$2n(n - 1)$	$O(s_t n + s_m n^2)$

## Using Batching to Improve Throughput Scalability

	Messages <i>(per trans.)</i>	Size <i>(per trans.)</i>	
PBFT	$2n(n - 1)$	$2n(n - 1)$	$\mathcal{O}(s_t n + s_m n^2)$
PBFT-n	$2n(n - 1)$	$2(n - 1)$	$\mathcal{O}(ns_t n + s_m n^2)$

## Using Batching to Improve Throughput Scalability

	Messages <i>(per trans.)</i>	Size <i>(per trans.)</i>	
PBFT	$2n(n - 1)$	$2n(n - 1)$	$\mathcal{O}(s_t n + s_m n^2)$
PBFT-n	$2n(n - 1)$	$2(n - 1)$	$\mathcal{O}(ns_t n + s_m n^2)$

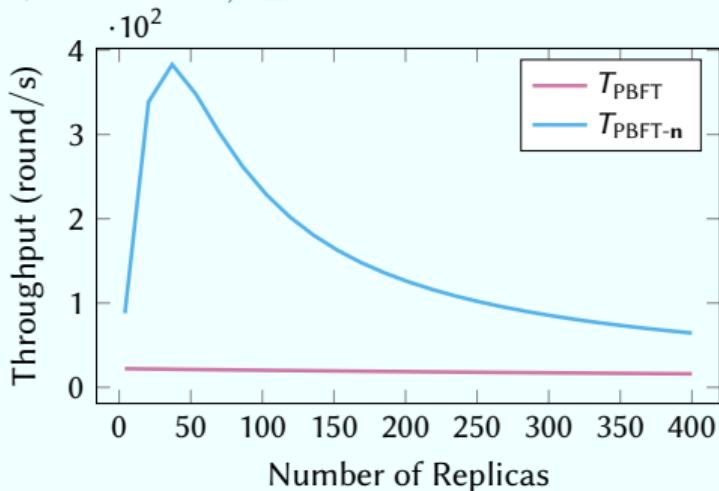
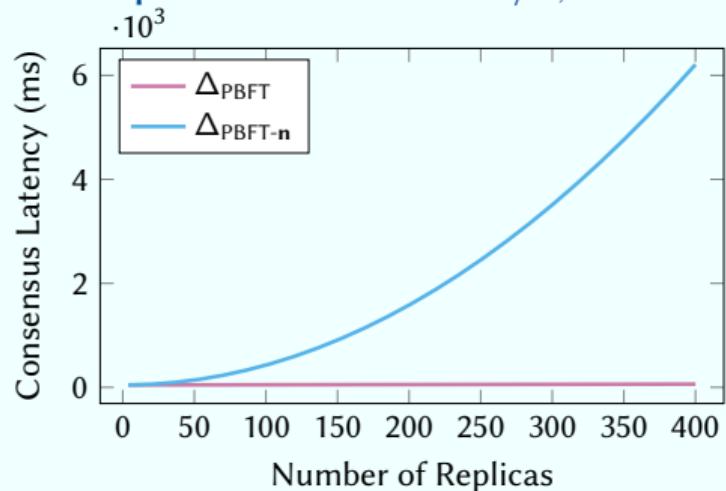
*Quadratic*

*Linear*

# Using Batching to Improve Throughput Scalability

	Messages <i>(per trans.)</i>	Size <i>(per trans.)</i>	
PBFT	$2n(n - 1)$	$2n(n - 1)$	$\mathcal{O}(s_t n + s_m n^2)$
PBFT-n	$2n(n - 1)$	$2(n - 1)$	$\mathcal{O}(ns_t n + s_m n^2)$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $\delta = 15 \text{ ms}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$



## Resource Utilization of Sequential PBFT

Assumption:  $n = 4$ ,  $B = 100 \text{ MiB/s}$ ,  $\delta = 15 \text{ ms}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$

$$\Delta_{\text{PBFT}} = \frac{(n - 1)s_t + 2(n - 1)s_m}{B} + 3\delta \approx 45.1 \text{ ms.}$$

## Resource Utilization of Sequential PBFT

Assumption:  $n = 4$ ,  $B = 100 \text{ MiB/s}$ ,  $\delta = 15 \text{ ms}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$

$$\Delta_{\text{PBFT}} = \underbrace{\frac{(n-1)s_t + 2(n-1)s_m}{B}}_{\begin{array}{c} \textit{message transfer} \\ \approx 0.1 \text{ ms} \end{array}} + 3\delta \approx 45.1 \text{ ms.}$$

## Resource Utilization of Sequential PBFT

Assumption:  $n = 4$ ,  $B = 100 \text{ MiB/s}$ ,  $\delta = 15 \text{ ms}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$

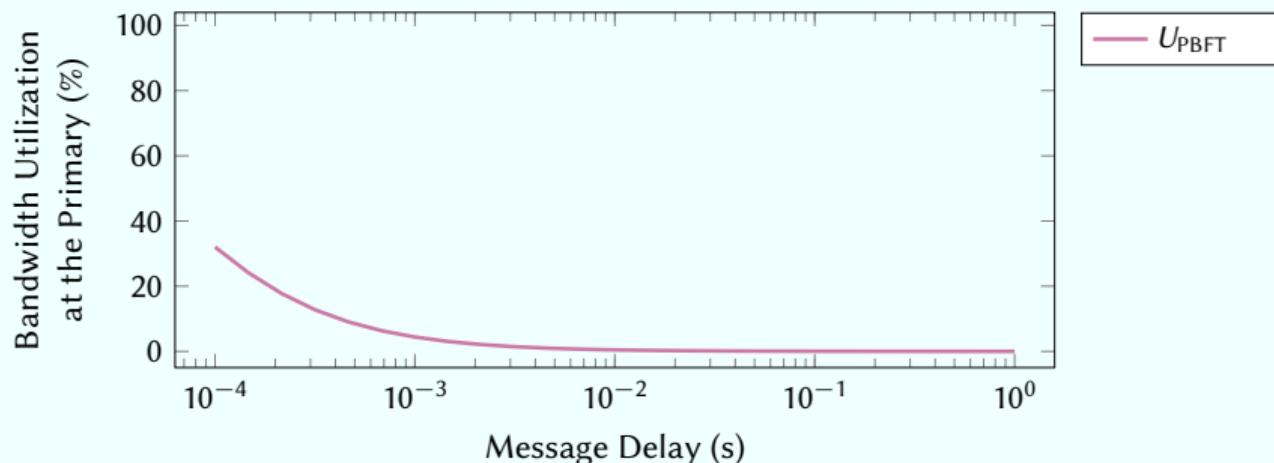
$$\Delta_{\text{PBFT}} = \underbrace{\frac{(n-1)s_t + 2(n-1)s_m}{B}}_{\begin{array}{c} \text{message transfer} \\ \approx 0.1 \text{ ms} \end{array}} + \underbrace{3\delta}_{\begin{array}{c} \text{waiting} \\ 45.0 \text{ ms} \end{array}} \approx 45.1 \text{ ms.}$$

# Resource Utilization of Sequential PBFT

Assumption:  $n = 4$ ,  $B = 100 \text{ MiB/s}$ ,  $\delta = 15 \text{ ms}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$

$$\Delta_{\text{PBFT}} = \underbrace{\frac{(n-1)s_t + 2(n-1)s_m}{B}}_{\text{message transfer}} + \underbrace{3\delta}_{\text{waiting}} \approx 45.1 \text{ ms.}$$

$\approx 0.1 \text{ ms}$        $45.0 \text{ ms}$



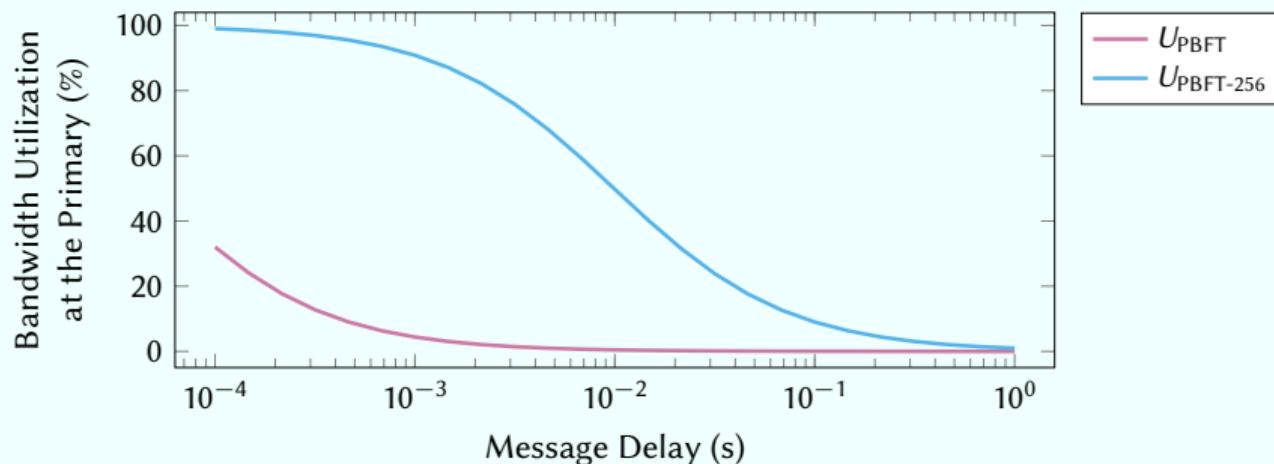
# Resource Utilization of Sequential PBFT

Assumption:  $n = 4$ ,  $B = 100 \text{ MiB/s}$ ,  $\delta = 15 \text{ ms}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$

$$\Delta_{\text{PBFT}} = \underbrace{\frac{(n-1)s_t + 2(n-1)s_m}{B}}_{\text{message transfer}} + 3\delta \approx 45.1 \text{ ms.}$$

$\approx 0.1 \text{ ms}$        $45.0 \text{ ms}$

*message transfer*      *waiting*



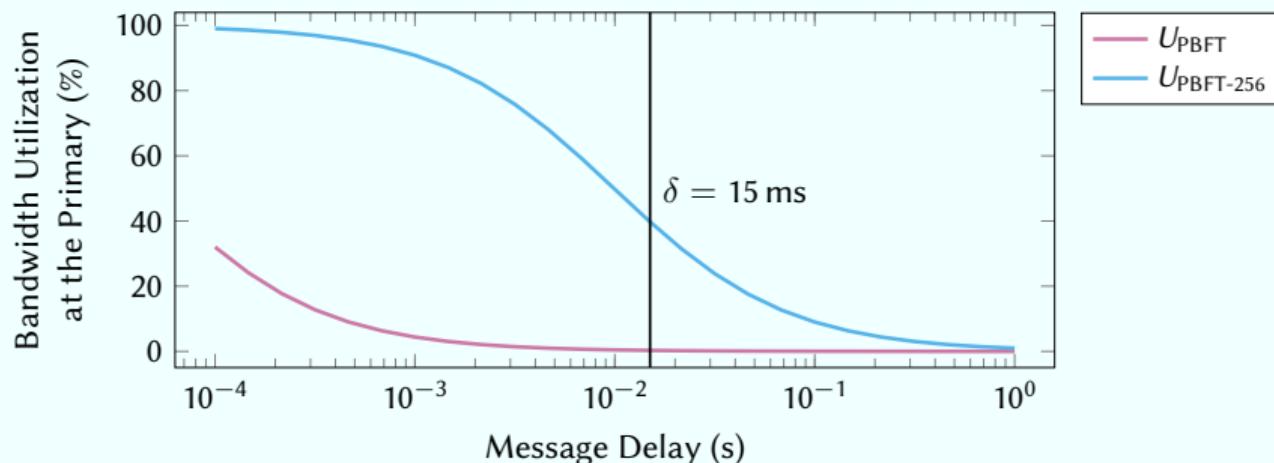
# Resource Utilization of Sequential PBFT

Assumption:  $n = 4$ ,  $B = 100 \text{ MiB/s}$ ,  $\delta = 15 \text{ ms}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$

$$\Delta_{\text{PBFT}} = \underbrace{\frac{(n-1)s_t + 2(n-1)s_m}{B}}_{\text{message transfer}} + 3\delta \approx 45.1 \text{ ms.}$$

$\approx 0.1 \text{ ms}$        $\approx 45.0 \text{ ms}$

*message transfer*      *waiting*



## Problem: Resource Utilization of PBFT

- ▶ Typically *less than 5%* bandwidth utilization at primary.
- ▶ With huge batches still *less than 40%* bandwidth utilization at primary.

## Problem: Resource Utilization of PBFT

- ▶ Typically *less than 5%* bandwidth utilization at primary.
- ▶ With huge batches still *less than 40%* bandwidth utilization at primary.

To maximize throughput: use *all* bandwidth at the primary.

## Problem: Resource Utilization of PBFT

- ▶ Typically *less than 5%* bandwidth utilization at primary.
- ▶ With huge batches still *less than 40%* bandwidth utilization at primary.

To maximize throughput: use *all* bandwidth at the primary.

## Out-of-order processing

Primary can proposes *future rounds* before current rounds are finished.

## Problem: Resource Utilization of PBFT

- ▶ Typically *less than 5%* bandwidth utilization at primary.
- ▶ With huge batches still *less than 40%* bandwidth utilization at primary.

To maximize throughput: use *all* bandwidth at the primary.

## Out-of-order processing

Primary can propose *future rounds* before current rounds are finished.

## Practical challenges

- ▶ Memory usage: replicas maintain meta-data for each *active* round.
- ▶ Byzantine behavior: exhaust the set of round numbers.

## Problem: Resource Utilization of PBFT

- ▶ Typically *less than 5%* bandwidth utilization at primary.
- ▶ With huge batches still *less than 40%* bandwidth utilization at primary.

To maximize throughput: use *all* bandwidth at the primary.

## Out-of-order processing

Primary can propose *future rounds* before current rounds are finished.

## Practical challenges

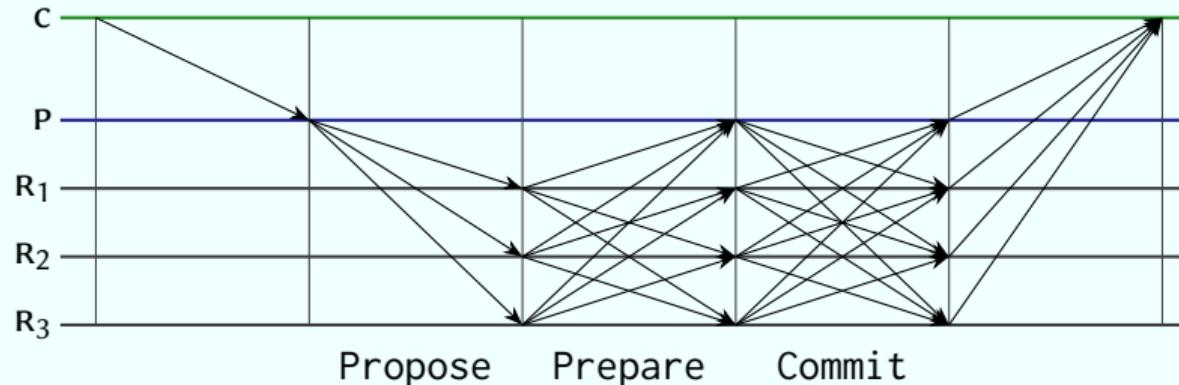
- ▶ Memory usage: replicas maintain meta-data for each *active* round.
- ▶ Byzantine behavior: exhaust the set of round numbers.

Limit proposals to an *active window* of valid rounds.

E.g., only proposals in 1000 rounds after the last finished round.

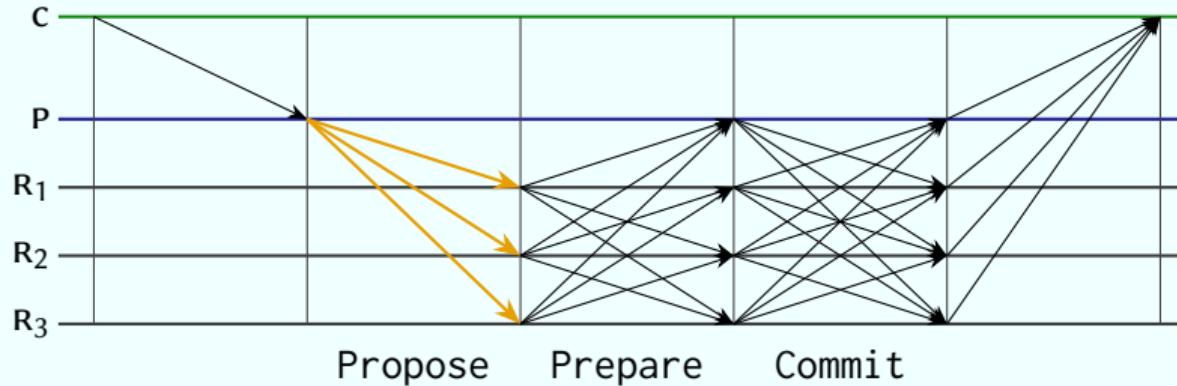
# The Single-Round Cost of PBFT (revised)

Assumption: Primary does most work ( $s_t > s_m$ )



# The Single-Round Cost of PBFT (revised)

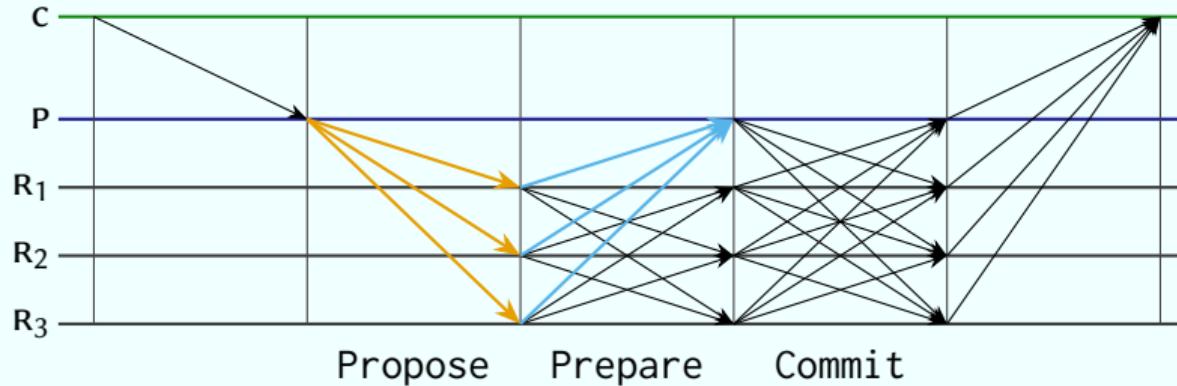
Assumption: Primary does most work ( $s_t > s_m$ )



- ▶ Send  $n - 1$  messages
- ▶  $s_t$  B each

# The Single-Round Cost of PBFT (revised)

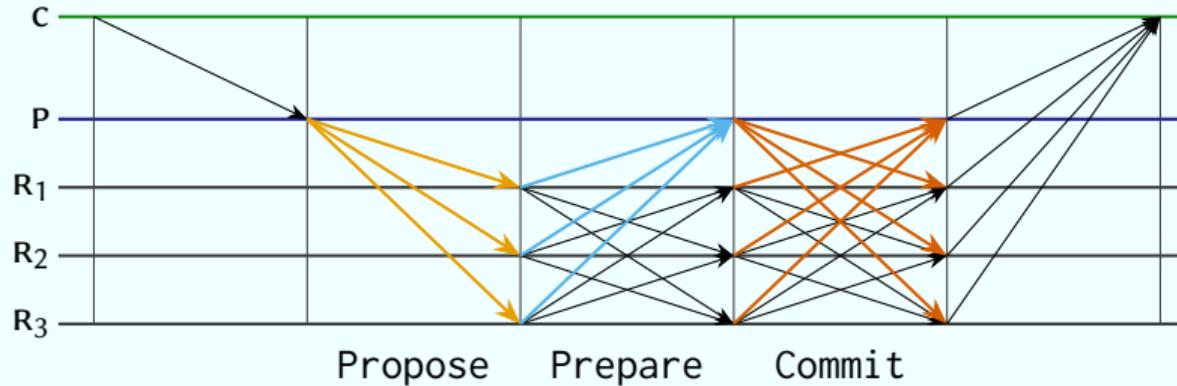
Assumption: Primary does most work ( $s_t > s_m$ )



- ▶ Receive  $n - 1$  messages
- ▶  $s_m$  B each

# The Single-Round Cost of PBFT (revised)

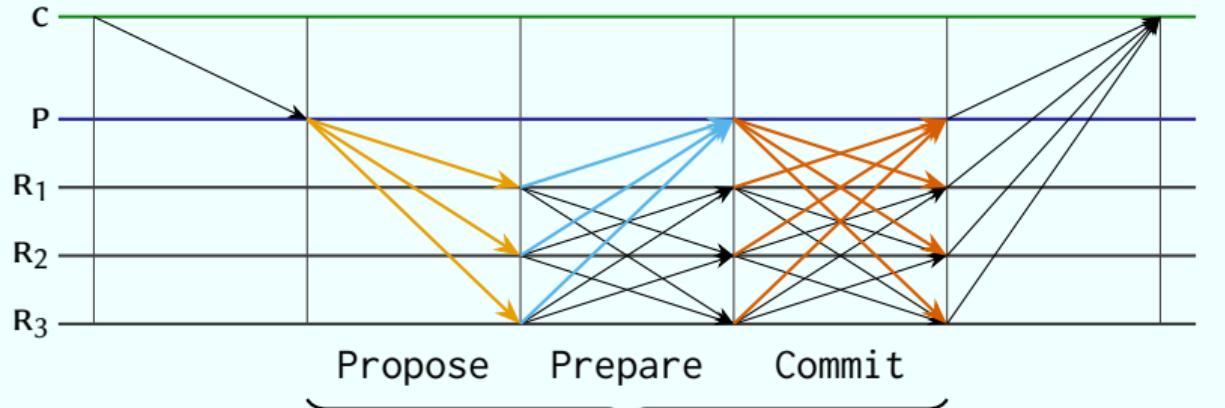
Assumption: Primary does most work ( $s_t > s_m$ )



- ▶ Send and receive  $n - 1$  messages
- ▶  $s_m$  B each

# The Single-Round Cost of PBFT (revised)

Assumption: Primary does most work ( $s_t > s_m$ )



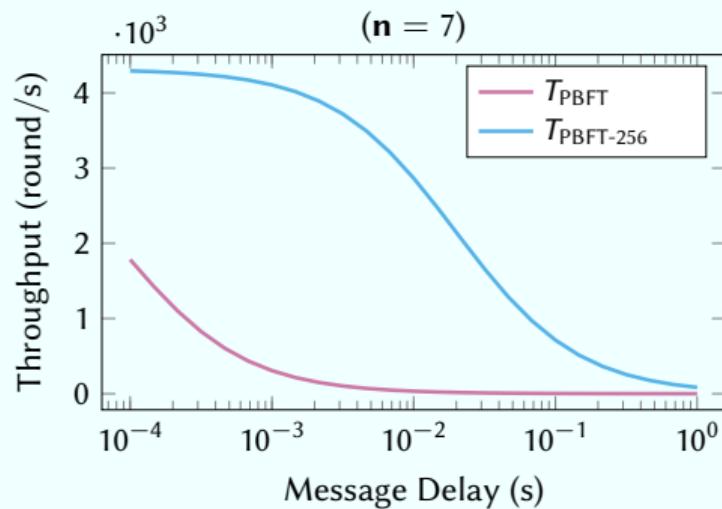
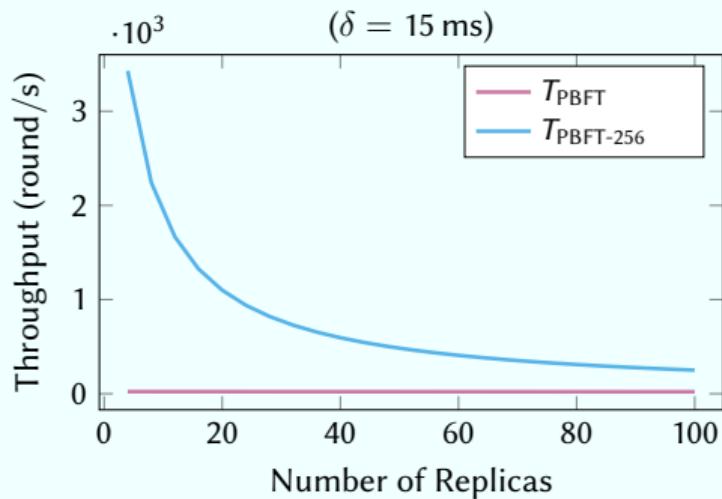
$$(n - 1)s_t + (n - 1)s_m + 2(n - 1)s_m = (n - 1)(s_t + 3s_m) \text{ B/round.}$$

# The Out-of-Order Throughput of PBFT

Assumption: Primary does most work ( $s_t > s_m$ )

$$T_{\text{ooo-PBFT}} = \frac{B}{(n - 1)(s_t + 3s_m)}.$$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$

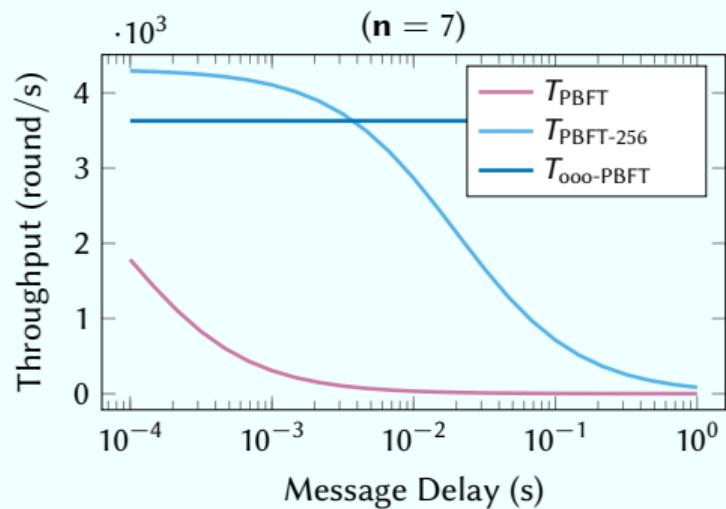
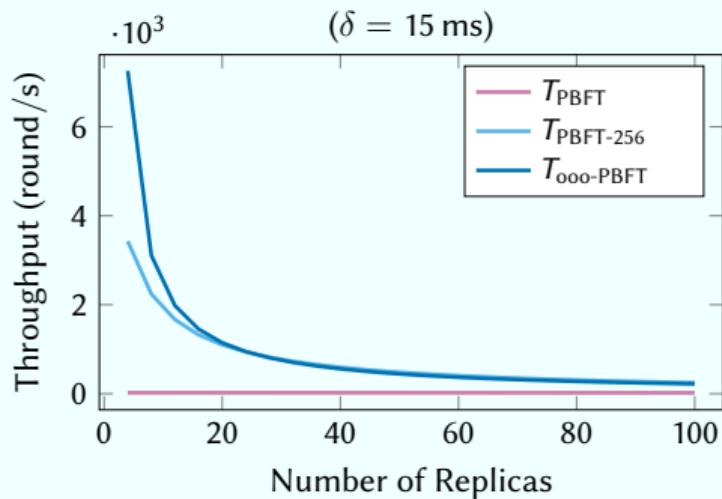


# The Out-of-Order Throughput of PBFT

Assumption: Primary does most work ( $s_t > s_m$ )

$$T_{\text{ooo-PBFT}} = \frac{B}{(n - 1)(s_t + 3s_m)}.$$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$

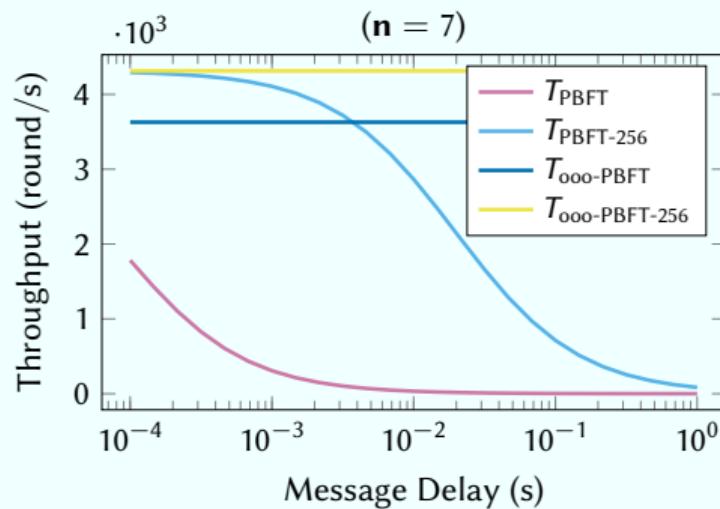
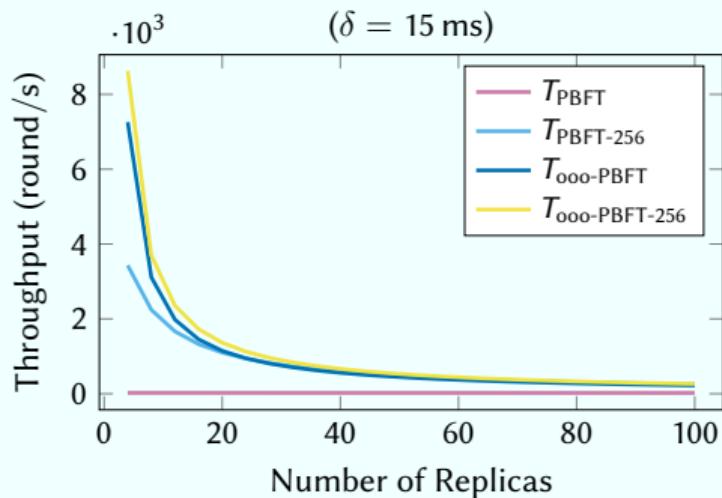


# The Out-of-Order Throughput of PBFT

Assumption: Primary does most work ( $s_t > s_m$ )

$$T_{\text{ooo-PBFT}} = \frac{B}{(\mathbf{n} - 1)(s_t + 3s_m)}; \quad T_{\text{ooo-PBFT-}m} = \frac{mB}{(\mathbf{n} - 1)(ms_t + 3s_m)}.$$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$



## Overlapping Communication Phases

Out-of-order processing is *complex* to implement.

## Overlapping Communication Phases

Out-of-order processing is *complex* to implement.

Consider a backup replica R.

- ▶ Last step of round  $\rho$ : Commit messages.
- ▶ First step of round  $\rho + 1$ : Prepare messages.

## Overlapping Communication Phases

Out-of-order processing is *complex* to implement.

Consider a backup replica R.

- ▶ Last step of round  $\rho$ : Commit messages.
- ▶ First step of round  $\rho + 1$ : Prepare messages.

Idea: Overlapping communication phases

Merge Commit message of  $\rho$  with the Prepare message of  $\rho + 1$ .

# Overlapping Communication Phases

Out-of-order processing is *complex* to implement.

Consider a backup replica R.

- ▶ Last step of round  $\rho$ : Commit messages.
- ▶ First step of round  $\rho + 1$ : Prepare messages.

## Idea: Overlapping communication phases

Merge Commit message of  $\rho$  with the Prepare message of  $\rho + 1$ .

- ▶ Make proposal of round  $\rho + 1$  *refer* to round  $\rho$ .
- ▶ Prepare for round  $\rho + 1$  *implies* Commit for round  $\rho$ .
- ▶ Primary proposes round  $\rho + 1$  *after* it finished the prepare-phase for round  $\rho$ .

# Overlapping Communication Phases

Out-of-order processing is *complex* to implement.

Consider a backup replica R.

- ▶ Last step of round  $\rho$ : Commit messages.
- ▶ First step of round  $\rho + 1$ : Prepare messages.

Idea: Overlapping communication phases

Merge Commit message of  $\rho$  with the Prepare message of  $\rho + 1$ .

- ▶ Make proposal of round  $\rho + 1$  *refer* to round  $\rho$ .
- ▶ Prepare for round  $\rho + 1$  *implies* Commit for round  $\rho$ .
- ▶ Primary proposes round  $\rho + 1$  *after* it finished the prepare-phase for round  $\rho$ .

Implies strict consecutive processing of rounds

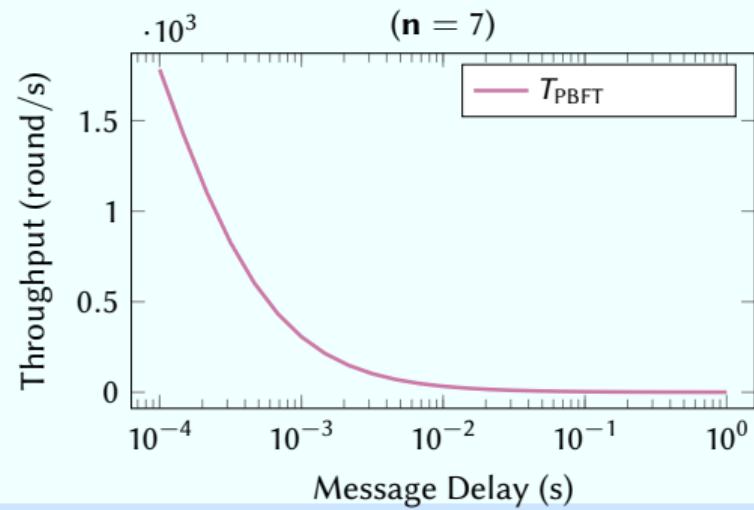
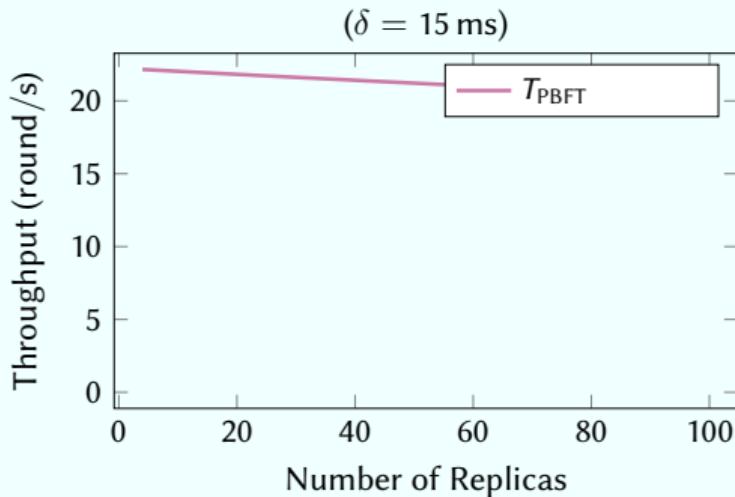
Overlapping *cannot* be combined with out-of-order processing!

# The Single-Round Cost of PBFT with Overlapping

$$\Delta_{\text{PBFT}} = \frac{(\mathbf{n} - 1)s_t + 2(\mathbf{n} - 1)s_m}{B} + 3\delta;$$

$$T_{\text{PBFT}} = \frac{1}{\Delta_{\text{PBFT}}}$$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$



# The Single-Round Cost of PBFT with Overlapping

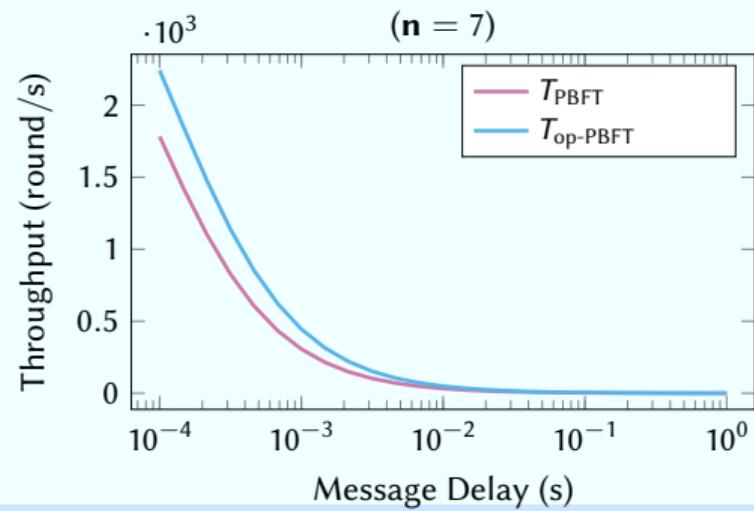
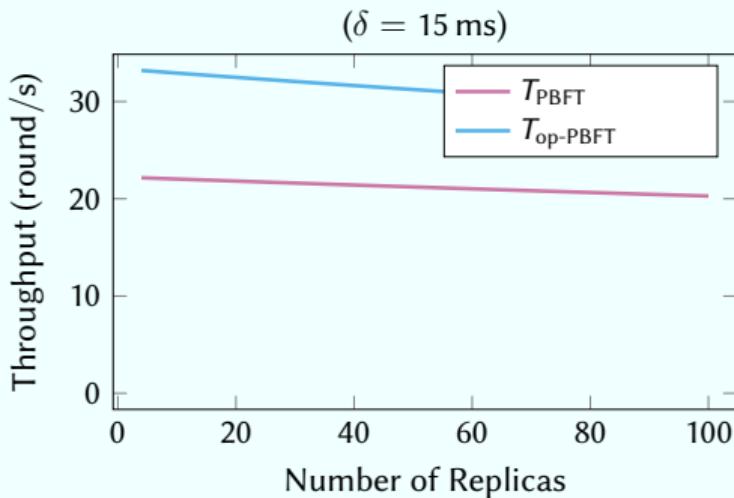
$$\Delta_{\text{PBFT}} = \frac{(\mathbf{n} - 1)s_t + 2(\mathbf{n} - 1)s_m}{B} + 3\delta;$$

$$\Delta_{\text{op-PBFT}} = \frac{(\mathbf{n} - 1)s_t + (\mathbf{n} - 1)s_m}{B} + 2\delta;$$

$$T_{\text{PBFT}} = \frac{1}{\Delta_{\text{PBFT}}}$$

$$T_{\text{op-PBFT}} = \frac{1}{\Delta_{\text{op-PBFT}}}.$$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$



# The Single-Round Cost of PBFT with Overlapping

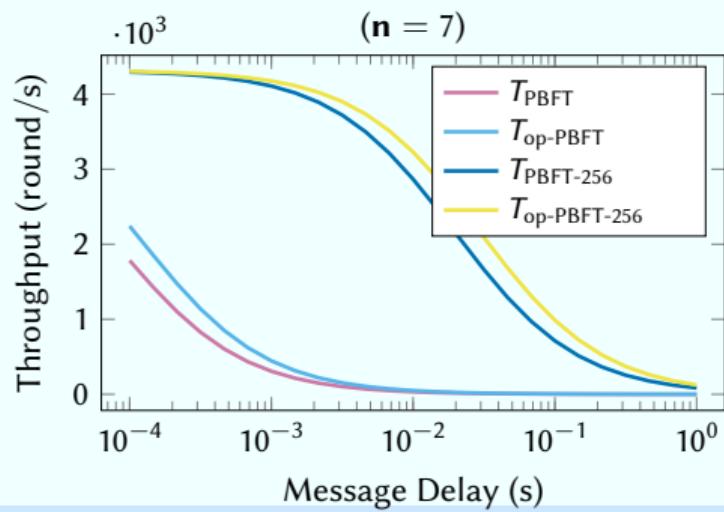
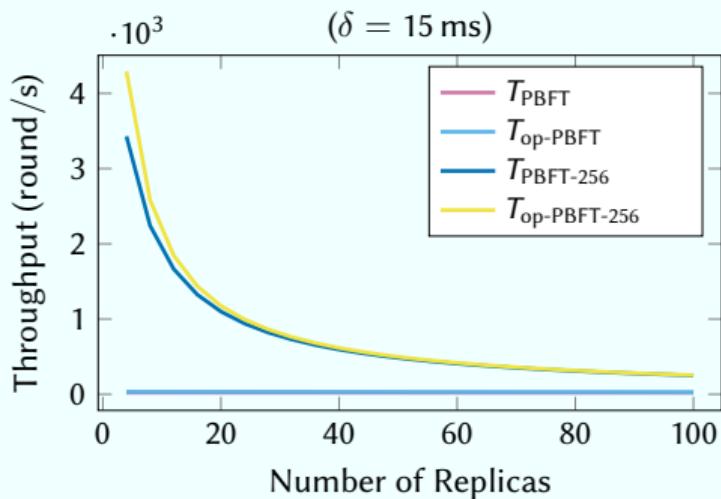
$$\Delta_{\text{PBFT}} = \frac{(\mathbf{n} - 1)s_t + 2(\mathbf{n} - 1)s_m}{B} + 3\delta;$$

$$\Delta_{\text{op-PBFT}} = \frac{(\mathbf{n} - 1)s_t + (\mathbf{n} - 1)s_m}{B} + 2\delta;$$

$$T_{\text{PBFT}} = \frac{1}{\Delta_{\text{PBFT}}}$$

$$T_{\text{op-PBFT}} = \frac{1}{\Delta_{\text{op-PBFT}}}.$$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$



# The Single-Round Cost of PBFT with Overlapping

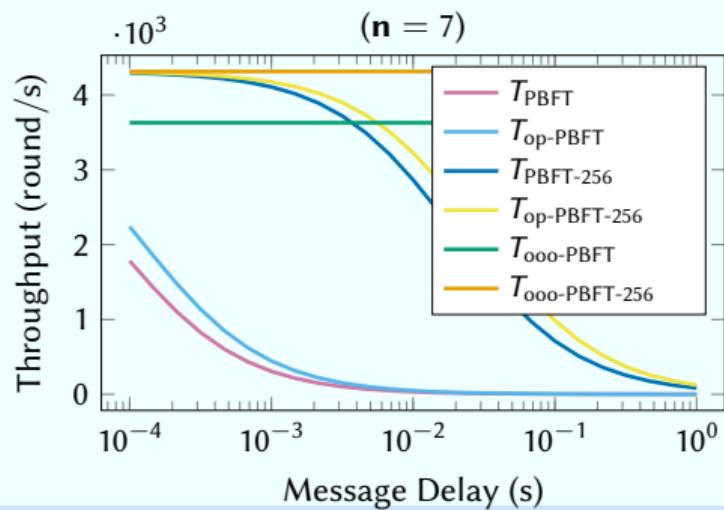
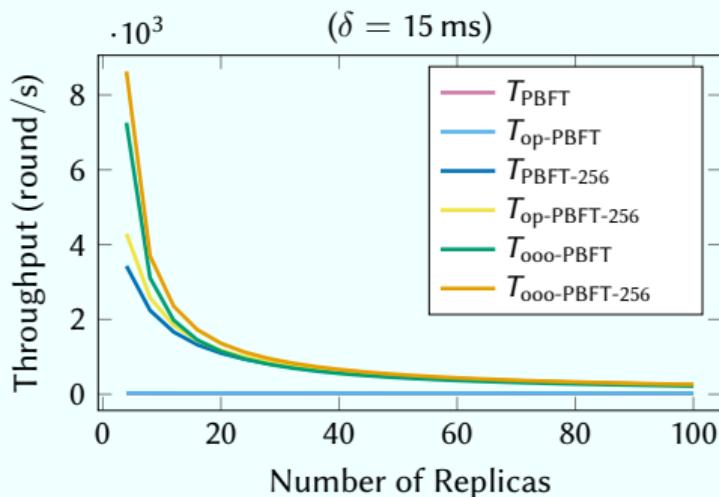
$$\Delta_{\text{PBFT}} = \frac{(\mathbf{n} - 1)s_t + 2(\mathbf{n} - 1)s_m}{B} + 3\delta;$$

$$\Delta_{\text{op-PBFT}} = \frac{(\mathbf{n} - 1)s_t + (\mathbf{n} - 1)s_m}{B} + 2\delta;$$

$$T_{\text{PBFT}} = \frac{1}{\Delta_{\text{PBFT}}}$$

$$T_{\text{op-PBFT}} = \frac{1}{\Delta_{\text{op-PBFT}}}.$$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$



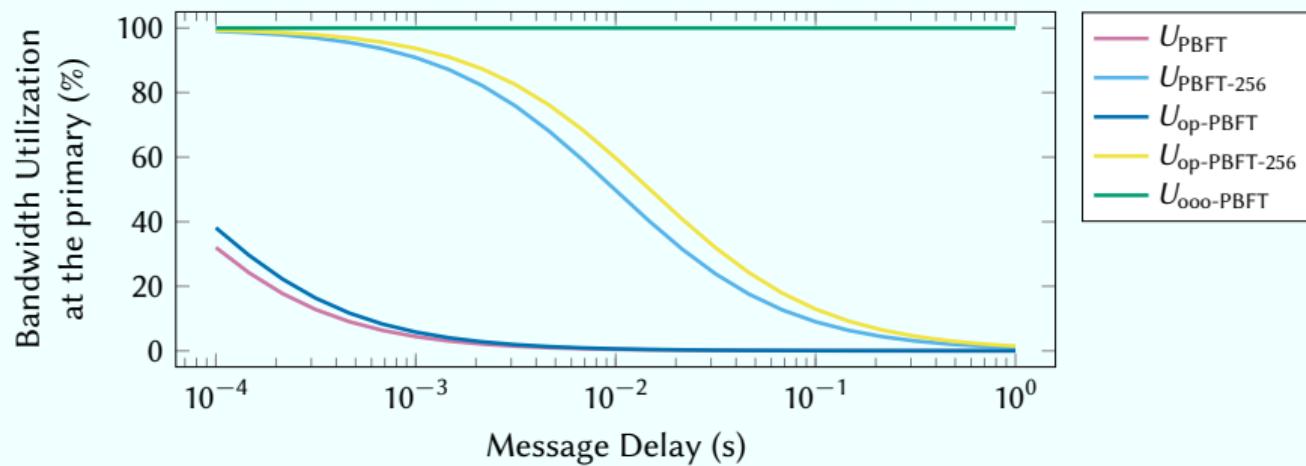
# Implementation techniques for PBFT: Summary

Batching introduces very high round latencies.

Out-of-order processing has high implementation complexity.

Overlapping only provides limited gains.

Assumption:  $n = 4$ ,  $B = 100 \text{ MiB/s}$ ,  $\delta = 15 \text{ ms}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$



# Primary-backup Consensus Beyond PBFT

A PBFT-like design is at the basis of *many* consensus protocols.

# Primary-backup Consensus Beyond PBFT

A PBFT-like design is at the basis of *many* consensus protocols.

## Technologies employed by PBFT-like consensus

**Threshold signatures** eliminate quadratic all-to-all communication.

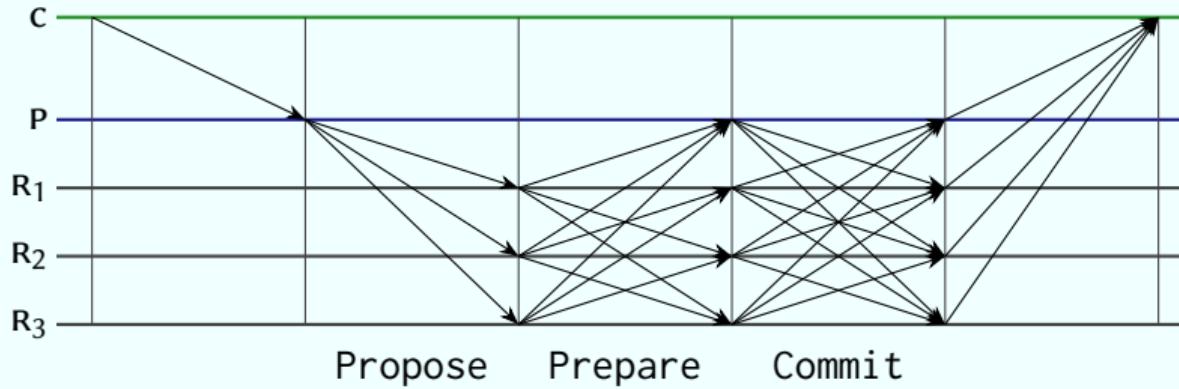
**Speculative execution** execute before strong recovery guarantees are met.

**Optimistic execution** fully optimize for when the primary is correct.

**Trusted components** use hardware components that cannot behave Byzantine.

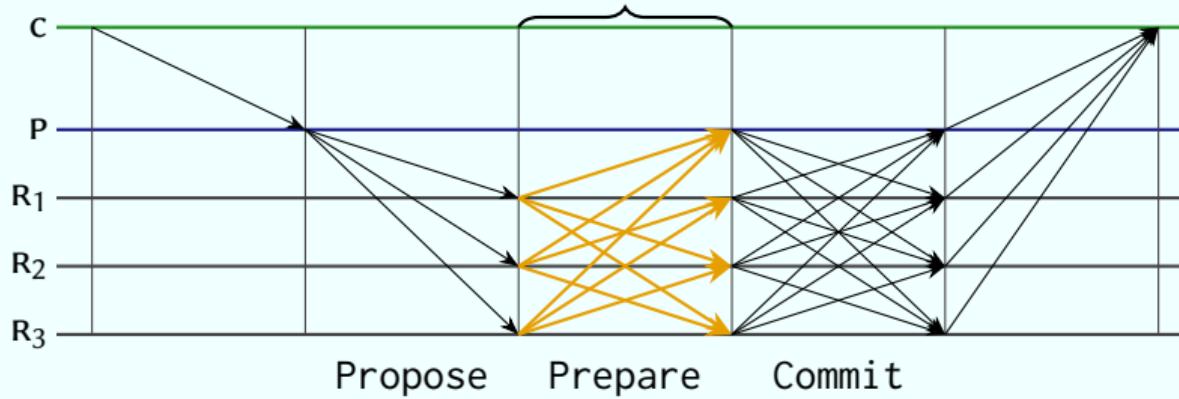
Here, we will only cover *threshold signatures*.

# All-to-All Communication in PBFT

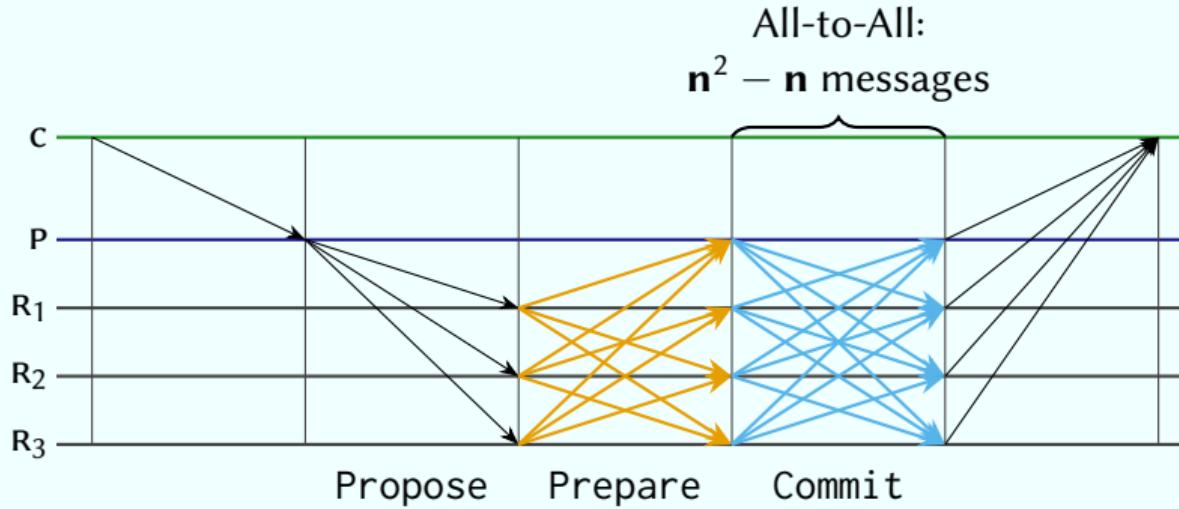


# All-to-All Communication in PBFT

Almost All-to-All:  
 $(n - 1)^2$  messages

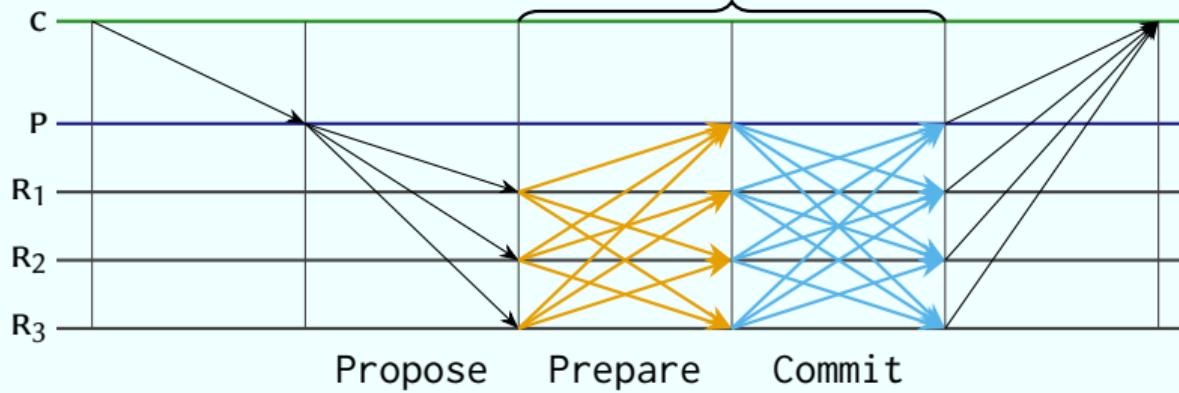


# All-to-All Communication in PBFT

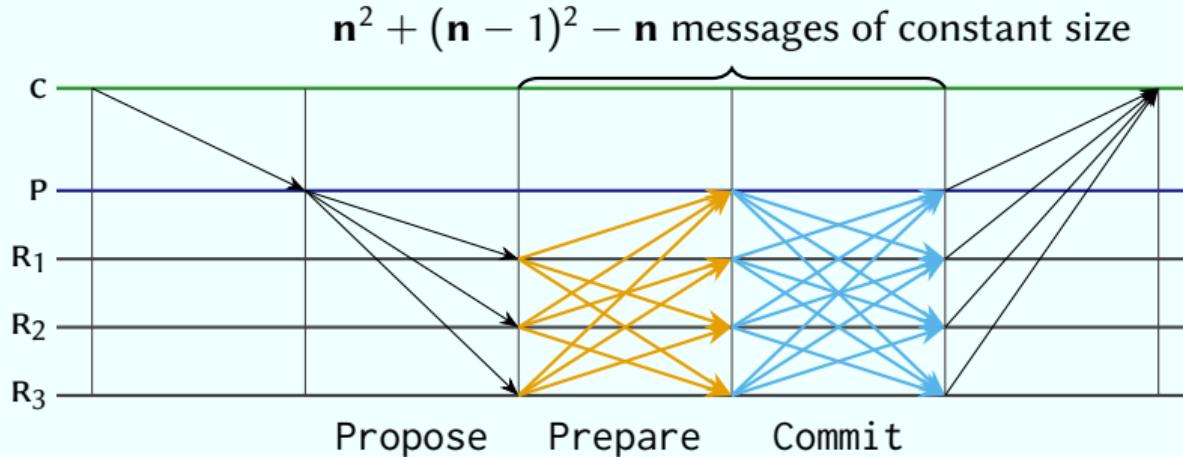


# All-to-All Communication in PBFT

$$n^2 + (n - 1)^2 - n \text{ messages of constant size}$$



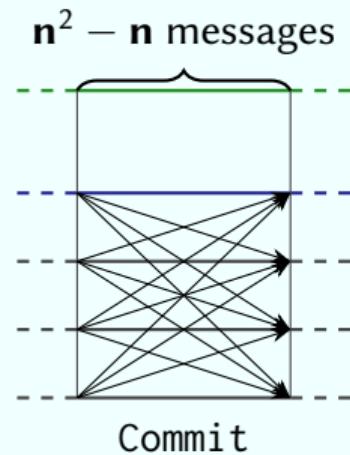
## All-to-All Communication in PBFT



Challenge: Reduce communication from  $\mathcal{O}(n^2)$  to  $\mathcal{O}(n)$  messages of constant size.

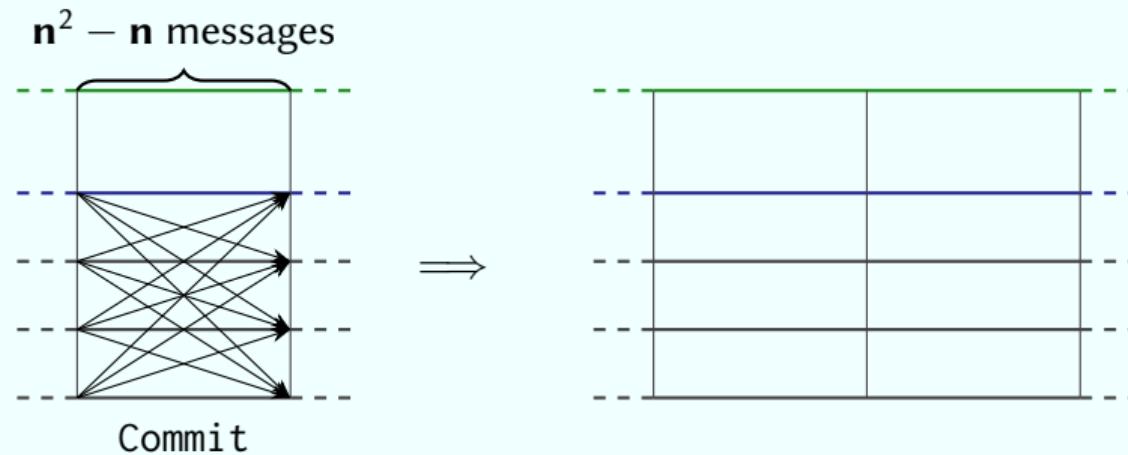
# Tackling All-to-All via All-to-one-to-All Aggregation

Consider the commit phase



# Tackling All-to-All via All-to-one-to-All Aggregation

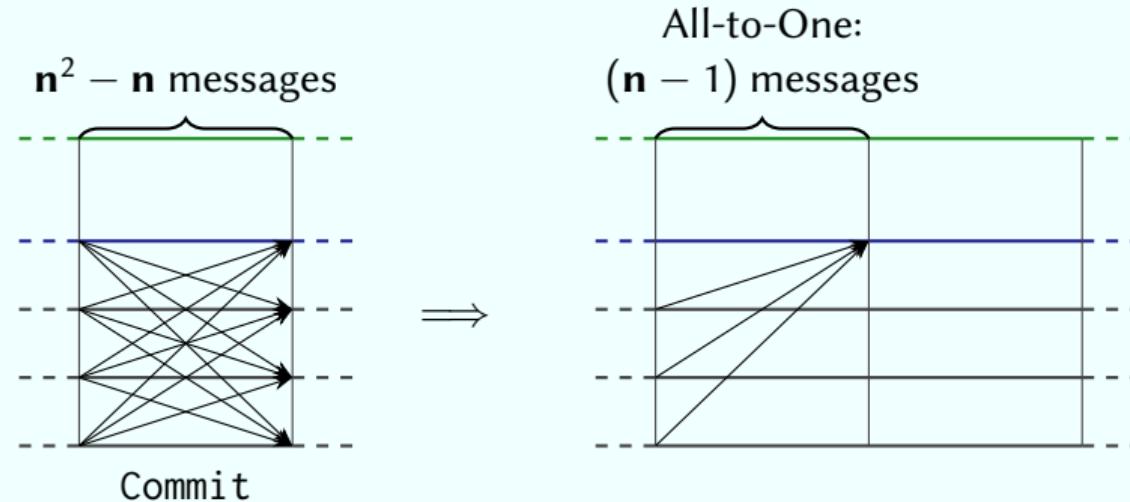
Consider the commit phase



Idea: All replicas send to *one aggregator* that then sends to all replicas.

# Tackling All-to-All via All-to-one-to-All Aggregation

Consider the commit phase

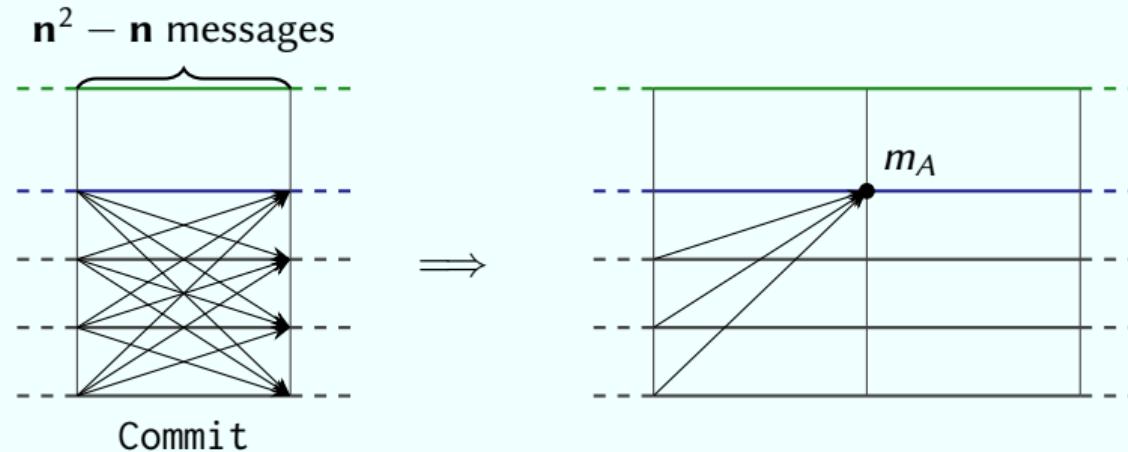


Idea: All replicas send to *one aggregator* that then sends to all replicas.

1. All replicas send their Commit messages to the aggregator.

# Tackling All-to-All via All-to-one-to-All Aggregation

Consider the commit phase

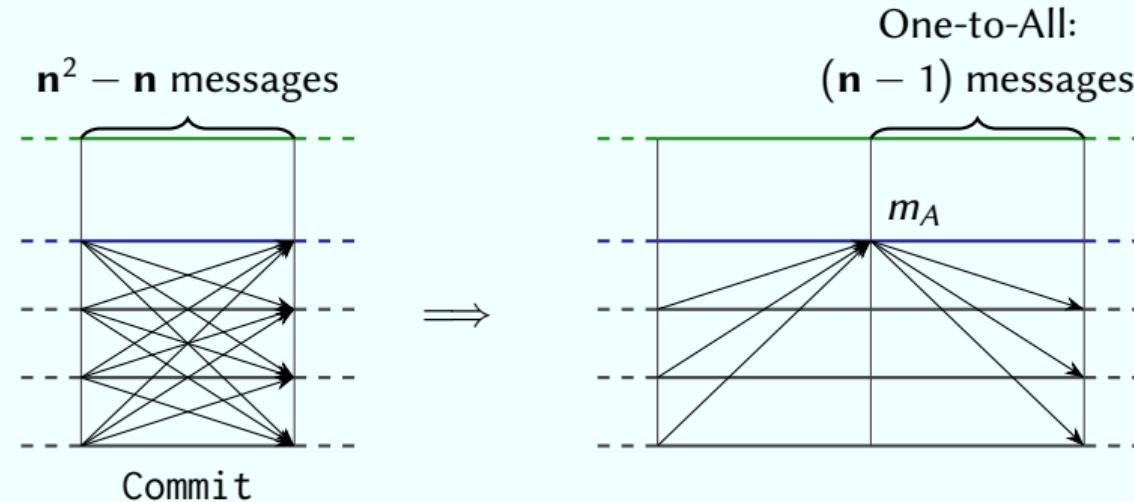


Idea: All replicas send to *one aggregator* that then sends to all replicas.

2. The aggregator combines  $n - f$  Commit messages into an aggregated message  $m_A$ .

# Tackling All-to-All via All-to-one-to-All Aggregation

Consider the commit phase

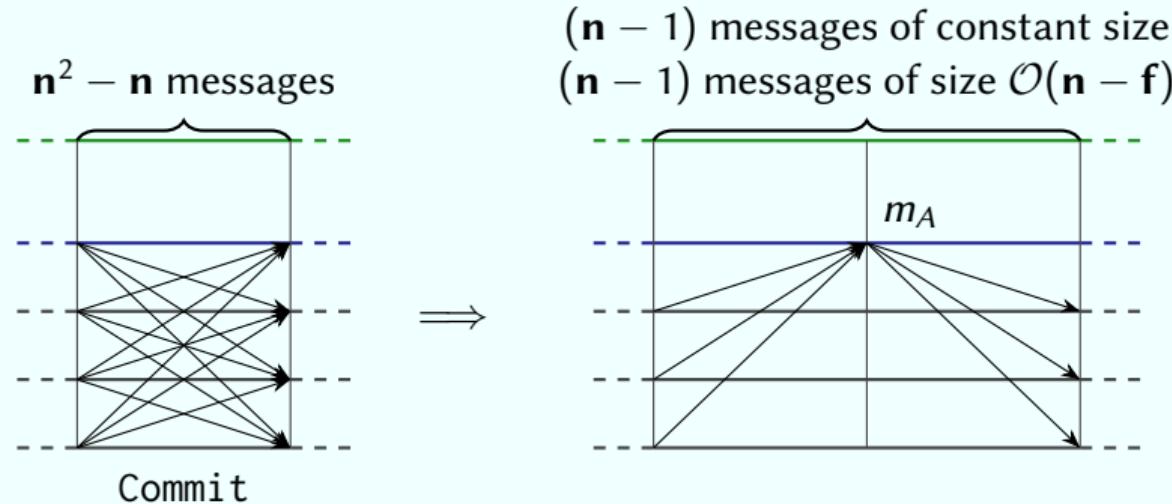


Idea: All replicas send to *one aggregator* that then sends to all replicas.

3. The aggregator sends  $m_A$  to all replicas.

# Tackling All-to-All via All-to-one-to-All Aggregation

Consider the commit phase

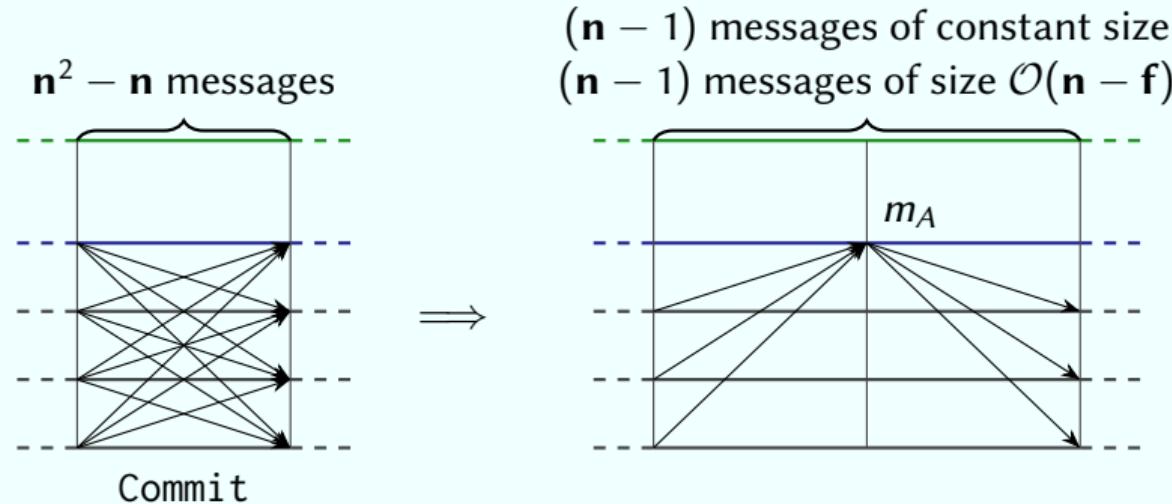


Idea: All replicas send to *one aggregator* that then sends to all replicas.

3. The aggregator sends  $m_A$  to all replicas *of size  $\mathcal{O}(n-f)$  each*.

# Tackling All-to-All via All-to-one-to-All Aggregation

Consider the commit phase



Idea: All replicas send to *one aggregator* that then sends to all replicas.

Effectively reduced communication from  $\mathcal{O}(n^2)$  to  $\mathcal{O}(n(n-f))$ .

## Improving Aggregation with Threshold Signatures

Problem: An aggregated message of size  $c$  will have size  $\mathcal{O}(c(n - f))$ .

- ▶ We have identical Commit messages from at-least  $n - f$  replicas.
- ▶ Goal: aggregate these into a single message of size  $\mathcal{O}(c)$  instead of  $\mathcal{O}(c(n - f))$ .

# Improving Aggregation with Threshold Signatures

Problem: An aggregated message of size  $c$  will have size  $\mathcal{O}(c(n - f))$ .

- ▶ We have identical Commit messages from at-least  $n - f$  replicas.
- ▶ Goal: aggregate these into a single message of size  $\mathcal{O}(c)$  instead of  $\mathcal{O}(c(n - f))$ .
- ▶ Crucially: we want to aggregate the digital signatures!

# Improving Aggregation with Threshold Signatures

Problem: An aggregated message of size  $c$  will have size  $\mathcal{O}(c(n - f))$ .

- ▶ We have identical Commit messages from at-least  $n - f$  replicas.
- ▶ Goal: aggregate these into a single message of size  $\mathcal{O}(c)$  instead of  $\mathcal{O}(c(n - f))$ .
- ▶ Crucially: we want to aggregate the digital signatures!

Solution: Using a  $n : f$ -threshold-signature scheme with public key  $K$

- ▶ Each replica has a unique private key.
- ▶ Replicas can produce partial signatures for value  $v$  using their private key.
- ▶ Using  $n - f$  partial signatures for  $v$ , one can produce a *threshold signature*.

# Improving Aggregation with Threshold Signatures

Problem: An aggregated message of size  $c$  will have size  $\mathcal{O}(c(n - f))$ .

- ▶ We have identical Commit messages from at-least  $n - f$  replicas.
- ▶ Goal: aggregate these into a single message of size  $\mathcal{O}(c)$  instead of  $\mathcal{O}(c(n - f))$ .
- ▶ Crucially: we want to aggregate the digital signatures!

Solution: Using a  $n : f$ -threshold-signature scheme with public key  $K$

- ▶ Each replica has a unique private key.
- ▶ Replicas can produce partial signatures for value  $v$  using their private key.
- ▶ Using  $n - f$  partial signatures for  $v$ , one can produce a *threshold signature*.

# Improving Aggregation with Threshold Signatures

Problem: An aggregated message of size  $c$  will have size  $\mathcal{O}(c(n - f))$ .

- ▶ We have identical Commit messages from at-least  $n - f$  replicas.
- ▶ Goal: aggregate these into a single message of size  $\mathcal{O}(c)$  instead of  $\mathcal{O}(c(n - f))$ .
- ▶ Crucially: we want to aggregate the digital signatures!

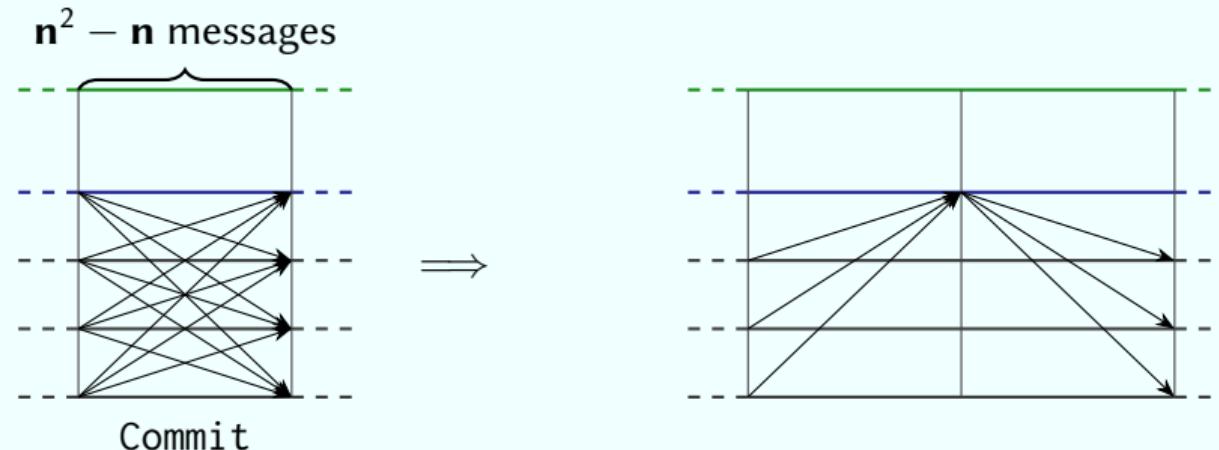
Solution: Using a  $n : f$ -threshold-signature scheme with public key  $K$

- ▶ Each replica has a unique private key.
- ▶ Replicas can produce partial signatures for value  $v$  using their private key.
- ▶ Using  $n - f$  partial signatures for  $v$ , one can produce a *threshold signature*.

Threshold signatures aggregate  $n - f$  distinct signatures into a *single constant-sized* value.

# All-to-one-to-All Aggregation with Threshold Signatures

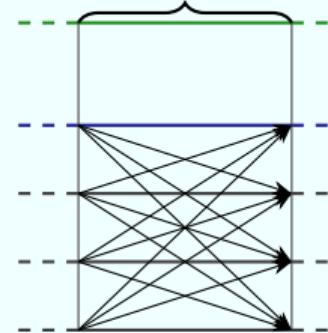
Consider the commit phase



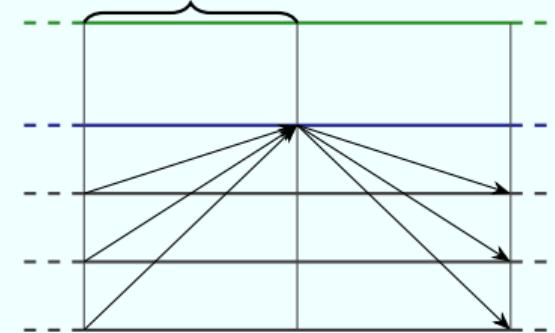
# All-to-one-to-All Aggregation with Threshold Signatures

Consider the commit phase

$n^2 - n$  messages



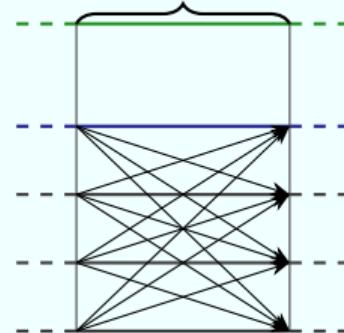
All-to-One:  
 $(n - 1)$  partial signatures



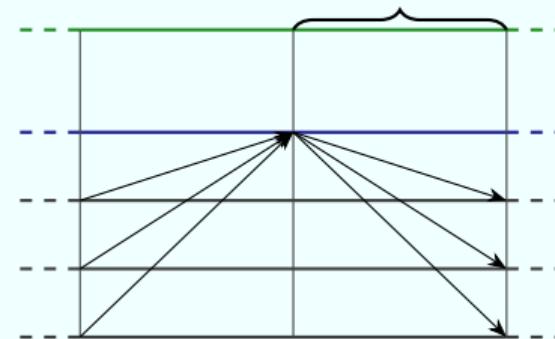
# All-to-one-to-All Aggregation with Threshold Signatures

Consider the commit phase

$n^2 - n$  messages



One-to-All:  
 $(n - 1)$  threshold signatures

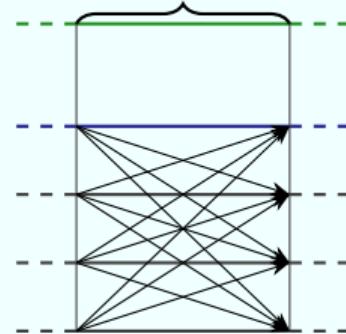


Commit

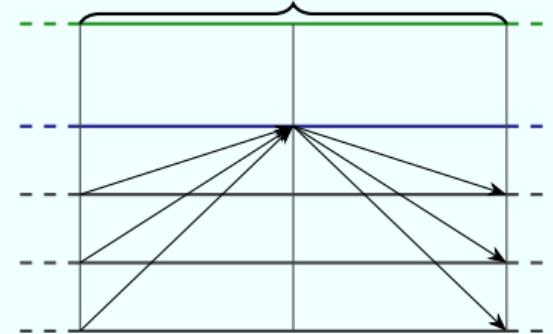
# All-to-one-to-All Aggregation with Threshold Signatures

Consider the commit phase

$n^2 - n$  messages

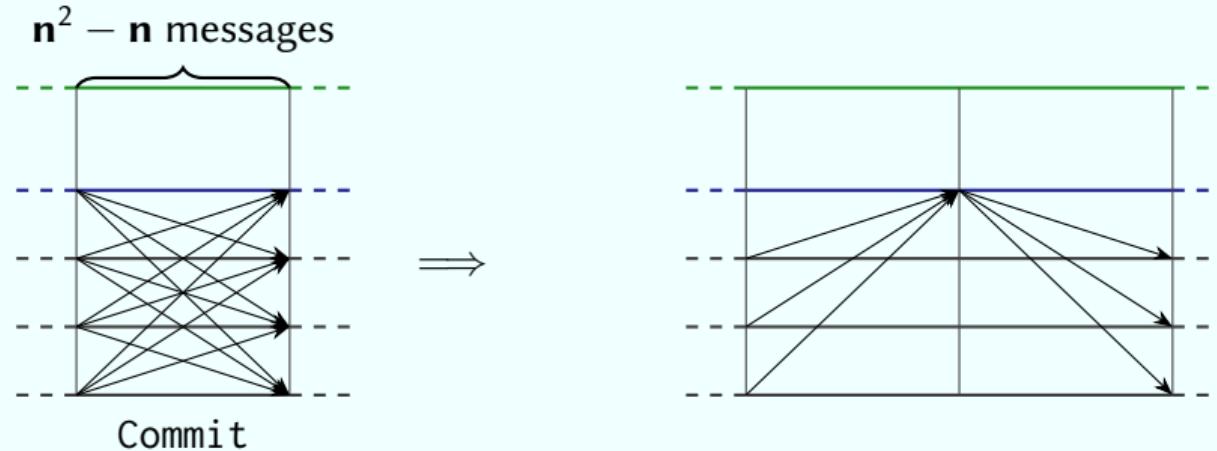


$(n - 1)$  partial signatures of constant size  
 $(n - 1)$  threshold signatures of constant size



# All-to-one-to-All Aggregation with Threshold Signatures

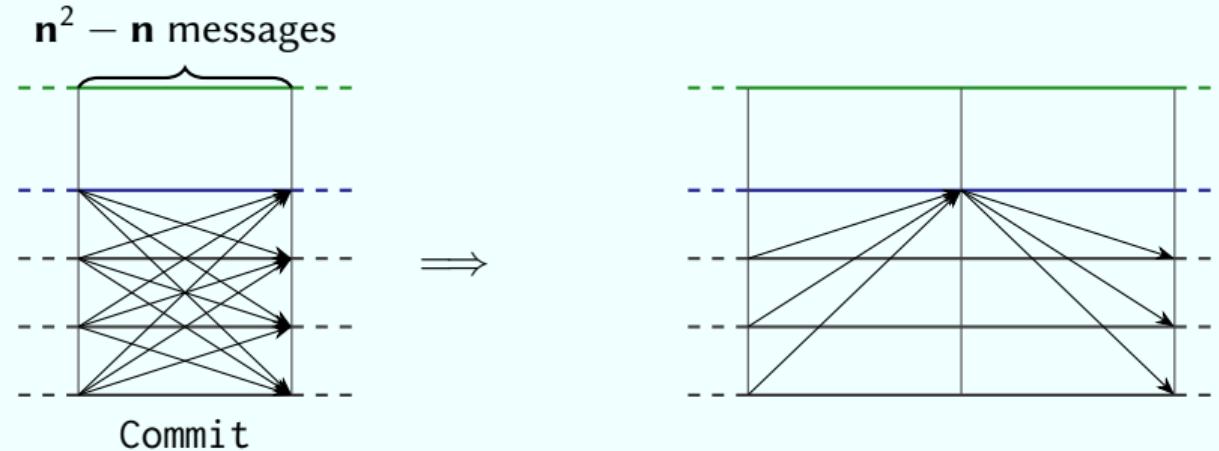
Consider the commit phase



Effectively reduced communication from  $\mathcal{O}(n^2)$  to  $\mathcal{O}(n)$ .

# All-to-one-to-All Aggregation with Threshold Signatures

Consider the commit phase



Effectively reduced communication from  $\mathcal{O}(n^2)$  to  $\mathcal{O}(n)$ .  
Similar change can be made to the prepare phase.

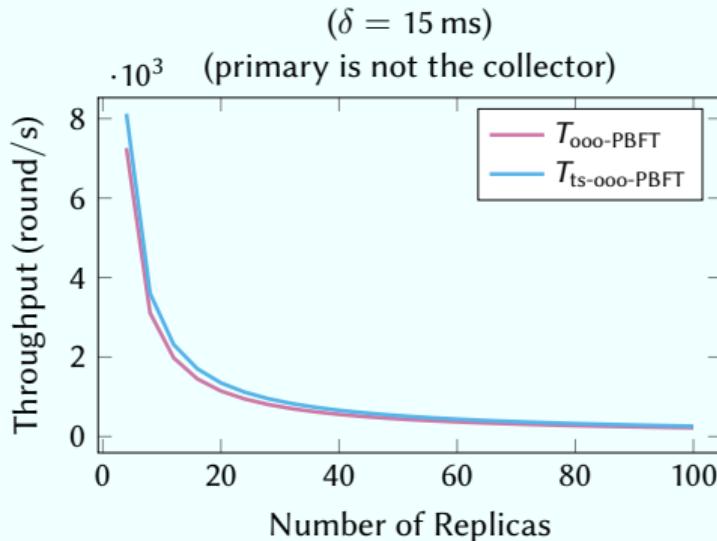
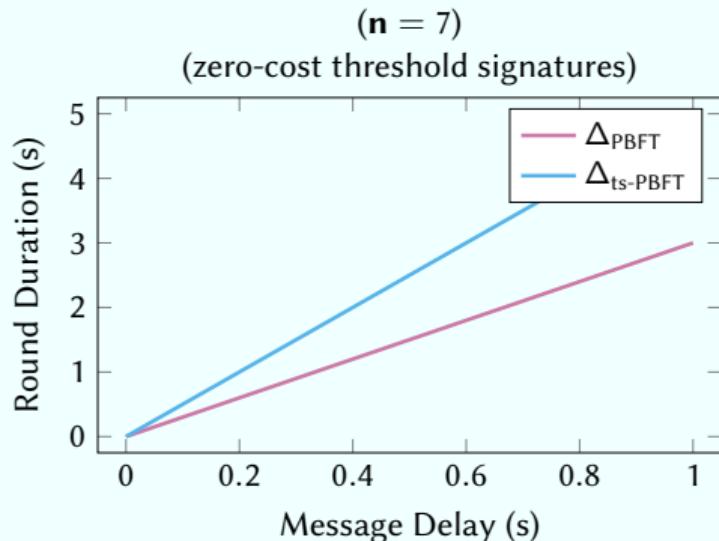
## Using Threshold Signatures in PBFT

- ▶ Both prepare and commit phase: from  $2(n - 1)^2$  to  $4(n - 1)$  messages.
- ▶ Consensus from *three* to *five* rounds: higher consensus and client latencies.
- ▶ High *computational cost* for the aggregator.
- ▶ Need recovery methods to deal with *faulty aggregators*.

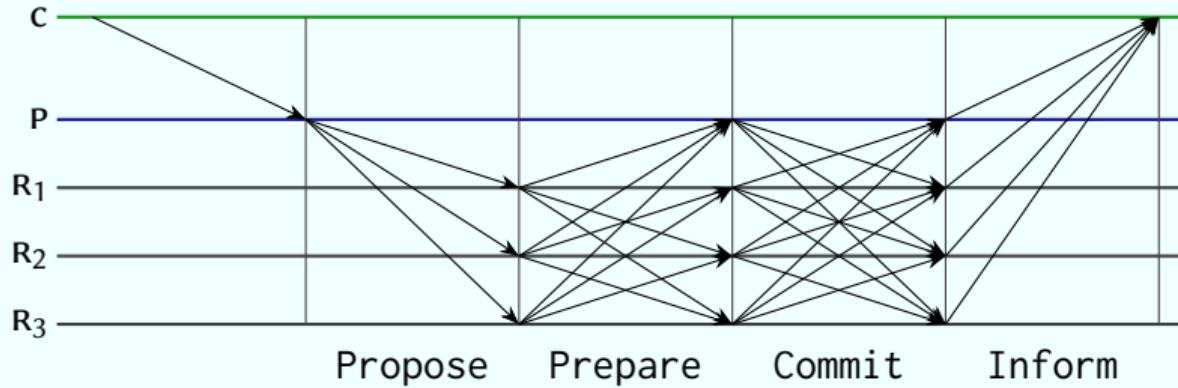
# Using Threshold Signatures in PBFT

- ▶ Both prepare and commit phase: from  $2(n - 1)^2$  to  $4(n - 1)$  messages.
- ▶ Consensus from *three* to *five* rounds: higher consensus and client latencies.
- ▶ High *computational cost* for the aggregator.
- ▶ Need recovery methods to deal with *faulty aggregators*.

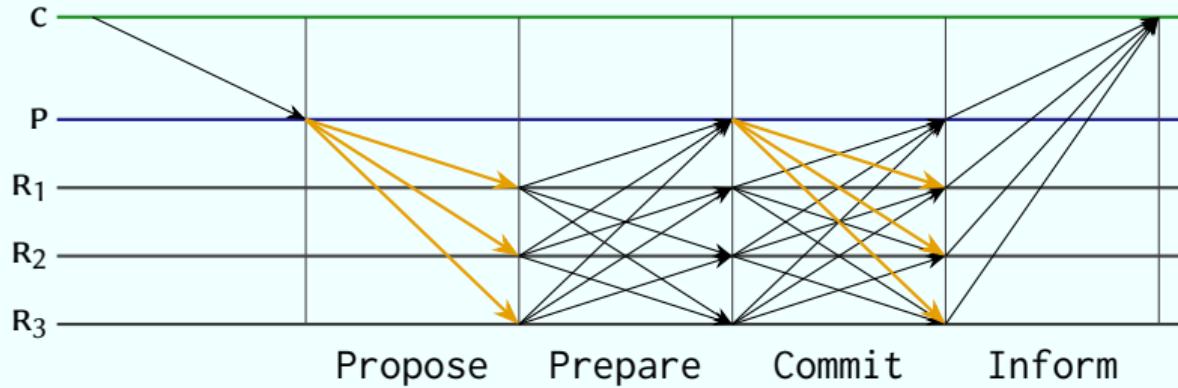
Assumption:  $B = 100 \text{ MiB/s}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$



## Limitations of Primary-Backup Consensus

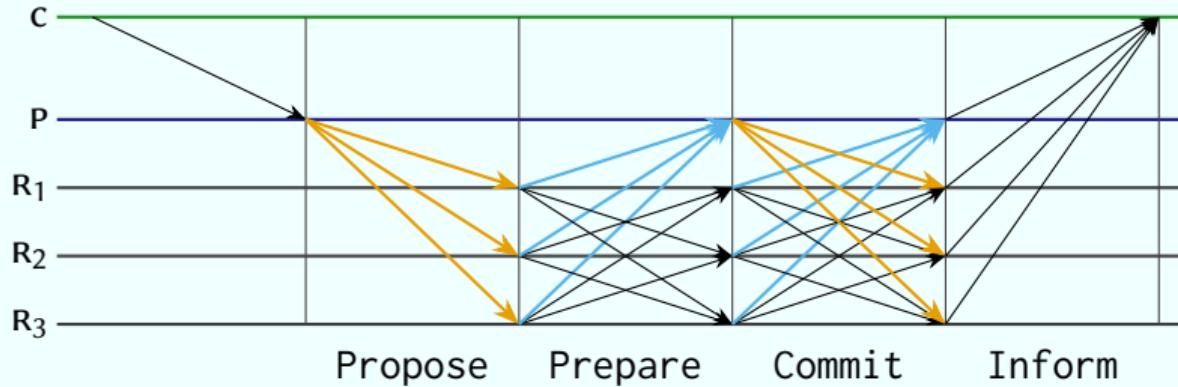


## Limitations of Primary-Backup Consensus



Primary Send  $(n - 1)$  Propose, send  $(n - 1)$  Commit.

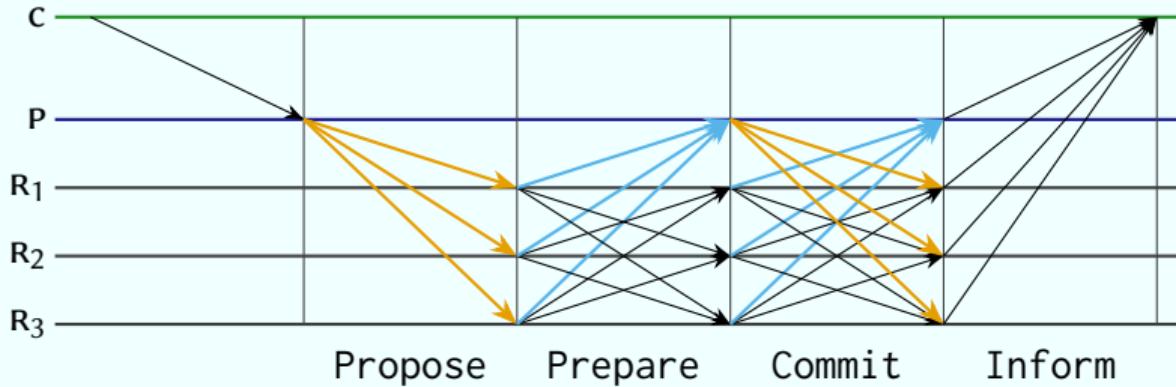
## Limitations of Primary-Backup Consensus



Primary Send  $(n - 1)$  Propose, send  $(n - 1)$  Commit.

Receive  $(n - 1)$  Prepare, receive  $(n - 1)$  Commit.

## Limitations of Primary-Backup Consensus

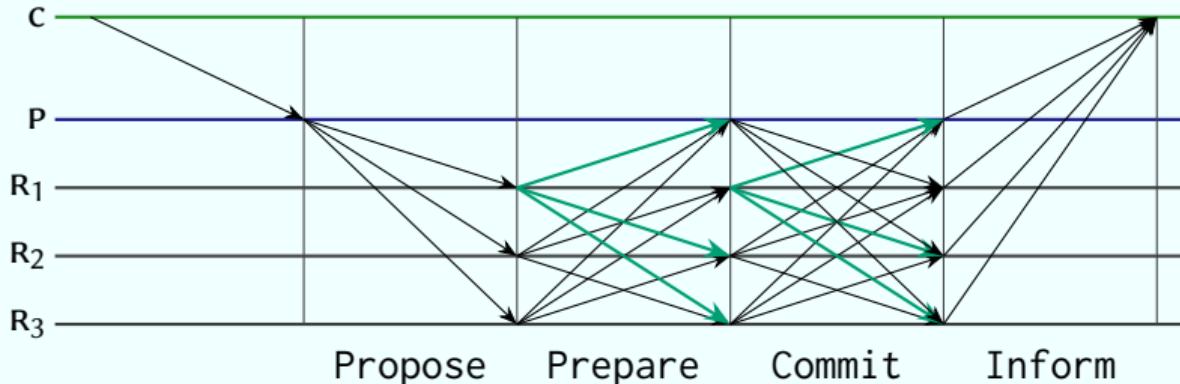


Primary Send  $(n - 1)$  Propose, send  $(n - 1)$  Commit.

Receive  $(n - 1)$  Prepare, receive  $(n - 1)$  Commit.

Total:  $m(n - 1)s_t + 3(n - 1)s_m$ .

## Limitations of Primary-Backup Consensus



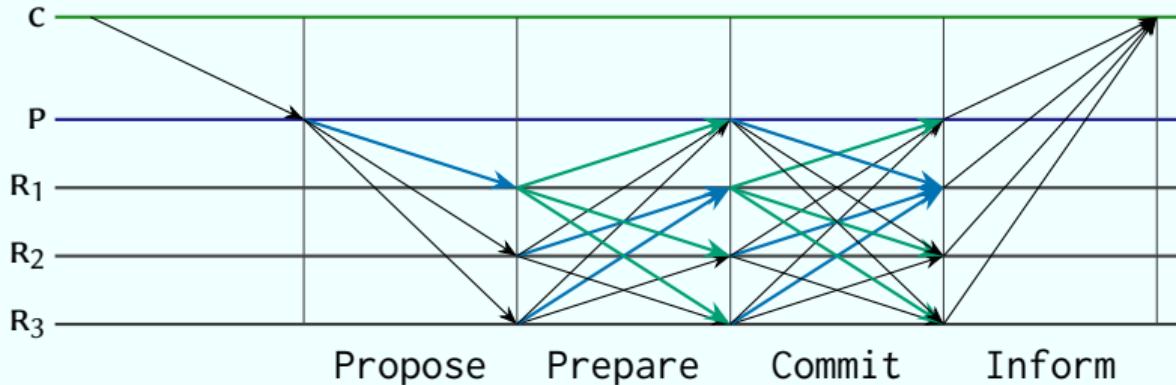
Primary Send  $(n - 1)$  Propose, send  $(n - 1)$  Commit.

Receive  $(n - 1)$  Prepare, receive  $(n - 1)$  Commit.

Total:  $m(n - 1)s_t + 3(n - 1)s_m$ .

Backup Send  $(n - 1)$  Prepare, send  $(n - 1)$  Commit.

# Limitations of Primary-Backup Consensus



**Primary** Send  $(n - 1)$  Propose, send  $(n - 1)$  Commit.

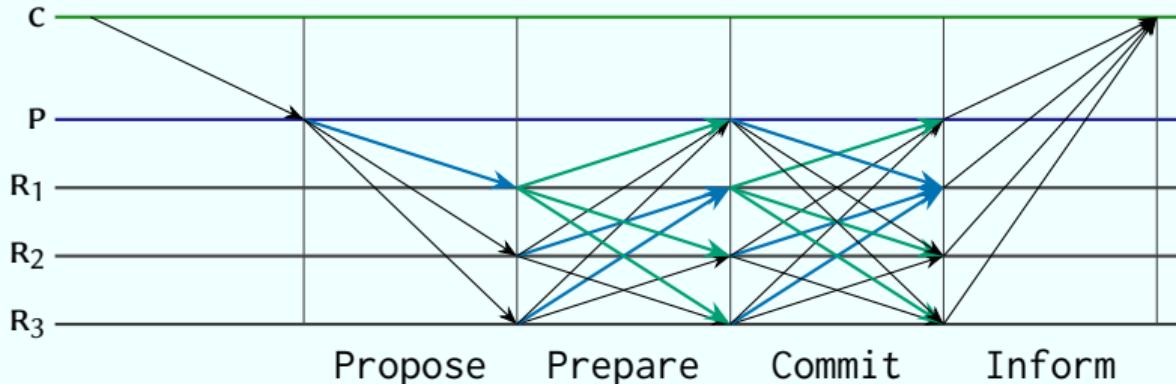
Receive  $(n - 1)$  Prepare, receive  $(n - 1)$  Commit.

Total:  $m(n - 1)s_t + 3(n - 1)s_m$ .

**Backup** Send  $(n - 1)$  Prepare, send  $(n - 1)$  Commit.

Receive one Propose, receive  $(n - 2)$  Prepare, receive  $(n - 1)$  Commit.

# Limitations of Primary-Backup Consensus



**Primary** Send  $(n - 1)$  Propose, send  $(n - 1)$  Commit.

Receive  $(n - 1)$  Prepare, receive  $(n - 1)$  Commit.

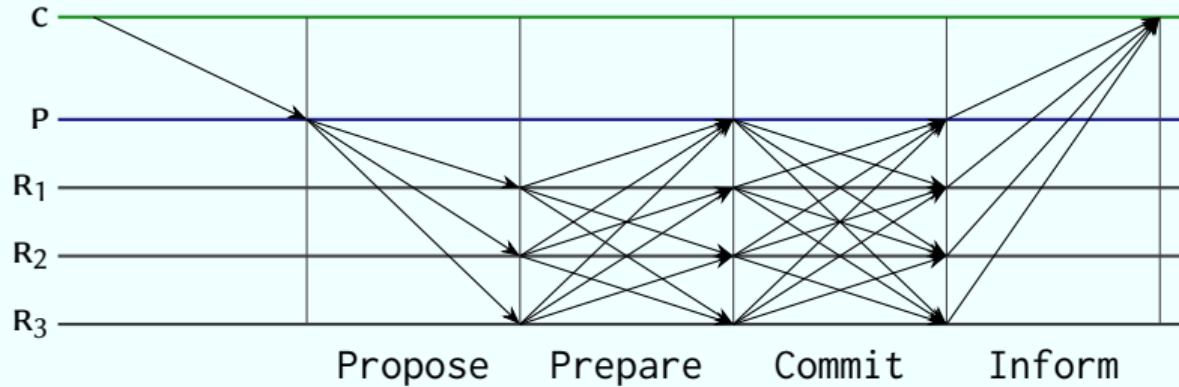
Total:  $m(n - 1)s_t + 3(n - 1)s_m$ .

**Backup** Send  $(n - 1)$  Prepare, send  $(n - 1)$  Commit.

Receive one Propose, receive  $(n - 2)$  Prepare, receive  $(n - 1)$  Commit.

Total:  $ms_t + 4(n - 1)s_m - s_m$ .

## Limitations of Primary-Backup Consensus



Bandwidth ratio between primary and backups

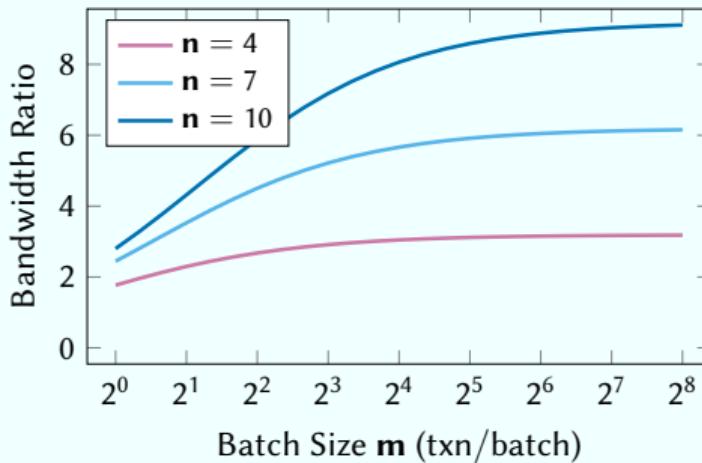
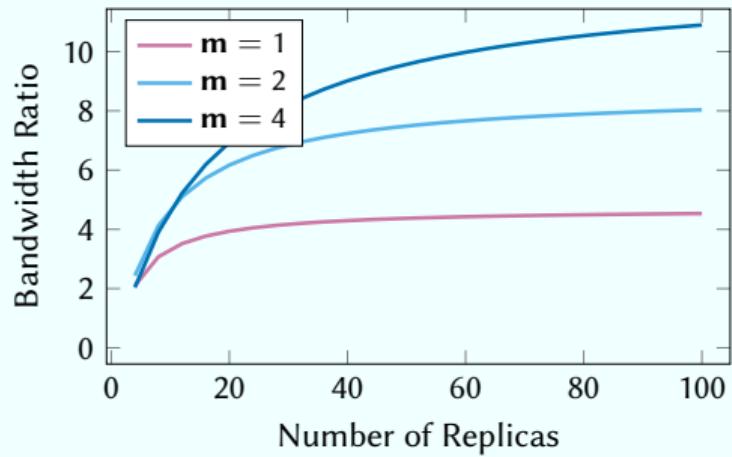
$$R_{\text{PBFT-}m} = \frac{m(n-1)s_t + 3(n-1)s_m}{ms_t + 4(n-1)s_m - s_m}.$$

# Limitations of Primary-Backup Consensus

Bandwidth ratio between primary and backups

$$R_{\text{PBFT-}m} = \frac{m(n-1)s_t + 3(n-1)s_m}{ms_t + 4(n-1)s_m - s_m}.$$

Assumption:  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$



## Maximum Throughput of Primary-Backup Consensus

Consider failure-free replication: Primary *only* proposes, no other communication.

## Maximum Throughput of Primary-Backup Consensus

Consider failure-free replication: Primary *only* proposes, no other communication.

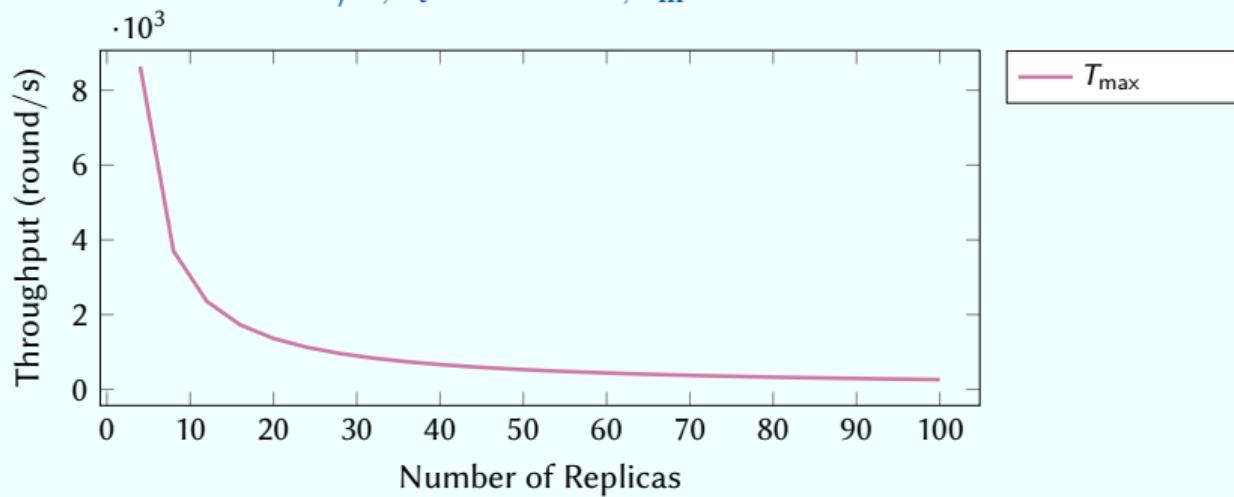
$$T_{\max} = \frac{B}{(\mathbf{n} - 1)s_t}.$$

# Maximum Throughput of Primary-Backup Consensus

Consider failure-free replication: Primary *only* proposes, no other communication.

$$T_{\max} = \frac{B}{(\mathbf{n} - 1)s_t}.$$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$

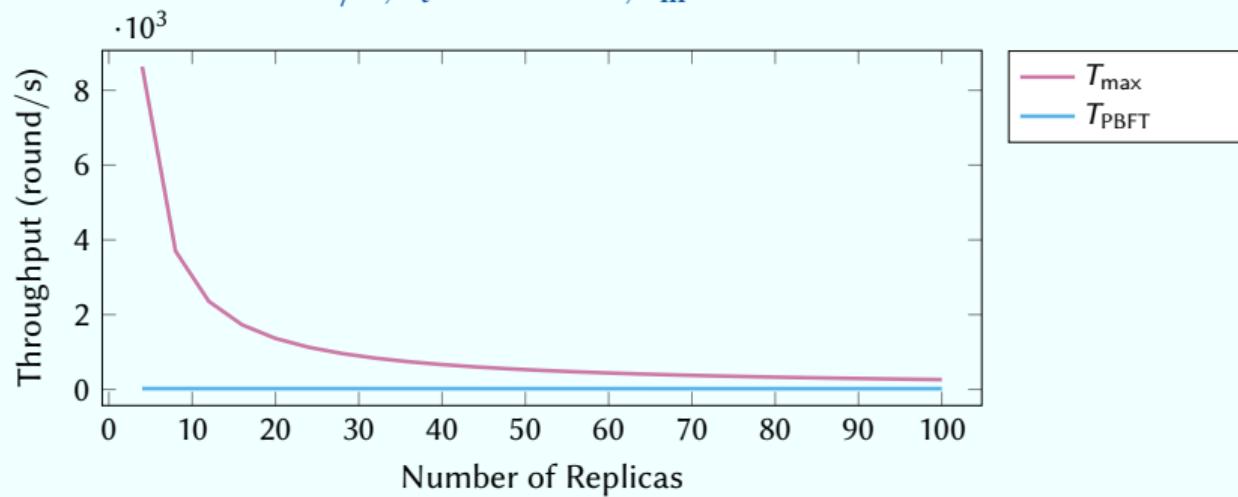


# Maximum Throughput of Primary-Backup Consensus

Consider failure-free replication: Primary *only* proposes, no other communication.

$$T_{\max} = \frac{B}{(\mathbf{n} - 1)s_t}.$$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$

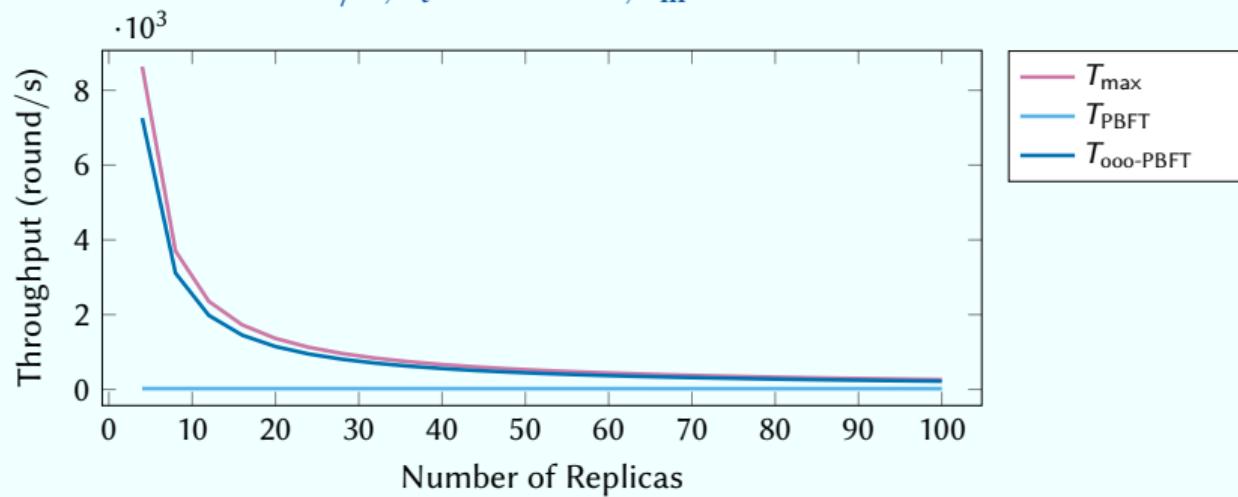


# Maximum Throughput of Primary-Backup Consensus

Consider failure-free replication: Primary *only* proposes, no other communication.

$$T_{\max} = \frac{B}{(\mathbf{n} - 1)s_t}.$$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$

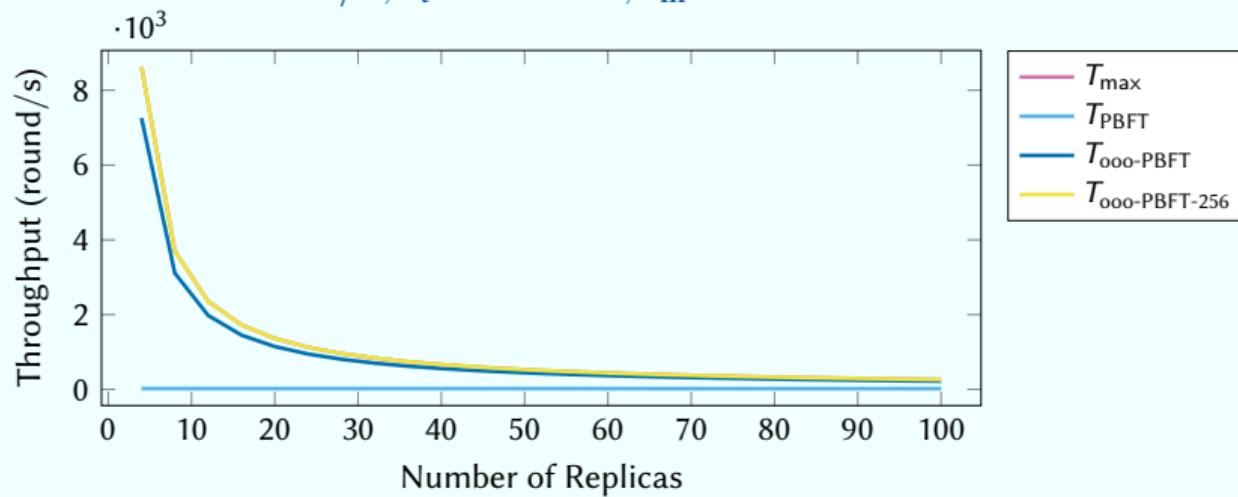


# Maximum Throughput of Primary-Backup Consensus

Consider failure-free replication: Primary *only* proposes, no other communication.

$$T_{\max} = \frac{B}{(\mathbf{n} - 1)s_t}.$$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$



## Concurrent Consensus

Idea: Multiple instances of PBFT, each with a distinct primary

$1 \leq z \leq n$  primaries:  $z$  simultaneous rounds of consensus that decide the next  $z$  requests.

## Concurrent Consensus

Idea: Multiple instances of PBFT, each with a distinct primary

$1 \leq z \leq n$  primaries:  $z$  simultaneous rounds of consensus that decide the next  $z$  requests.

Consider the communication of one of the  $z$  primaries.

## Concurrent Consensus

Idea: Multiple instances of PBFT, each with a distinct primary

$1 \leq z \leq n$  primaries:  $z$  simultaneous rounds of consensus that decide the next  $z$  requests.

Consider the communication of one of the  $z$  primaries.

- As *primary* of its own instance:

Send  $(n - 1)$  Propose, send  $(n - 1)$  Commit.

Receive  $(n - 1)$  Prepare, receive  $(n - 1)$  Commit.

Total:  $m(n - 1)s_t + 3(n - 1)s_m$ .

# Concurrent Consensus

Idea: Multiple instances of PBFT, each with a distinct primary

$1 \leq z \leq n$  primaries:  $z$  simultaneous rounds of consensus that decide the next  $z$  requests.

Consider the communication of one of the  $z$  primaries.

- As *primary* of its own instance:

Send  $(n - 1)$  Propose, send  $(n - 1)$  Commit.

Receive  $(n - 1)$  Prepare, receive  $(n - 1)$  Commit.

Total:  $m(n - 1)s_t + 3(n - 1)s_m$ .

- As *backup* of the other  $z - 1$  instances ( $(z - 1)$  times):

Send  $(n - 1)$  Prepare, send  $(n - 1)$  Commit.

Receive one Propose, receive  $(n - 2)$  Prepare, receive  $(n - 1)$  Commit.

Total:  $(z - 1)(ms_t + 4(n - 1)s_m - s_m)$ .

# Concurrent Consensus

Idea: Multiple instances of PBFT, each with a distinct primary

$1 \leq z \leq n$  primaries:  $z$  simultaneous rounds of consensus that decide the next  $z$  requests.

Consider the communication of one of the  $z$  primaries.

- As *primary* of its own instance:

Send  $(n - 1)$  Propose, send  $(n - 1)$  Commit.

Receive  $(n - 1)$  Prepare, receive  $(n - 1)$  Commit.

Total:  $m(n - 1)s_t + 3(n - 1)s_m$ .

- As *backup* of the other  $z - 1$  instances ( $(z - 1)$  times):

Send  $(n - 1)$  Prepare, send  $(n - 1)$  Commit.

Receive one Propose, receive  $(n - 2)$  Prepare, receive  $(n - 1)$  Commit.

Total:  $(z - 1)(ms_t + 4(n - 1)s_m - s_m)$ .

$zmB$

$$T_{c\text{-ooo-PBFT-}(z, m)} = \frac{zmB}{(m(n - 1)s_t + 3(n - 1)s_m) + ((z - 1)(ms_t + 4(n - 1)s_m - s_m))}.$$

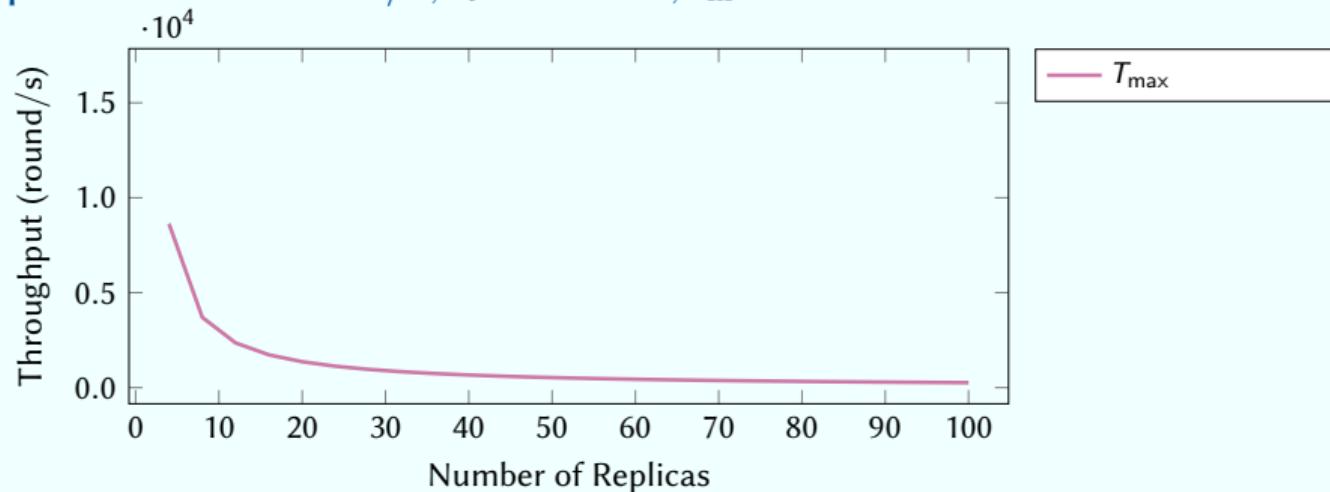
# Concurrent Consensus

Idea: Multiple instances of PBFT, each with a distinct primary

$1 \leq z \leq n$  primaries:  $z$  simultaneous rounds of consensus that decide the next  $z$  requests.

$$T_{c\text{-ooo-PBFT-}(z,m)} = \frac{zmB}{(m(n-1)s_t + 3(n-1)s_m) + ((z-1)(ms_t + 4(n-1)s_m - s_m))}.$$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$



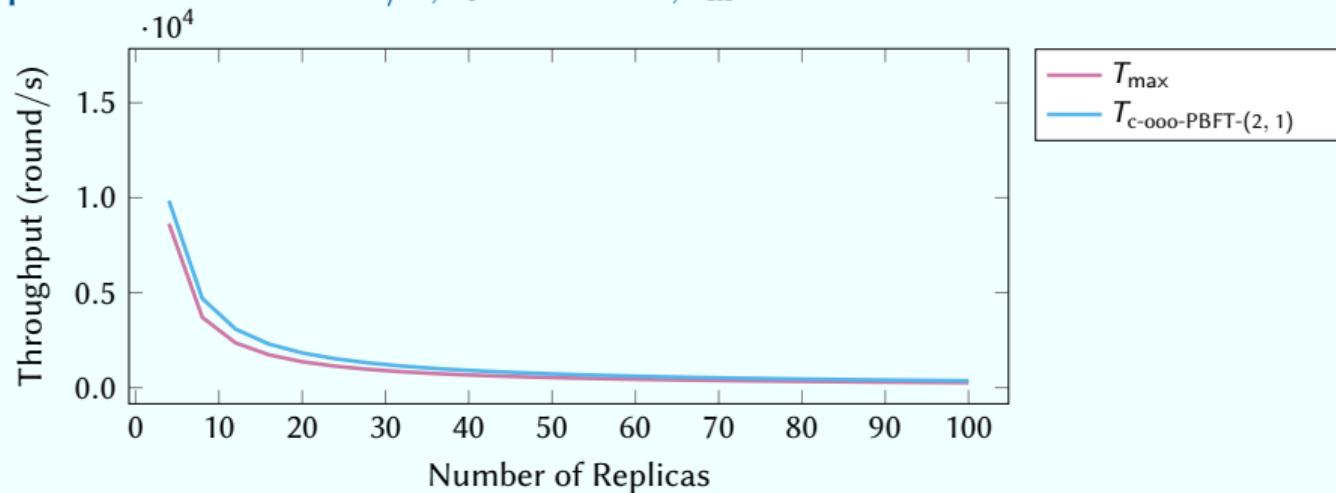
# Concurrent Consensus

Idea: Multiple instances of PBFT, each with a distinct primary

$1 \leq z \leq n$  primaries:  $z$  simultaneous rounds of consensus that decide the next  $z$  requests.

$$T_{c\text{-ooo-PBFT-}(z,m)} = \frac{zmB}{(m(n-1)s_t + 3(n-1)s_m) + ((z-1)(ms_t + 4(n-1)s_m - s_m))}.$$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$



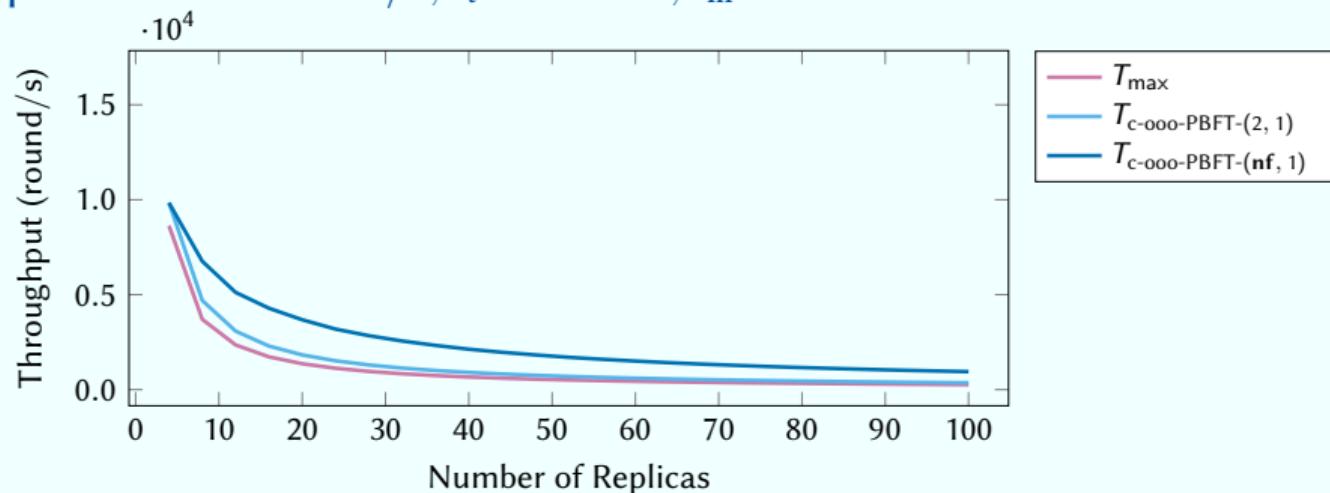
# Concurrent Consensus

Idea: Multiple instances of PBFT, each with a distinct primary

$1 \leq z \leq n$  primaries:  $z$  simultaneous rounds of consensus that decide the next  $z$  requests.

$$T_{c\text{-ooo-PBFT-}(z,m)} = \frac{zmB}{(m(n-1)s_t + 3(n-1)s_m) + ((z-1)(ms_t + 4(n-1)s_m - s_m))}.$$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$



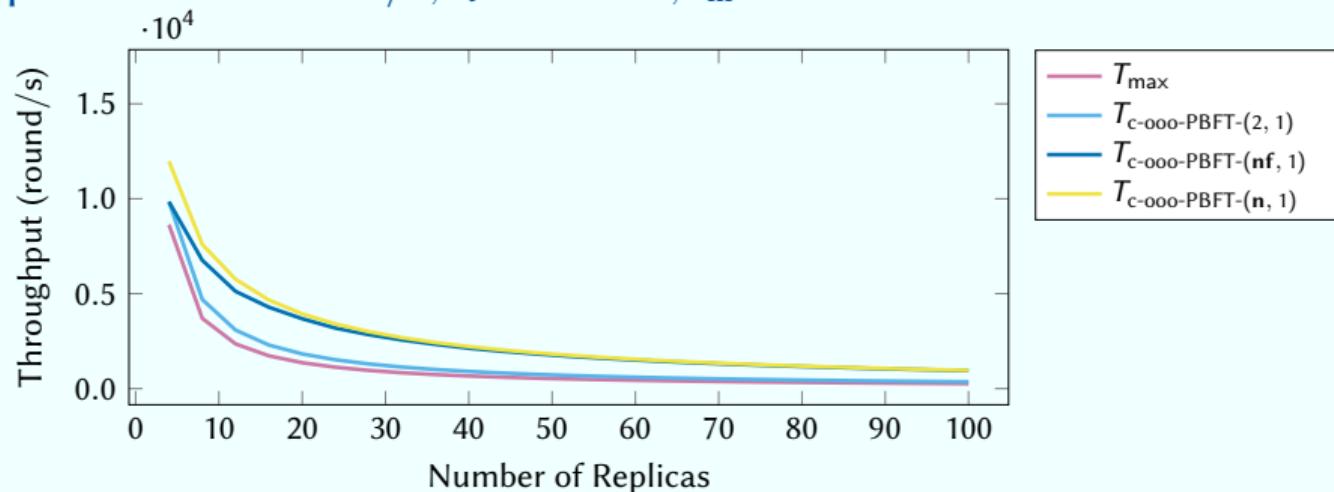
# Concurrent Consensus

Idea: Multiple instances of PBFT, each with a distinct primary

$1 \leq z \leq n$  primaries:  $z$  simultaneous rounds of consensus that decide the next  $z$  requests.

$$T_{c\text{-ooo-PBFT-}(z, m)} = \frac{zmB}{(m(n-1)s_t + 3(n-1)s_m) + ((z-1)(ms_t + 4(n-1)s_m - s_m))}.$$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$



# Concurrent Consensus

Idea: Multiple instances of PBFT, each with a distinct primary

$1 \leq z \leq n$  primaries:  $z$  simultaneous rounds of consensus that decide the next  $z$  requests.

$$T_{c\text{-ooo-PBFT-}(z, m)} = \frac{zmB}{(m(n-1)s_t + 3(n-1)s_m) + ((z-1)(ms_t + 4(n-1)s_m - s_m))}.$$

Assumption:  $B = 100 \text{ MiB/s}$ ,  $s_t = 4048 \text{ B}$ ,  $s_m = 256 \text{ B}$

