

An In-Depth Look of BFT Consensus in Blockchain: Challenges and Opportunities



Suyash Gupta



Jelle Hellings



Sajjad Rahnema



Mohammad Sadoghi

MokaBlox LLC

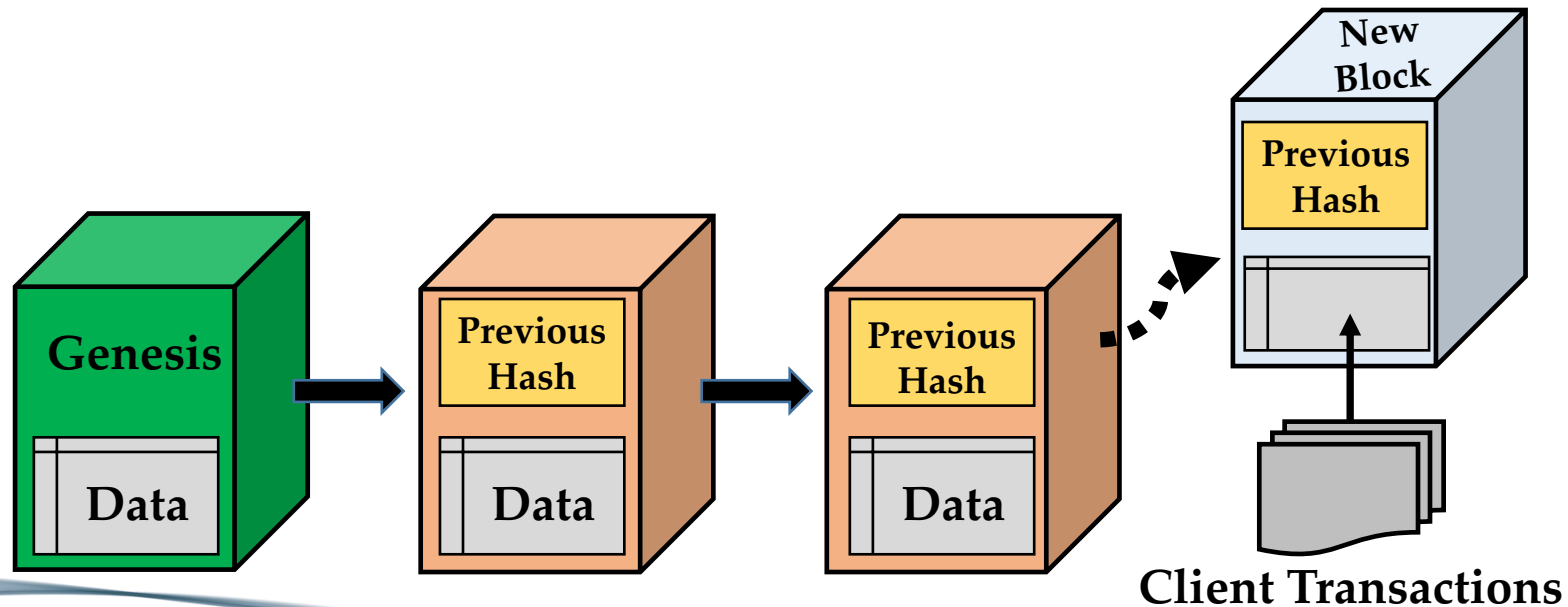
Exploratory Systems Lab

University of California Davis

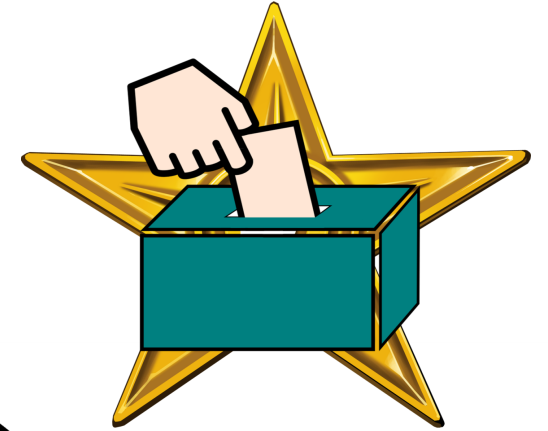
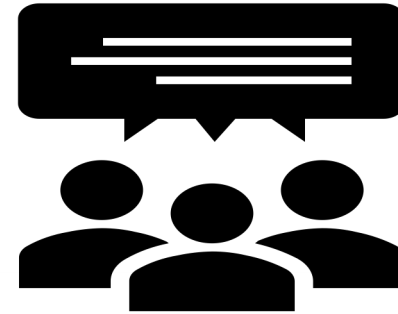
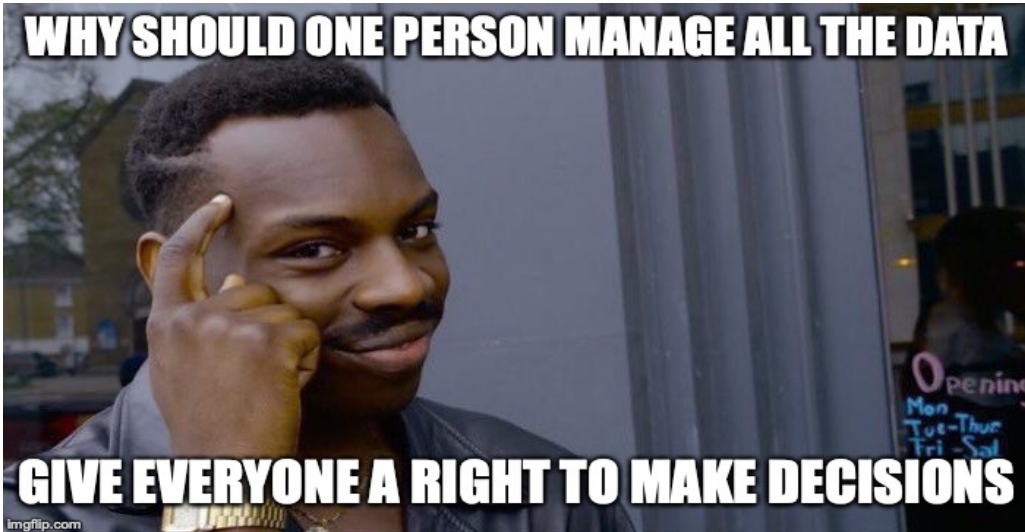


What is Blockchain?

- A linked list of blocks.
- Each block contains hash of the previous block.
- A block contains information about some client transactions.

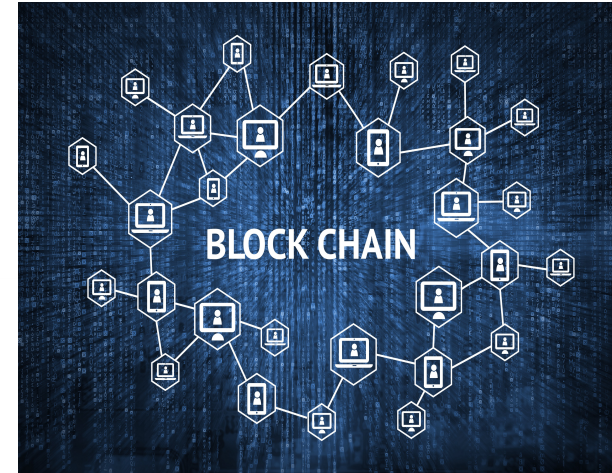


Why Blockchain?



Components of a Blockchain System

- Replicas → Store all the data.
- Client → Sends transactions to process.
- Consensus Protocol → Helps ordering transactions.
- Cryptographic Constructs → Authenticate replicas and clients.
- Ledger → Records transactions.

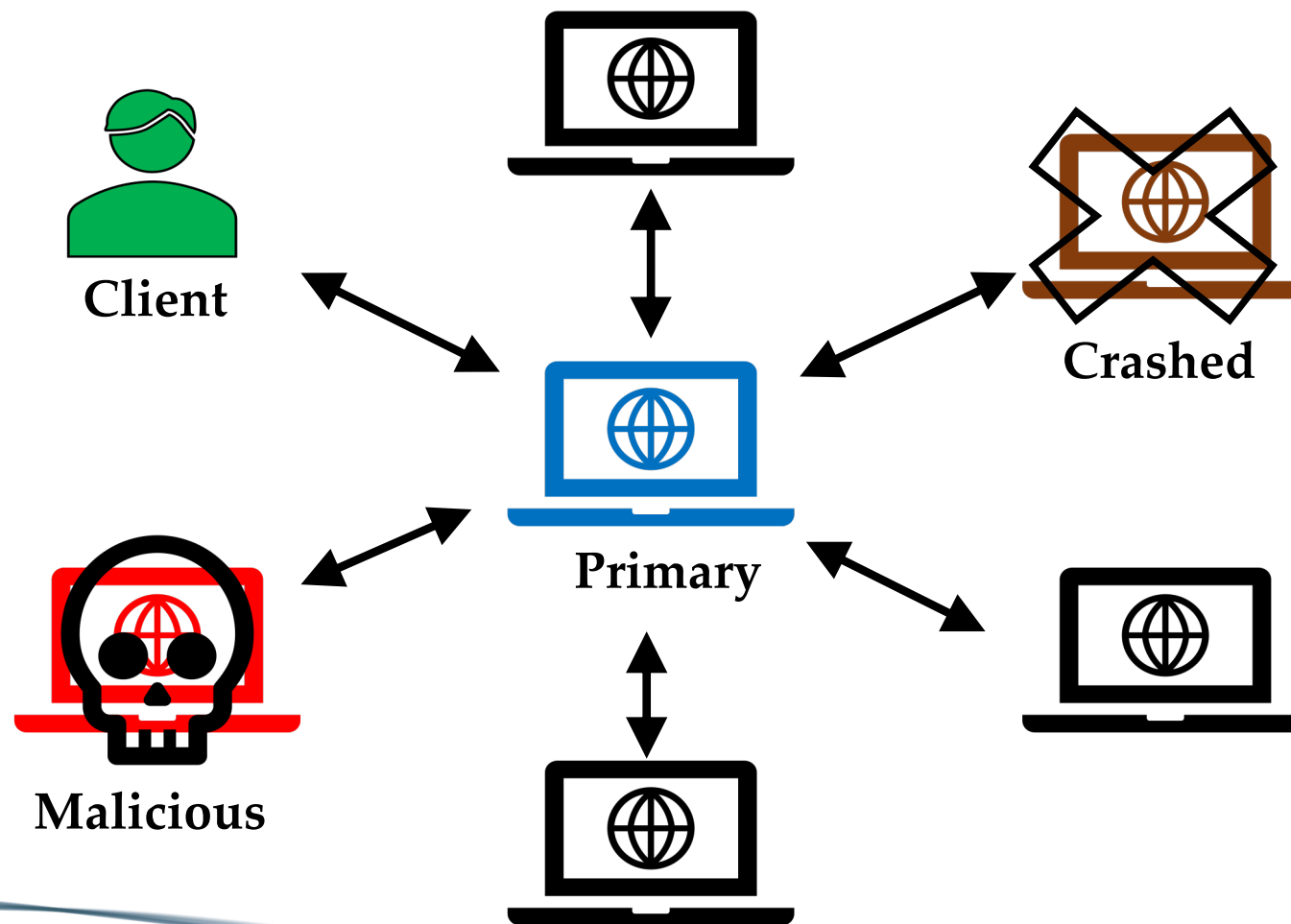


ResilientDB



ExpoLab
Creativity Unfolded

Consensus



Types of Blockchain Systems

- **Permissionless → Open Access**
 - Anyone can participate.
 - Identities of the replicas are unknown.
 - Applications include crypto-currency and money exchange.
- **Permissioned → Restricted Access**
 - Only a select group of replicas, although untrusted can participate.
 - Identities of the replica are known a priori.
 - Applications include health-care and energy trading.



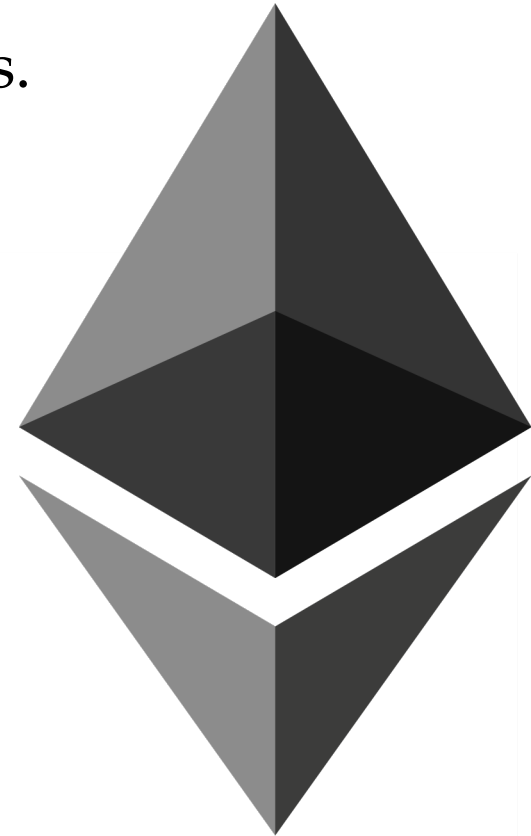
BITCOIN

- First Crypto-currency → a monetary application.
- Uses Nakamoto consensus → Proof-of-Work beneath the skin.
- Supports permissionless access.
- Requires solving hard cryptographic puzzles.
- Any replica that wants to create a new block proves that it did solve the puzzle.
- Difficulty of the puzzle helps prevent malicious attacks.



ETHEREUM

- Another Crypto-currency → a token used in variety of applications.
- Uses Proof-of-Work but plans to start using Proof-of-Stake.
- Supports permissionless access.
- Allows programmers to design their transactions or “*smart contracts*”.
- Hard dependency on Ethereum Virtual machine (EVM).
- Envisions design of Permissioned applications.



Terrorists Turn to Bitcoin for Funding, and They're Learning Fast



Login

Search Q

Disrupt Berlin 2019

THE VERGE

TECH ▾ REVIEWS ▾ SCIENCE ▾ CREATORS ▾ ENTERTAINMENT ▾ VIDEO ▾ FEAT

REPORT TECH CYBERSECURITY

Why the Ethereum Classic hack is a bad omen for the blockchain

The 51 percent attack is real, and it's easier than ever

By Russell Brandom | Jan 9, 2019, 8:47am EST

Binance says more than \$40 million in bitcoin stolen in 'large scale' hack

Zack Whittaker, Catherine Shu / 6:10 pm PDT • May 7, 2019

Comment



COINTELEGRAPH

The future of money

	BTC	ETH	XRP	BCH	LTC
	↓ \$7,424	\$152	\$0.22	\$216	\$47.98
	-1.98%	+0.03%	-0.75%	-1.00%	+1.24%

News ▾ Features ▾ Price Analysis ▾ Market Tools ▾ Cryptopedia ▾ Industry ▾



Ledger It's Black Friday all week! 30% off sitew



By William Suberg

NOV 16, 2019

Bitcoin Cash Hard Fork Sees Miners 'Waste' Money on 14 Invalid Blocks

coindesk

Bitcoin 24h	Ethereum 24h	XRP 24h
\$7,355.40 -2.58%	\$151.13 -0.62%	\$0.223687

Story from Tech →

Bitcoin Cash Miners Undo Attacker's Transactions With '51% Attack'

May 24, 2019 at 21:17 UTC • Updated May 25, 2019 at 10:39 UTC



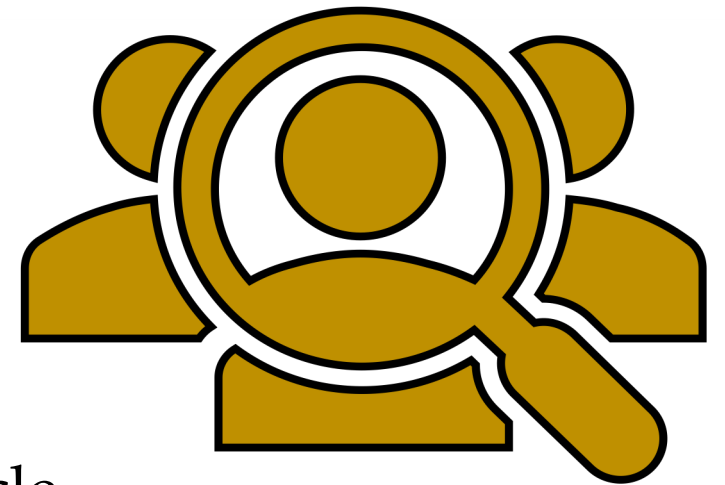
ResilientDB



ExpoLab
Creativity Unfolded

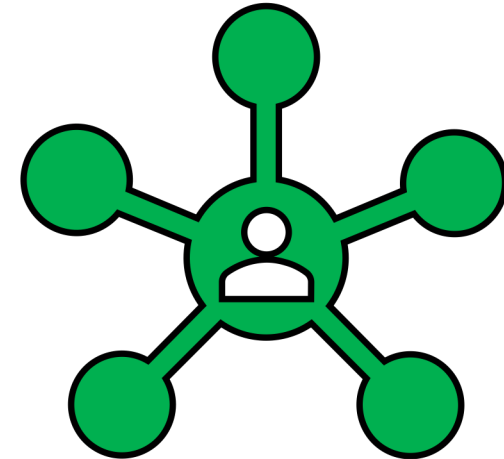
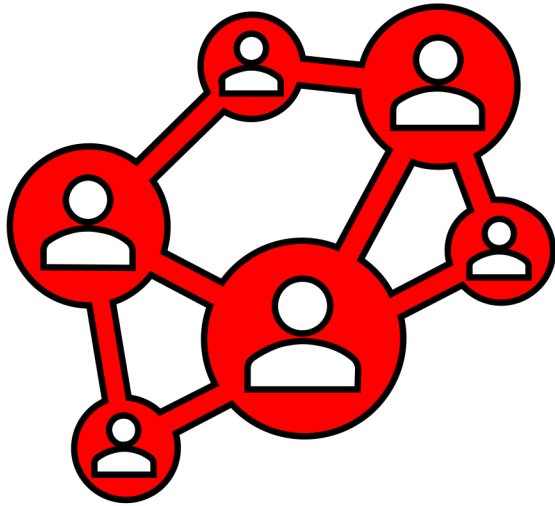
Permissioned Blockchain Systems

- Require identities of the participating replicas to be known a priori.
- Replicas still untrusted → Consensus through traditional BFT protocols.
- Computationally in-expensive.
- More reliance on *communication* primitives.
- Prevent chain forks.
- Suitable for needs of an industry → JP Morgan, IBM, Oracle
- Advent path for *Blockchain Databases*.





Transactions, Agreement and Consensus



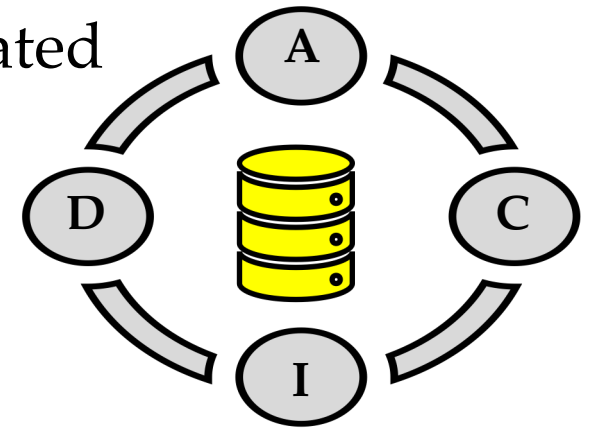
The Omniscient Transaction

- A transformation from a *consistent* state to another consistent state.
- A *contract* between two or more parties.
- A collection of *Read* or *Write* operations.
- Types of transactions: nested, compensating, multi-operation etc.



ACID Properties

- **Atomicity:** A transaction either completes fully or none of its changes take place.
- **Consistency:** The transaction must obey legal protocols
- **Isolation:** The intermediate state of a transaction is invisible to other transactions
- **Durability:** Once a transaction is committed, it cannot be abrogated



Consistency vs Availability

- An ongoing struggle that causes *performance tradeoffs*.
- **Availability** → Database needs to be always available for use.
 - **Solution?** Replication
 - **Issues?** Faults, Failures and Attacks.
- **Consistency** → Database needs to be correct.
 - **Solution?** All replicas should have same state.
 - **Issues?** Expensive.



ResilientDB



ExpoLab
Creativity Unfolded

A Deep Dive into BFT Consensus

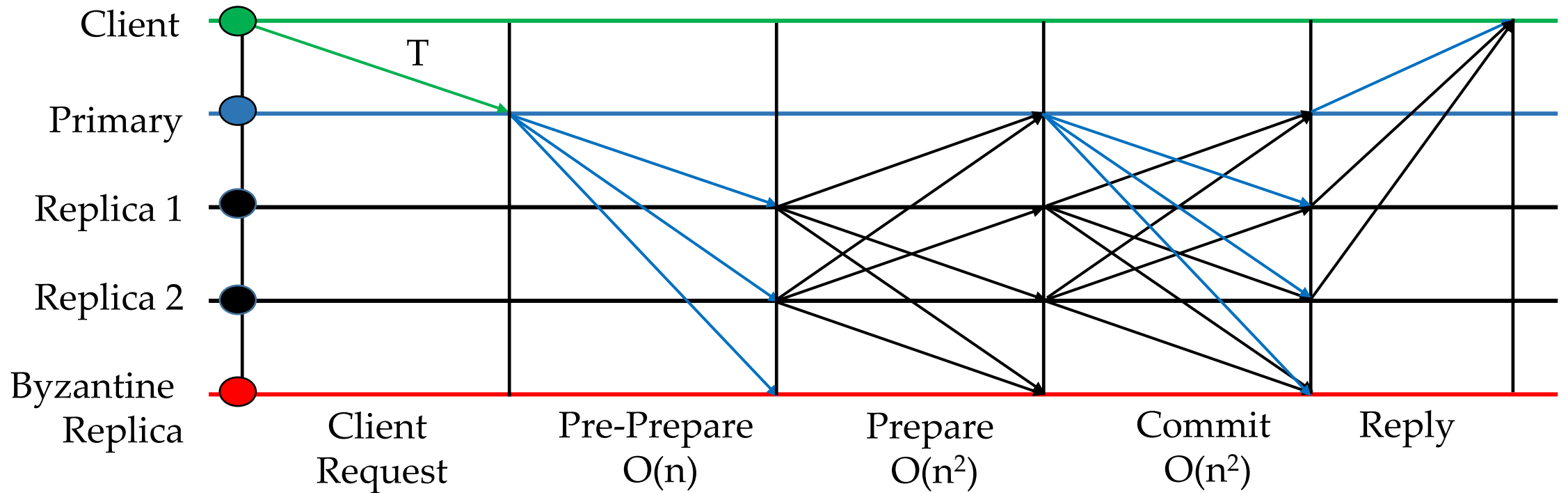


PBFT: Practical Byzantine Fault Tolerance

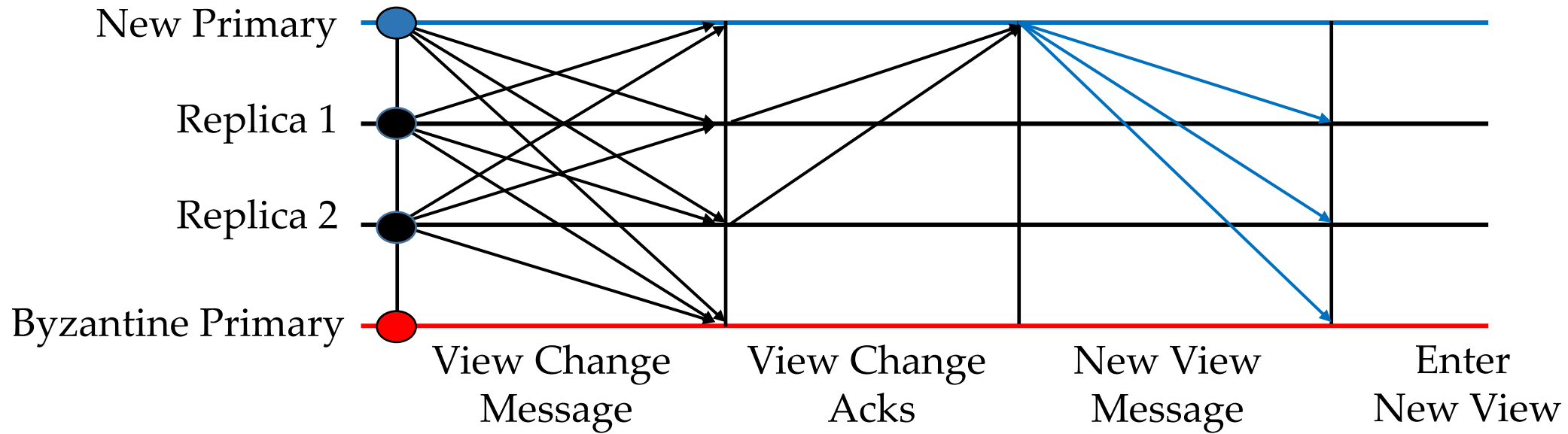
- First practical Byzantine Fault Tolerant Protocol.
- Tolerates up to f failure out of $3f+1$ replicas
- Three phases of which two require quadratic communication complexity.
- Safety is always guaranteed and Liveness is guaranteed in periods of partial synchrony.
- View-Change protocol for replacing malicious primary



PBFT Failure-Free Flow



PBFT Primary Failure (View Change)



Requirements of Existing BFT Protocols

- 1) Require three phases of communication, of which two necessitate quadratic communication (PBFT).
- 2) Expect no failures or dependence on clients (Zyzzzyva).
- 3) Incur high client latencies due to many phases of communication (PBFT, HotStuff).
- 4) Require threshold signatures, which are computationally expensive (HotStuff).
- 5) Require more than $3f+1$ replicas (Q/U, HQ).
- 6) Need trusted components (AHL, Attested Append-only memory).



Proof-of-Execution (PoE): Reaching Consensus through Fault-Tolerant Speculation

- *Speculative Execution* to reduce the client latency.
- *Out-of-Order message processing* for transactions.
- *Three Linear Phases*.
- *No Dependence* on Clients or requirement of expensive cryptographic primitives.
- *No Requirement* of a *Twin-Path* protocol.

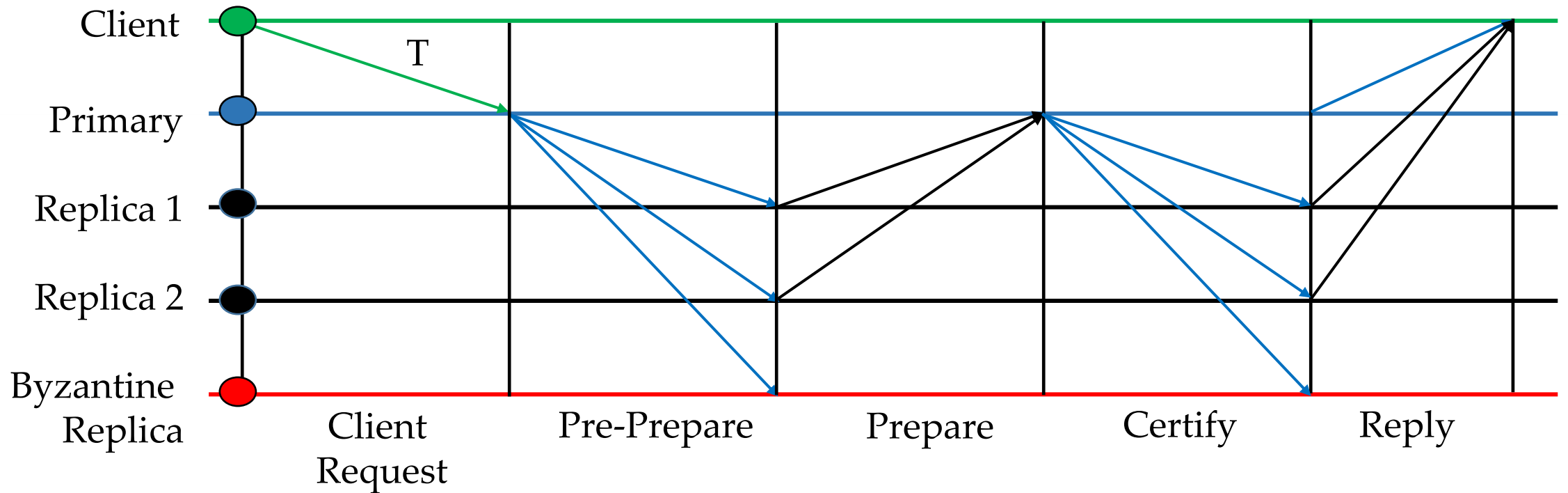


PoE vs Other Protocols

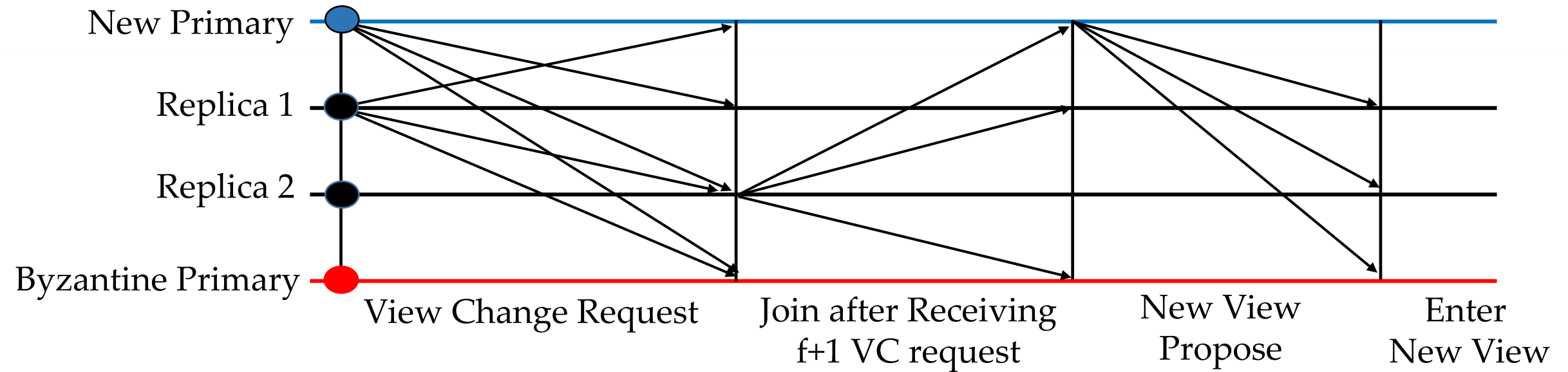
Protocol	Phases	Messages	Resilience	Requirements
ZYZZYVA	1	$\mathcal{O}(n)$	0	reliable clients and unsafe
PoE (our paper)	3	$\mathcal{O}(3n)$	f	sign. agnostic
PBFT	3	$\mathcal{O}(n + 2n^2)$	f	
HOTSTUFF	4	$\mathcal{O}(n + 3n^2)$	f	
HOTSTUFF-TS	8	$\mathcal{O}(8n)$	f	threshold sign.
SBFT	5	$\mathcal{O}(5n)$	0	threshold sign. and twin path



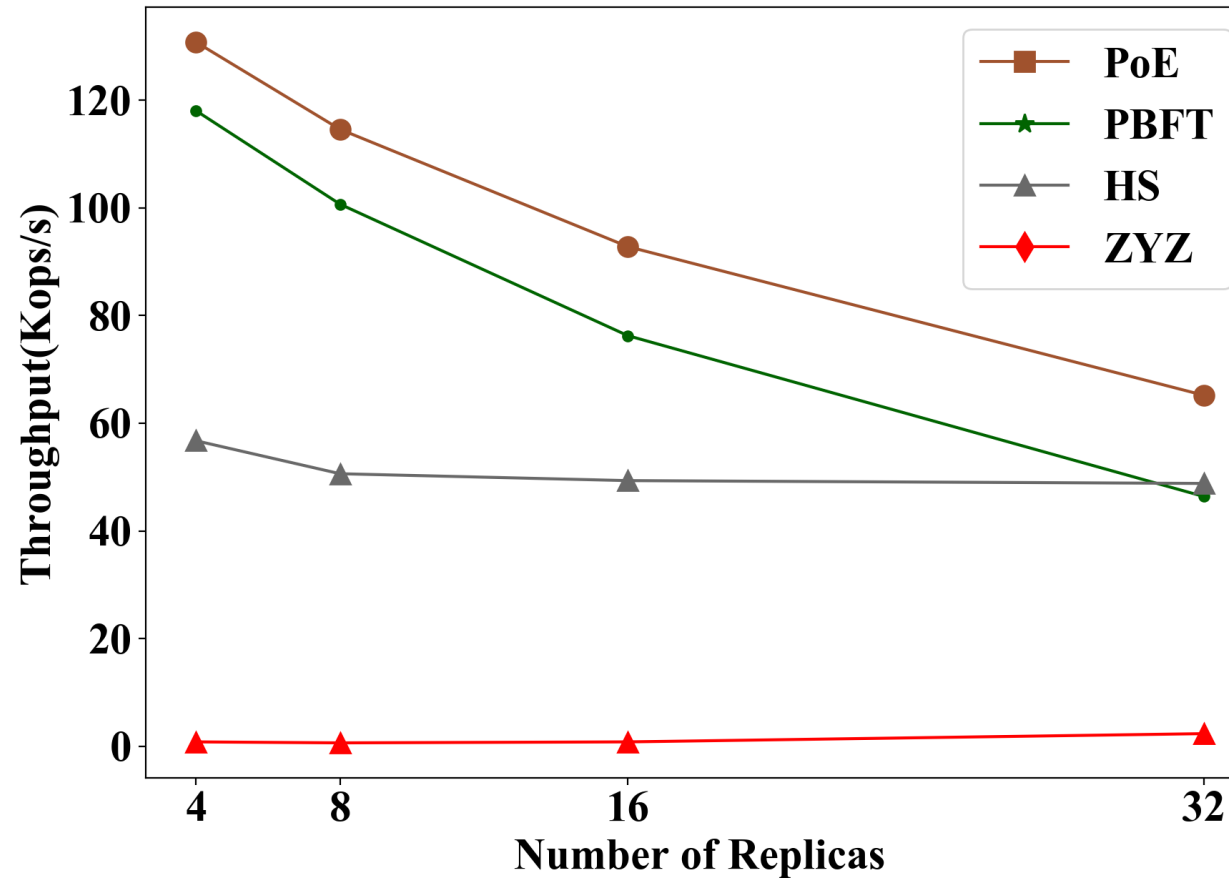
PoE Failure-Free Flow



PoE View Change Protocol



PoE Scalability under Single Failure



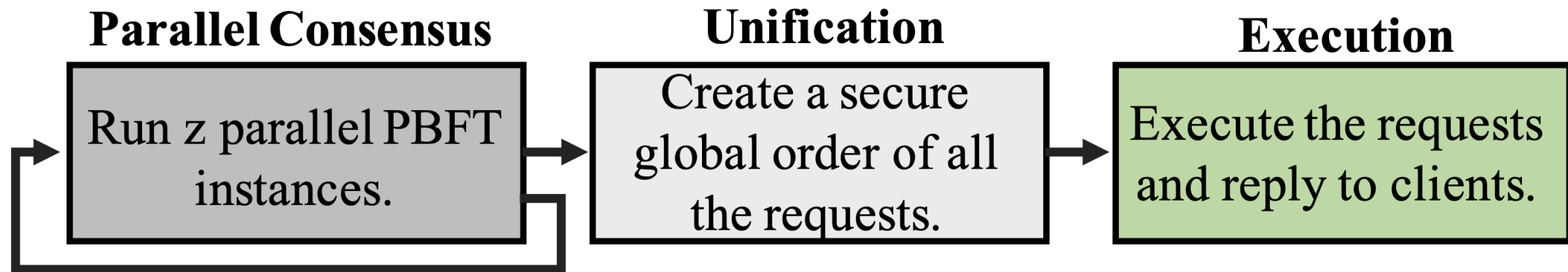
Scaling Blockchain Databases through Parallel Resilient Consensus Paradigm

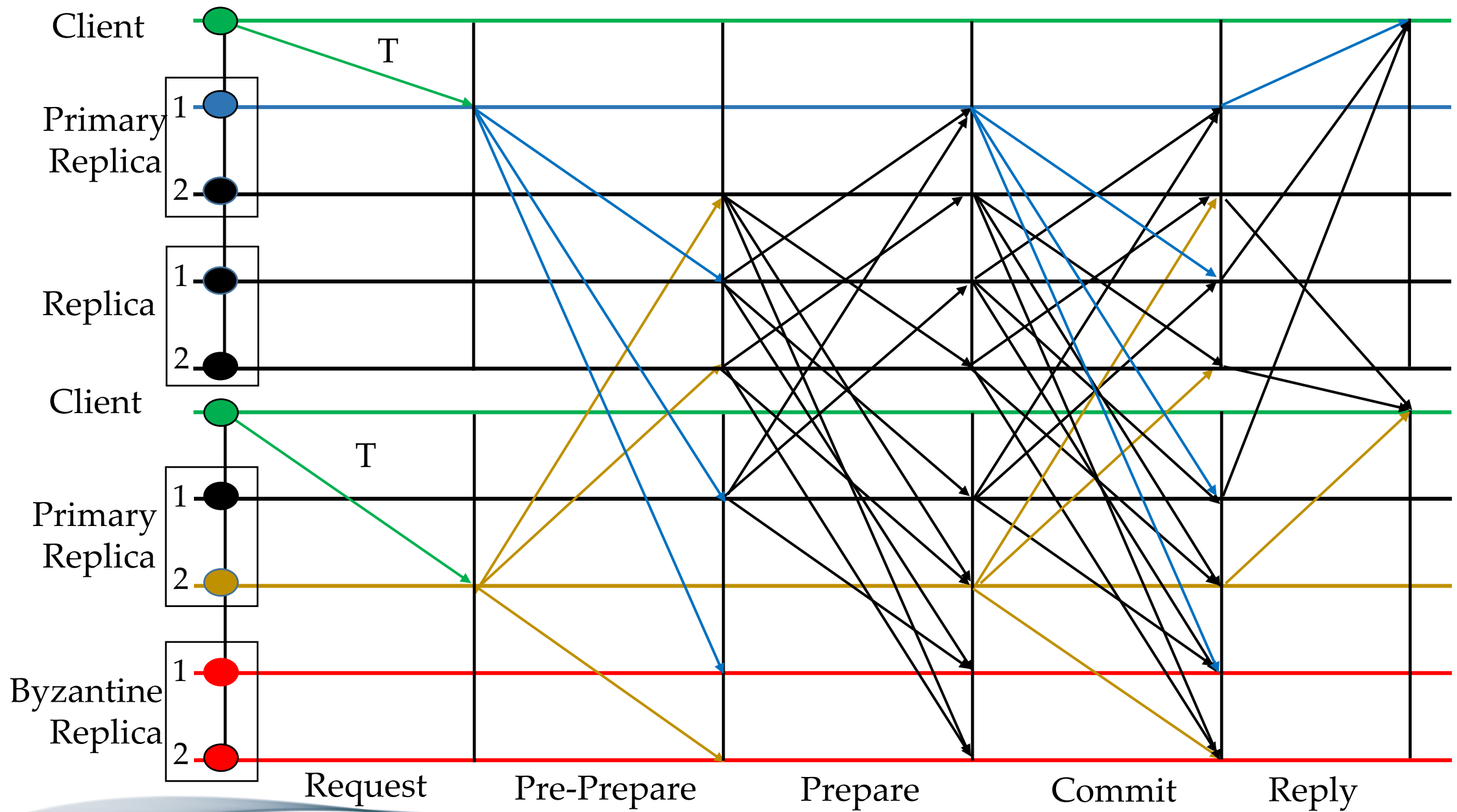
- Why should BFT protocols rely on just *one* primary replica?
- Malicious primary can *throttle* the system throughput.
- Malicious primary requires *replacement* → fall in throughput.



Multiple Byzantine Fault-Tolerance (MultiBFT) Paradigm

- Designate multiple replicas as Primaries!
- Run multiple parallel consensus on each replica.





MultiBFT with 2 parallel instances on each replica



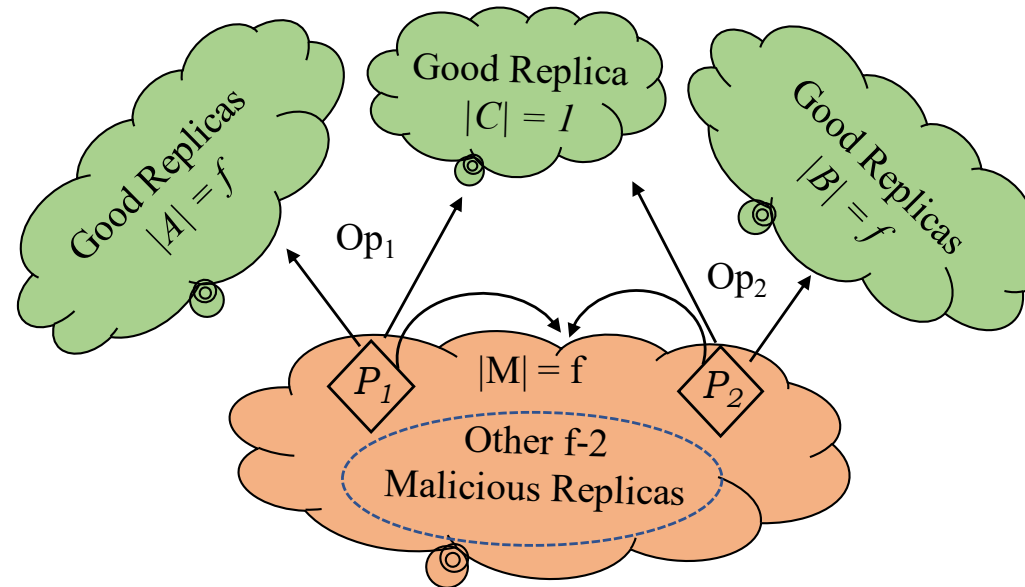
ResilientDB



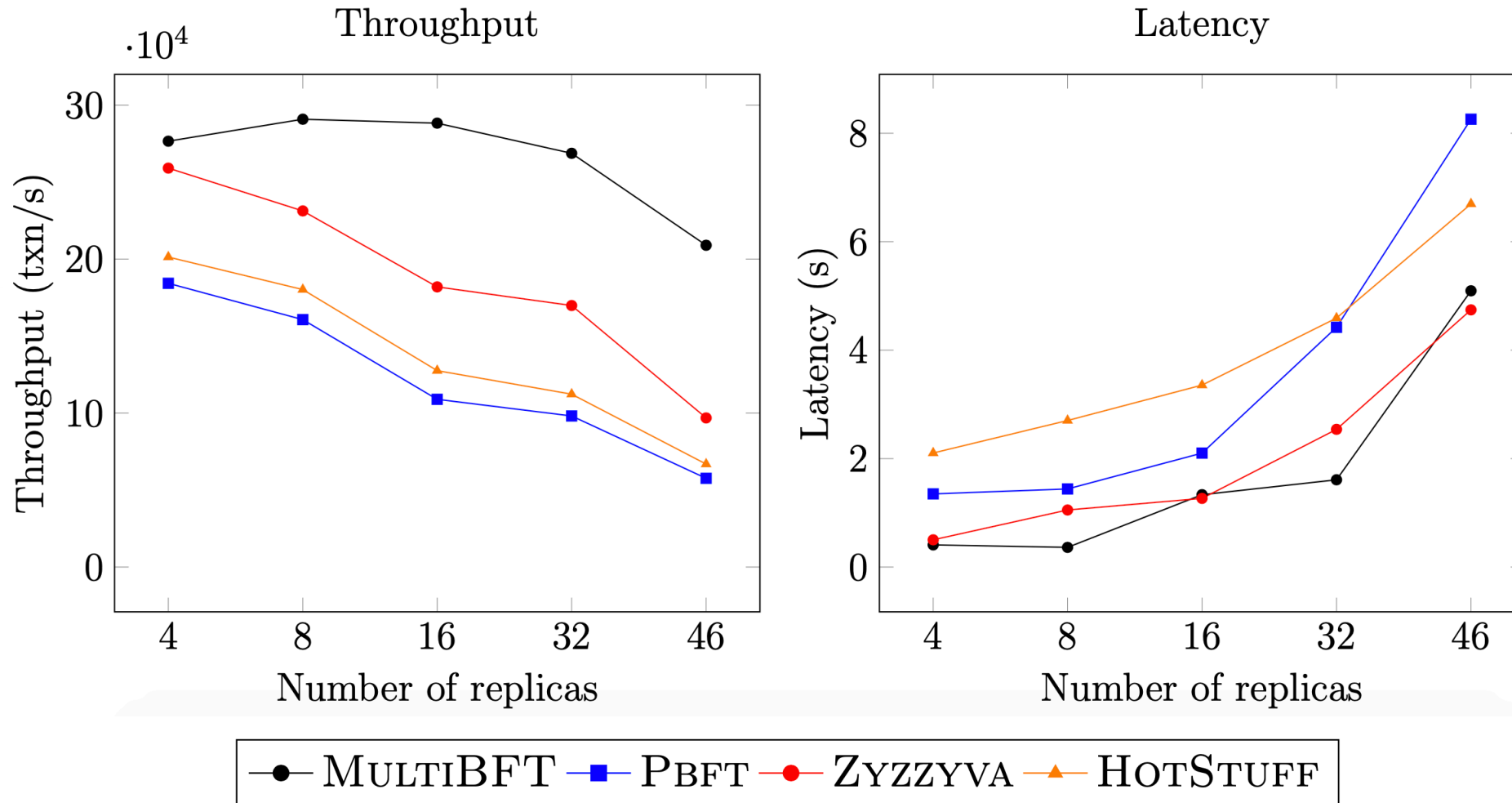
ExpoLab
Creativity Unfolded

Malicious Primaries Collusion

- Multiple malicious primaries can prevent liveness!
- Solution → Optimistic Recovery through State Exchange.



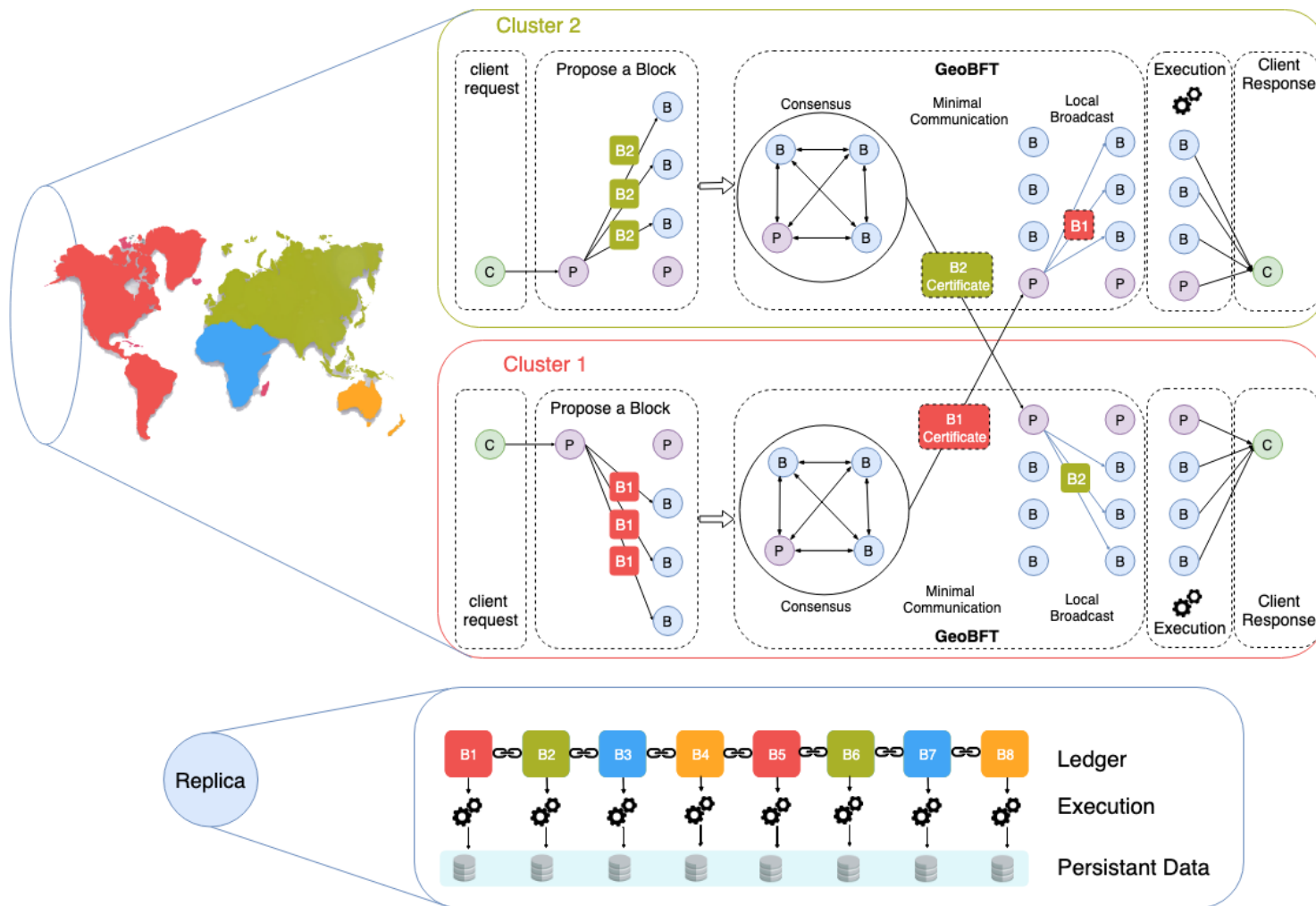
MultiBFT Scalability



Global Scale Resilient Blockchain Fabric*

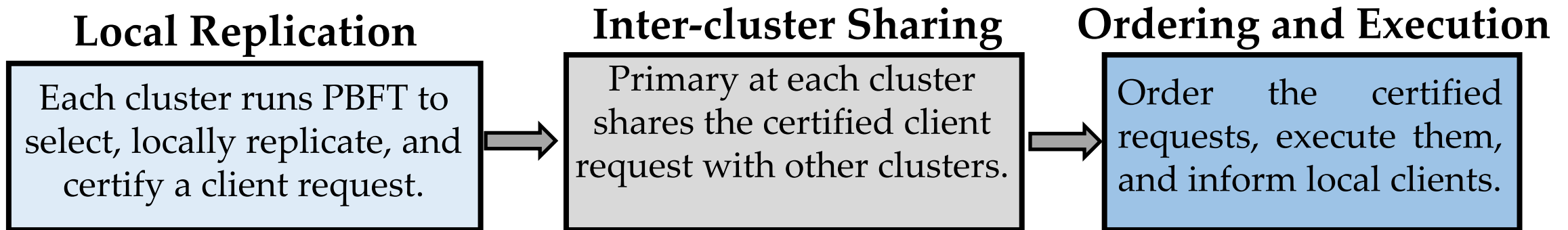
- Traditional BFT protocols do not scale to geographically large distances.
- Blockchain requires decentralization → replicas can be far apart → expensive communication!
- The underlying BFT consensus protocol should be topology-aware.

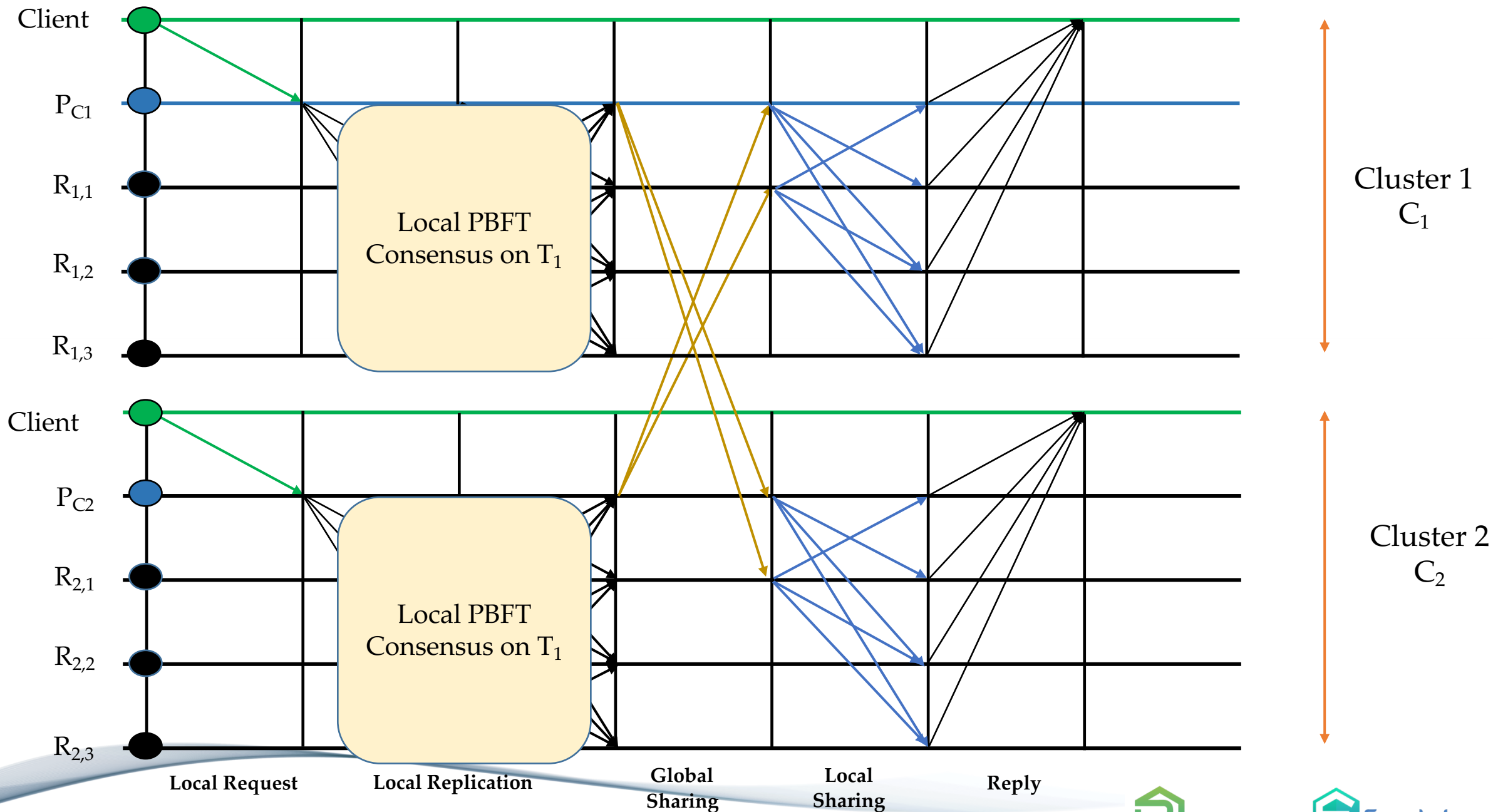
Vision Geo-Scale Byzantine Fault-Tolerance



GeoBFT Protocol

GeoBFT is a topology-aware protocol, which groups replicas into clusters. Each cluster runs the PBFT consensus protocol, in parallel and independently.



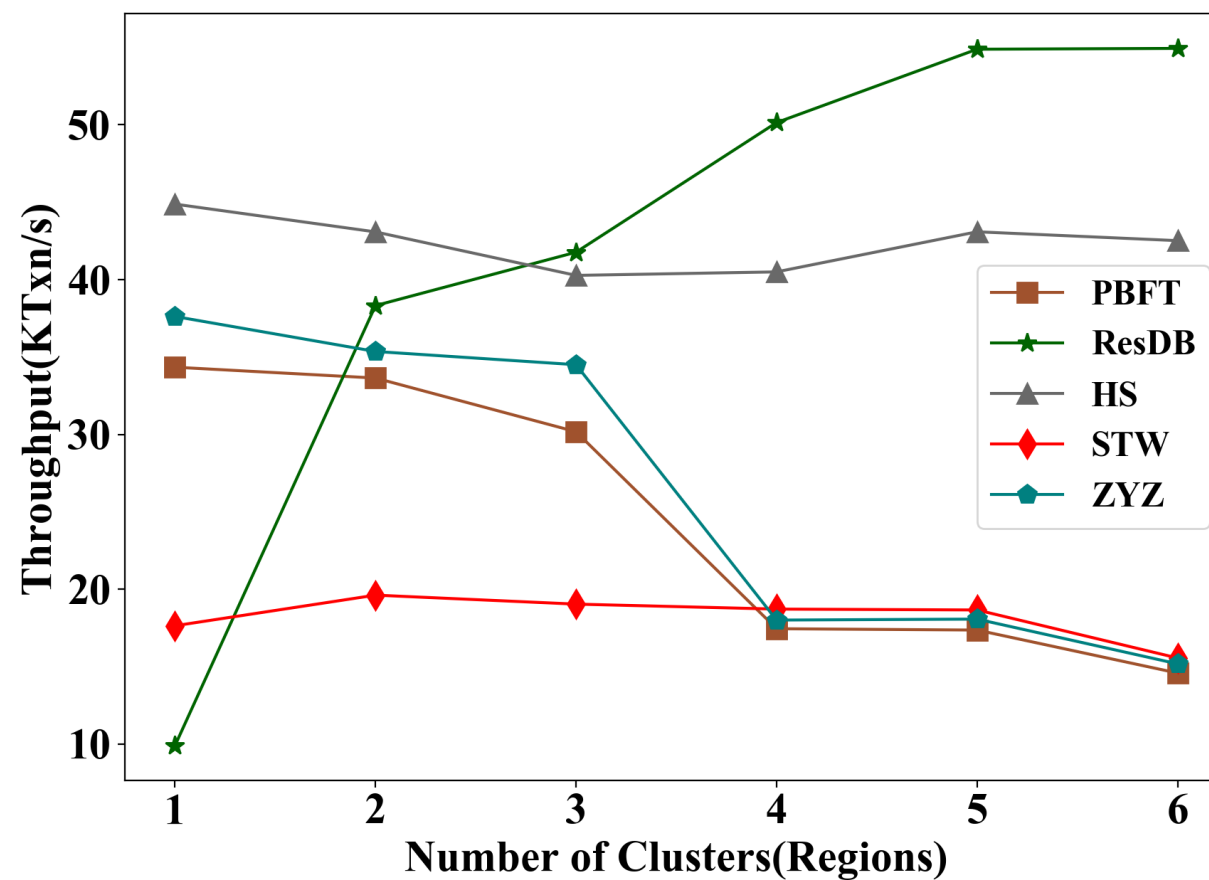


GeoBFT Takeaways

- To ensure common ordering → linear communication among the clusters is required.
- Primary replica at each cluster sends a secure certificate to $f+1$ replicas of every other cluster.
- Certificates guarantee common order for execution.
- If primary sends invalid certificates → will be detected as malicious.



GeoBFT Scalability



ResilientDB: High Throughput Yielding, Scalable Permissioned Blockchain Fabric

Visit at: <https://resilientdb.com/>



Why Should You Chose ResilientDB?

- 1) Bitcoin and Ethereum offer low throughputs of *10 txns/s*.
- 2) Existing Permissioned Blockchain Databases still have low throughputs (*20K txns/s*).
- 3) Prior works blame BFT consensus as *expensive*.
- 4) System Design is mostly *overlooked*.
- 5) ResilientDB adopts *well-researched* database and system practices.

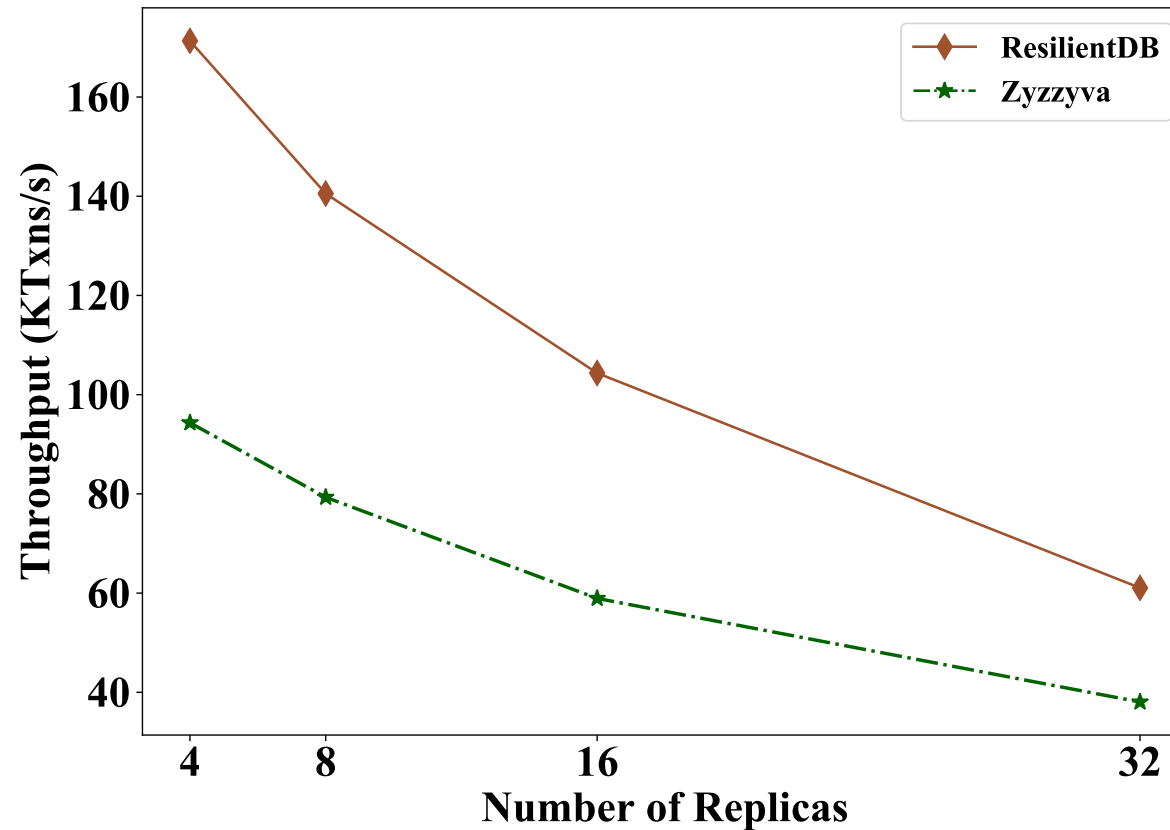


Dissecting Existing Permissioned Blockchains

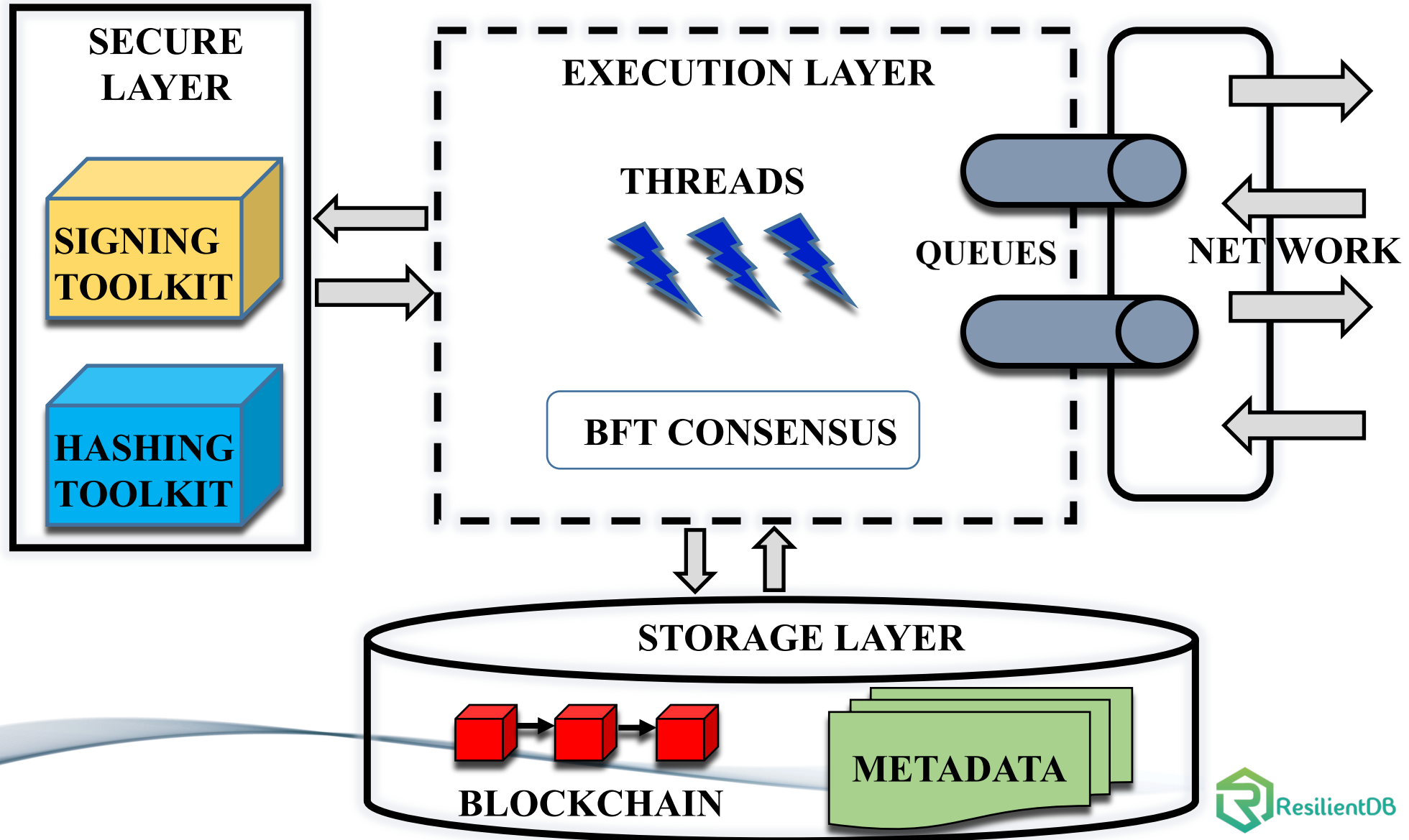
- 1) Single-threaded Monolithic Design
- 2) Successive Phases of Consensus
- 3) Integrated Ordering and Execution
- 4) Strict Ordering
- 5) Off-Chain Memory Management
- 6) Expensive Cryptographic Practices



Can a well-crafted system based on a classical BFT protocol outperform a modern protocol?



ResilientDB Architecture

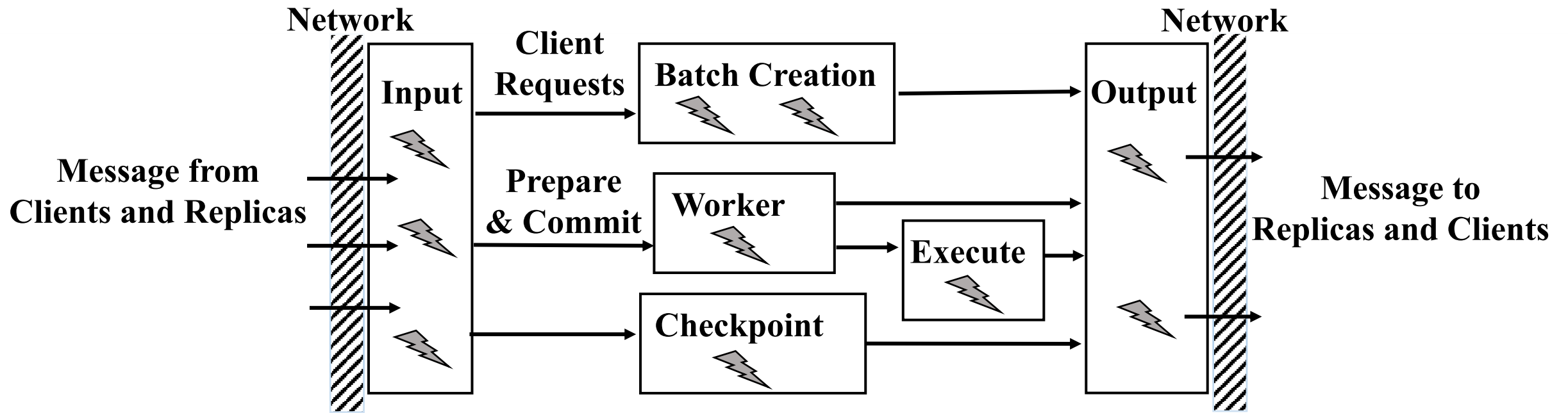


ResilientDB

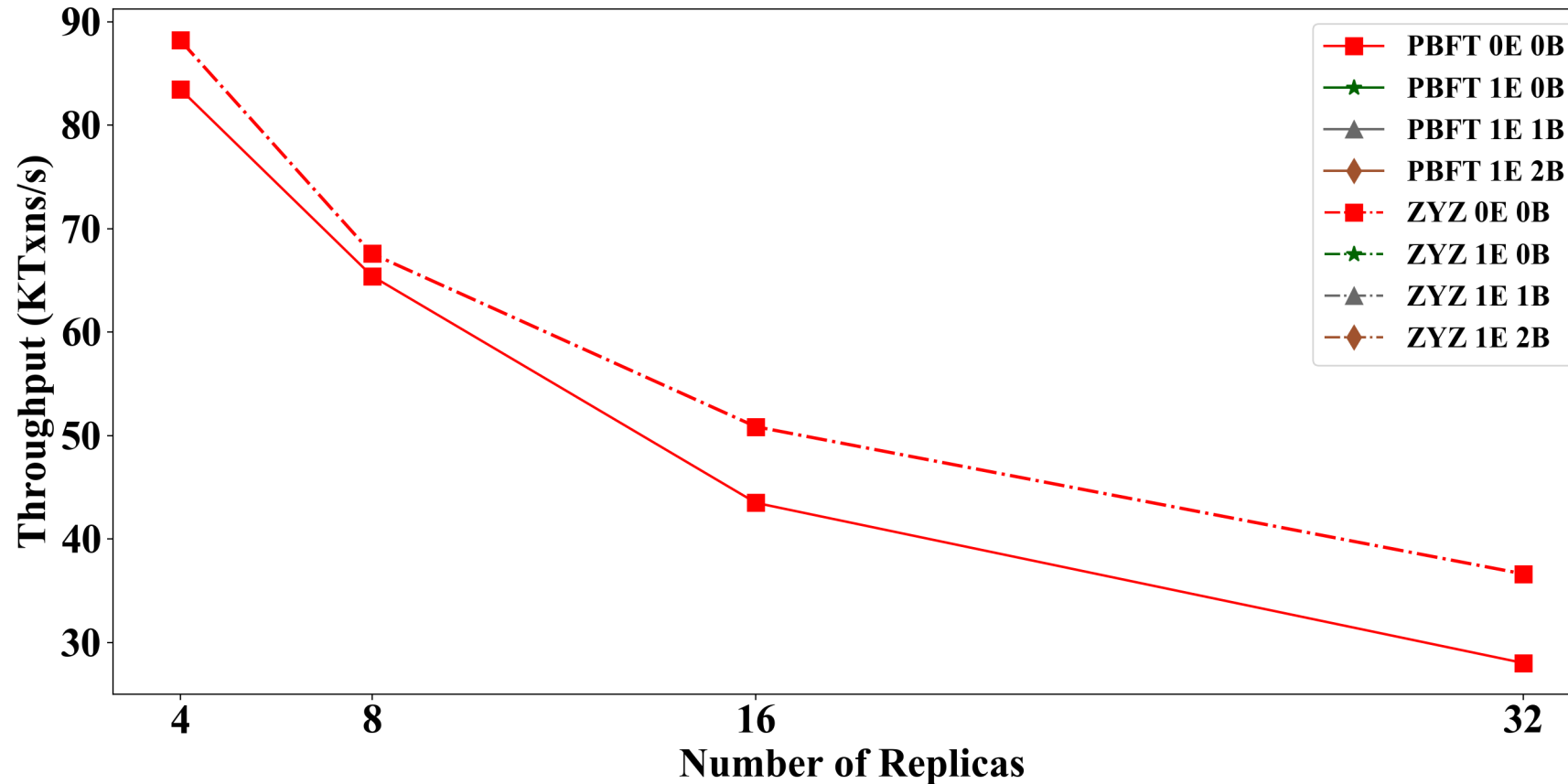


ExpoLab
Creativity Unfolded

ResilientDB Multi-Threaded Deep Pipeline

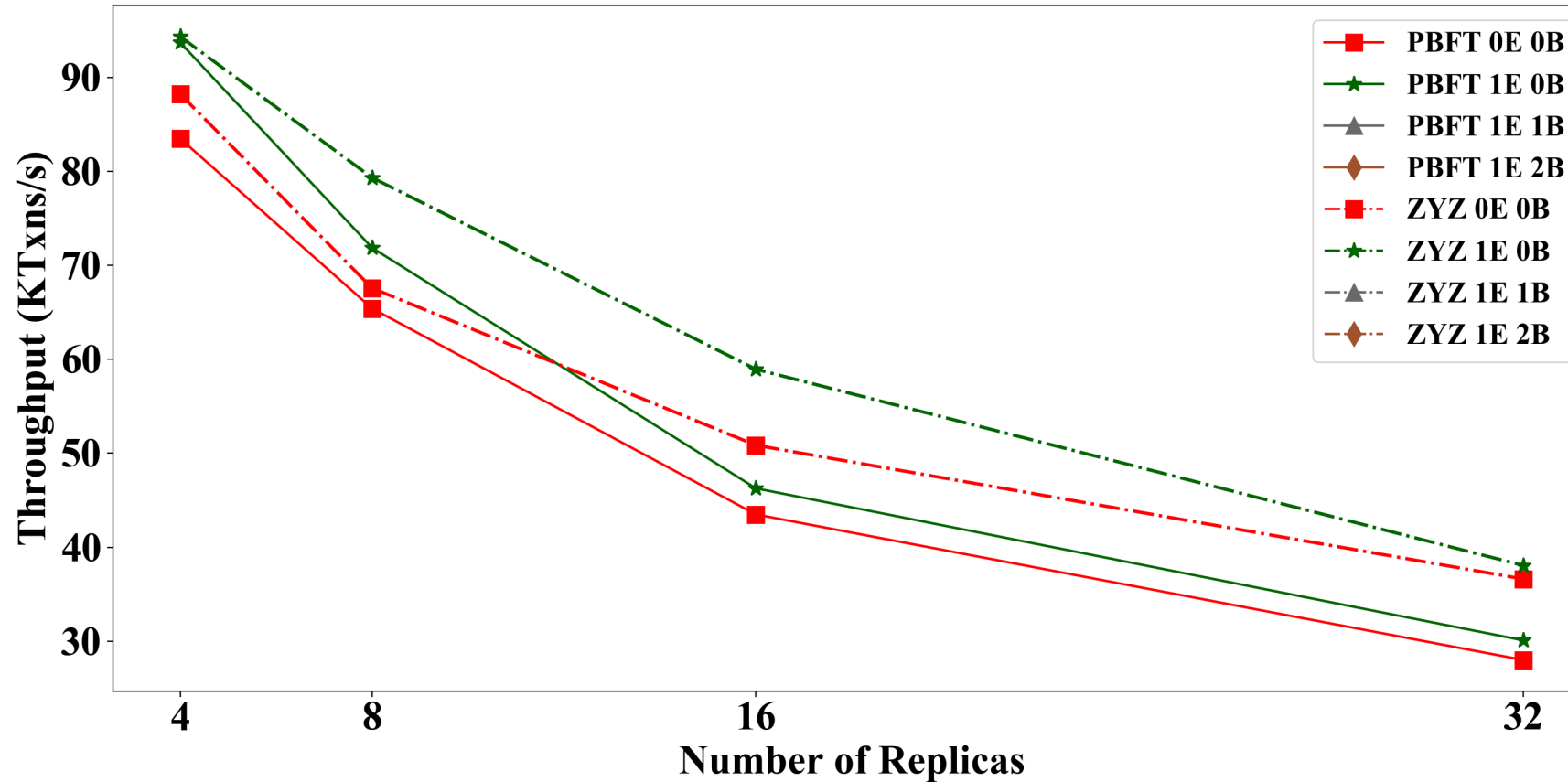


Insight 1: Multi-Threaded pipeline Gains



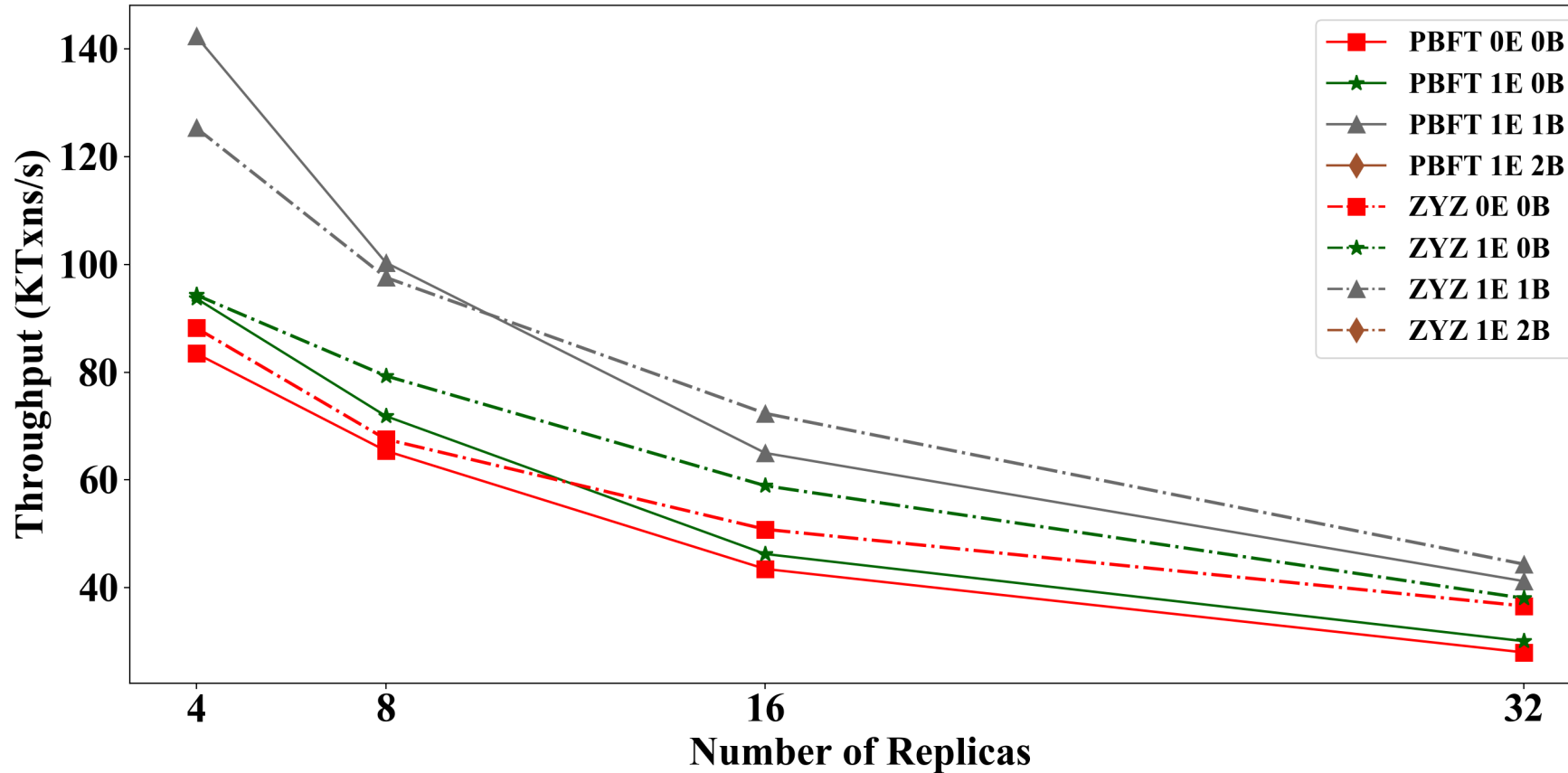
Parallelizing and Pipelining tasks across worker, execution (E) and batch-threads (B).

Insight 1: Multi-Threaded pipeline Gains



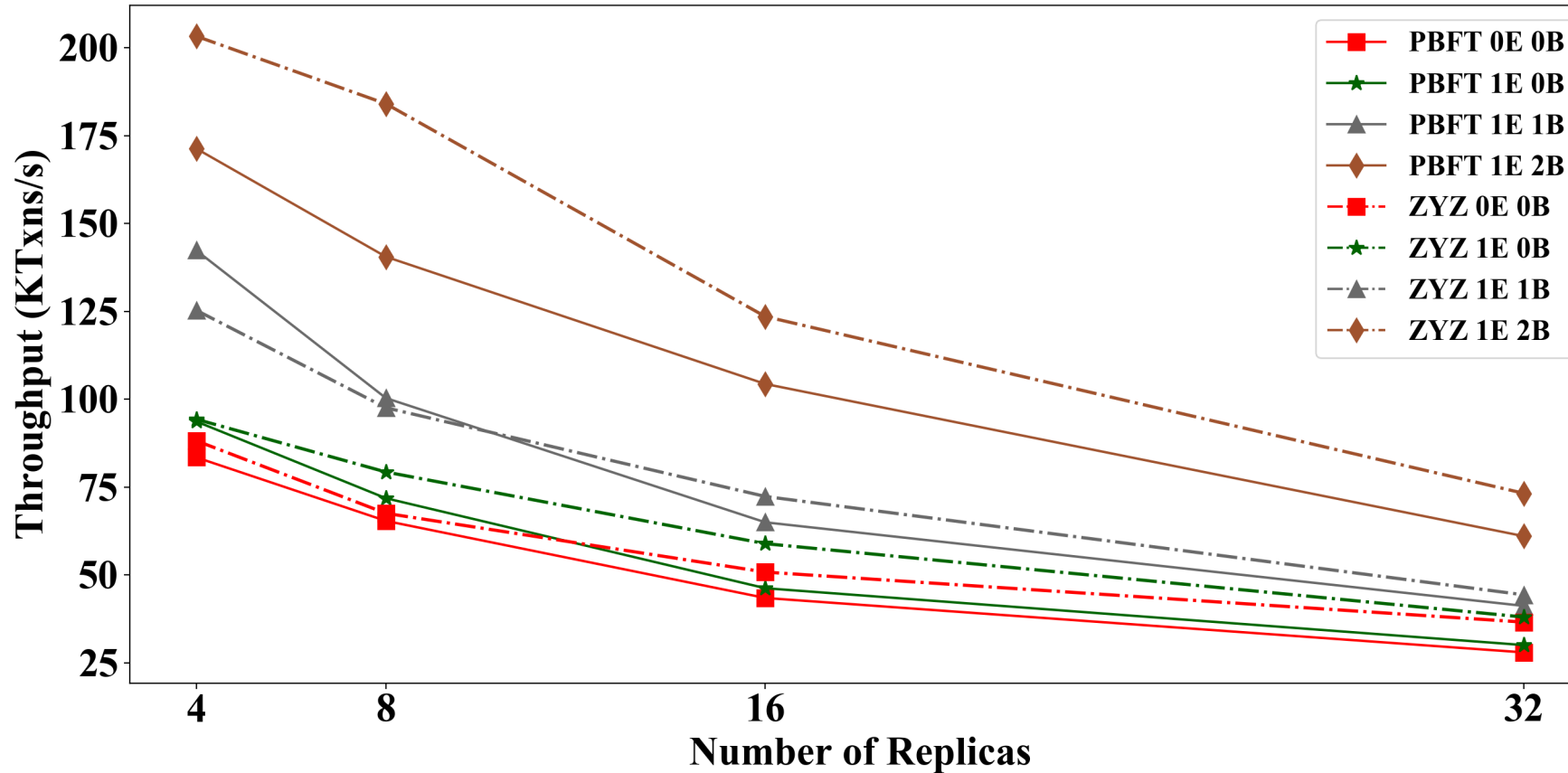
Parallelizing and Pipelining tasks across worker, execution (E) and batch-threads (B).

Insight 1: Multi-Threaded pipeline Gains



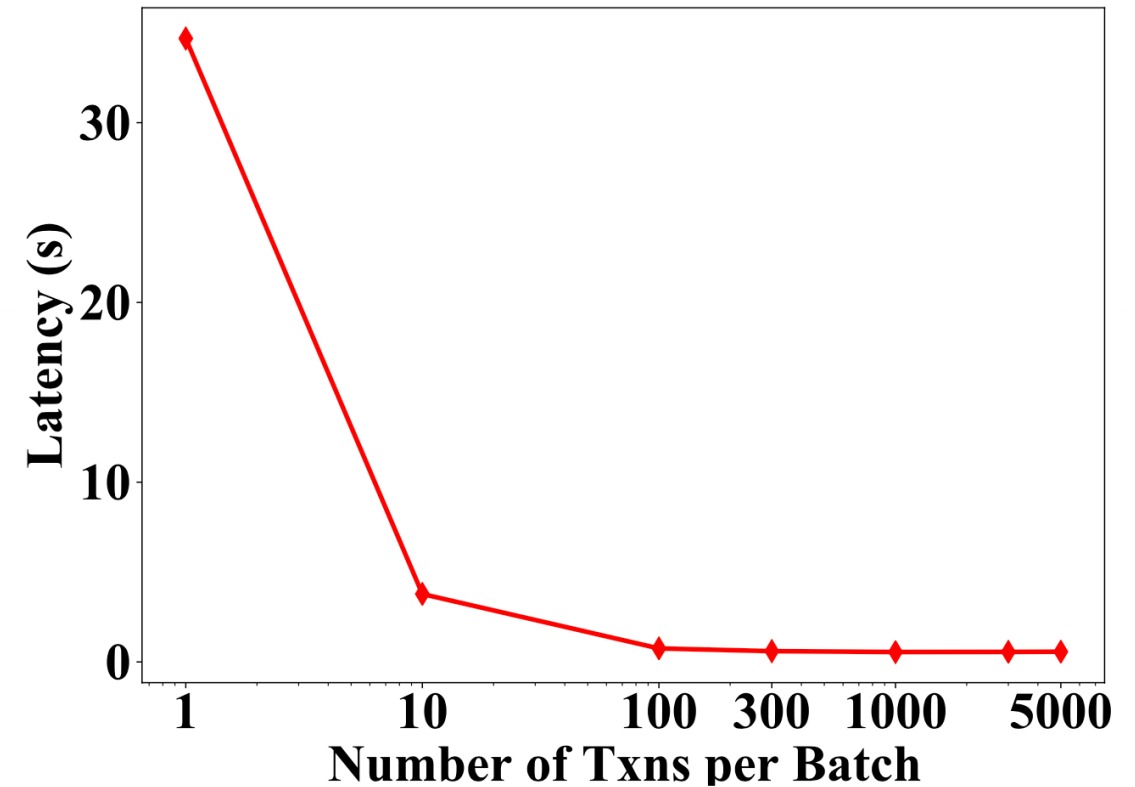
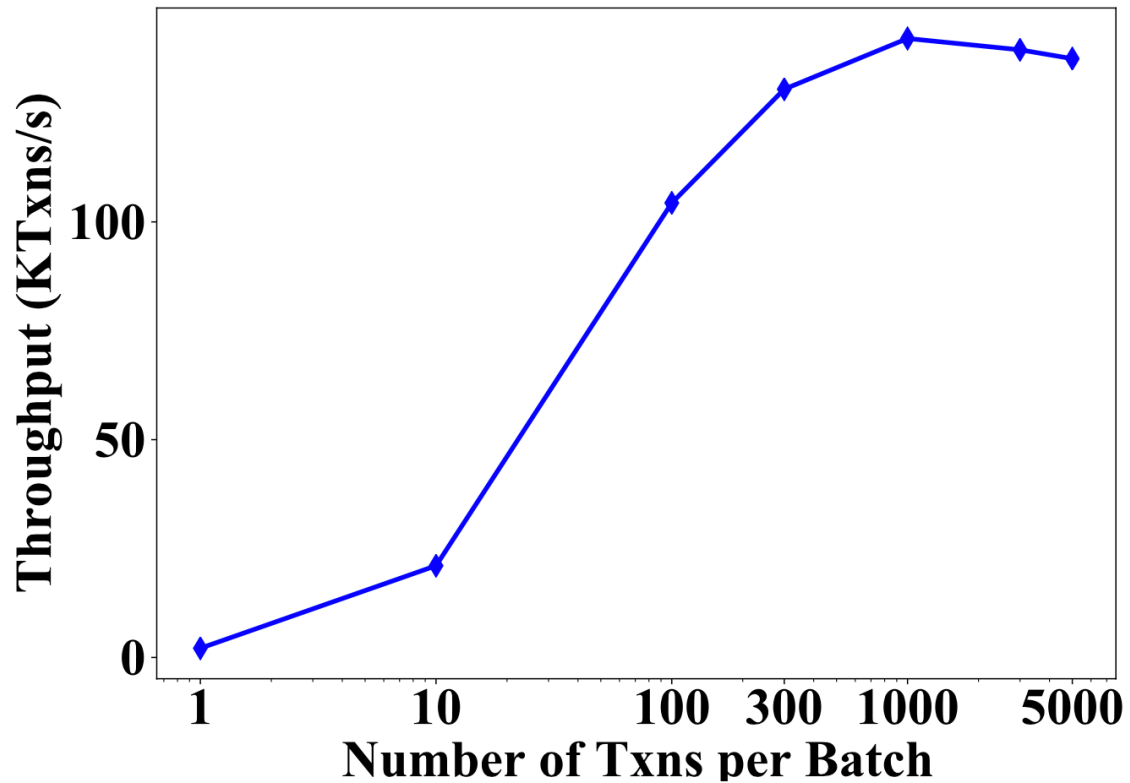
Parallelizing and Pipelining tasks across worker, execution (E) and batch-threads (B).

Insight 1: Multi-Threaded pipeline Gains



Parallelizing and Pipelining tasks across worker, execution (E) and batch-threads (B).

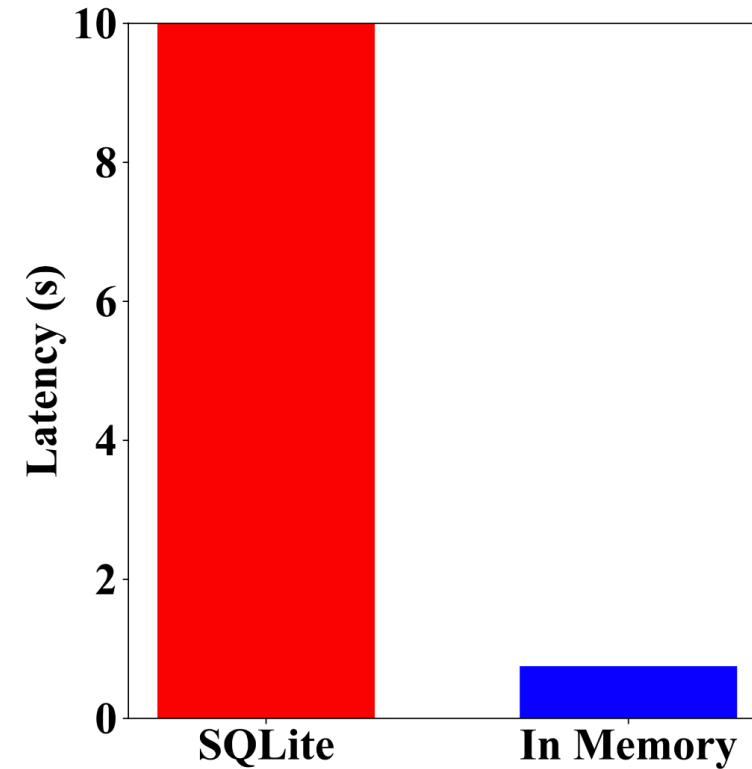
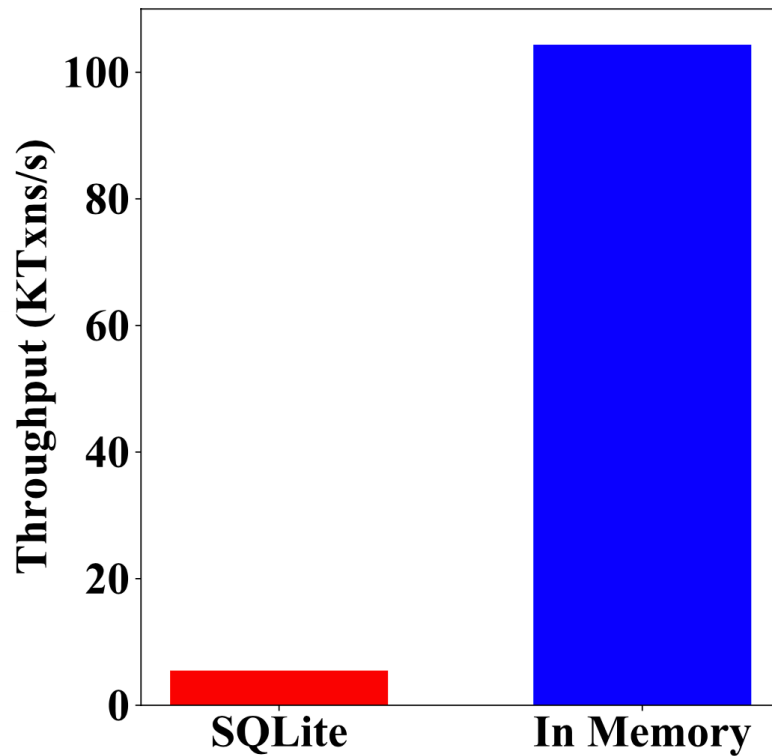
Insight 2: Optimal Batching Gains



More transactions batched together → increase in throughput
→ reduced phases of consensus.

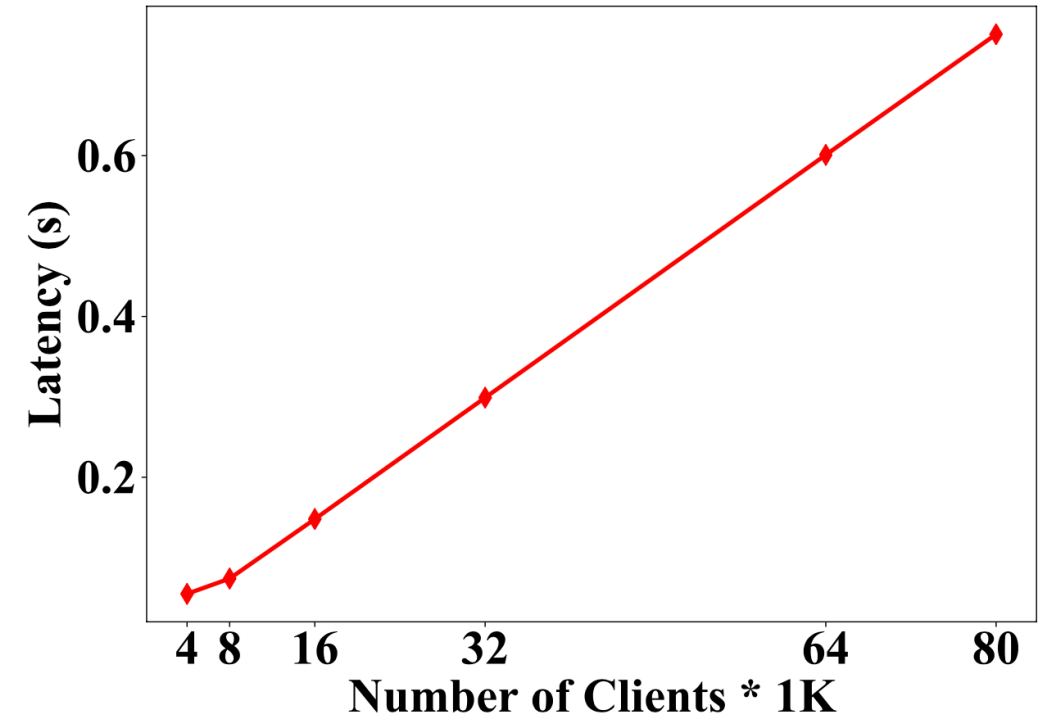
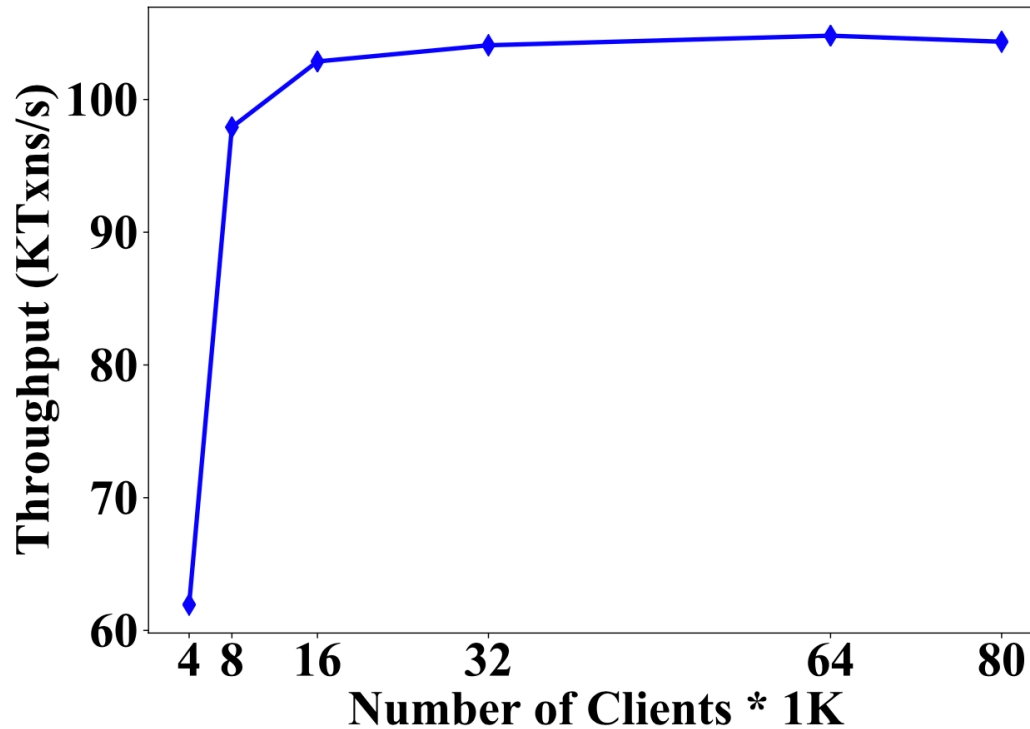


Insight 3: Memory Storage Gains



In-memory blockchain storage → reduces access cost.

Insight 4: Number of Clients



Too many clients → increases average latency.

ResilientDB: Hands On

Visit at: <https://github.com/resilientdb/resilientdb>



How to Run ResilientDB?

- Go to <https://github.com/resilientdb/resilientdb> and Fork it!
- Install Docker-CE and Docker-Compose (Links on git)
- Use the Script "*resilientDB-docker*" as following:

```
./resilientDB-docker --clients=1 --replicas=4
```

```
./resilientDB-docker -d [default 4 replicas and 1 client]
```

- Result will be printed on STDOUT and stored in *res.out* file.



How to Run ResilientDB?

resilientdb / resilientdb

Watch

5

Unstar

11

Fork

13

<> Code

Issues 1

Pull requests 0

Actions

Projects 0

Wiki

Security

Insights

46 commits

1 branch

0 packages

2 releases

4 contributors

MIT

Branch: master

New pull request

Create new file

Upload files

Find file

Clone or download

 gupta-suyash readme updated

Latest commit f2302e6 3 days ago

Docker CE

What is Docker?

*an open-source project that automates the deployment of software applications inside **containers** by providing an additional layer of abstraction and automation of **OS-level virtualization** on Linux.*

- Run a distributed program on one machine
- Simulate with lightweight virtual machines



Docker CE

What is Docker?

*an open-source project that automates the deployment of software applications inside **containers** by providing an additional layer of abstraction and automation of **OS-level virtualization** on Linux.*

- Run a distributed program on one machine
- Simulate with lightweight virtual machines



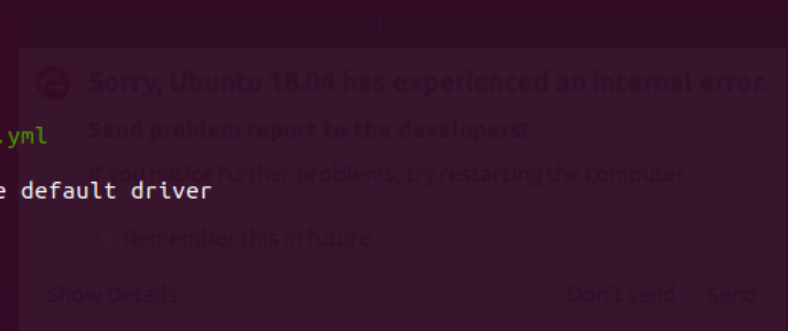
Resilient DB

./resilientDB-docker -d

- Remove old Containers
- Create new Containers
- Create IP address settings
- Install dependencies
- Compile Code
- Run binary files
- Gather the results

```
sajjad@sajjad-xps:~/WS/expo/resilientdb|master ⚡
> ./resilientDB-docker -d
Number of Replicas: 4
Number of Clients: 1
Stopping previous containers...
Stopping s3 ... done
Stopping s1 ... done
Stopping s4 ... done
Stopping c1 ... done
Stopping s2 ... done
Removing s3 ... done
Removing s1 ... done
Removing s4 ... done
Removing c1 ... done
Removing s2 ... done
Removing network resilientdb_default
Successfully stopped
Creating docker compose file ...
Docker compose file created --> docker-compose.yml
Starting the containers...
Creating network "resilientdb_default" with the default driver
Creating s4 ... done
Creating c1 ... done
Creating s1 ... done
Creating s2 ... done
Creating s3 ... done
ifconfig file exists... Deleting File
Deleted
Server sequence --> IP
c1 --> 172.21.0.3
s1 --> 172.21.0.4
s2 --> 172.21.0.6
s3 --> 172.21.0.2
s4 --> 172.21.0.5
Put Client IP at the bottom
ifconfig.txt Created!

Checking Dependencies...
Installing dependencies..
/home/sajjad/WS/expo/resilientdb
Dependencies has been installed
```



Resilient DB

- Throughput
 - Transaction per second
- Average Latency
 - The from client request to client reply
- Working Thread idleness
 - The time that thread is waiting
- WT0: Consensus Messages
- WT1 and WT2: Batch Threads
- WT3: checkpointing Thread
- WT4: Execute Thread

```
Throughputs:
0: 38525
1: 38530
2: 38558
3: 38551
4: 38564
Latencies:
latency 4: 0.505870

idle times:
Idleness of node: 0
Worker THD 0: 116.227
Worker THD 1: 62.0772
Worker THD 2: 62.2130
Worker THD 3: 105.098
Worker THD 4: 74.9193
Idleness of node: 1
Worker THD 0: 39.3157
Worker THD 1: 0.00000
Worker THD 2: 0.00000
Worker THD 3: 104.700
Worker THD 4: 74.8603
Idleness of node: 2
Worker THD 0: 35.0847
Worker THD 1: 0.00000
Worker THD 2: 0.00000
Worker THD 3: 102.415
Worker THD 4: 78.1078
Idleness of node: 3
Worker THD 0: 38.4452
Worker THD 1: 0.00000
Worker THD 2: 0.00000
Worker THD 3: 107.512
Worker THD 4: 77.6965
Memory:
0: 172 MB
1: 156 MB
2: 155 MB
3: 156 MB
4: 812 MB

avg thp: 4: 38541
avg lt : 1: .505
Code Ran successfully ---> res.out
```


PBFT: Practical Byzantine Fault Tolerance

Client Request

- Client/client_main.cpp
- System/client_thread.cpp
- ClientQueryBatch Class
- Process ClientBatch in primary

```
client > C++ client_main.cpp > ...
31 int main(int argc, char *argv[])
32 {
33     printf("Running client...\n\n");
34     // 0. initialize global data structure
35     parser(argc, argv);
36     assert(g_node_id >= g_node_cnt);
37     uint64_t seed = get_sys_clock();
38     srand(seed);
39     printf("Random seed: %ld\n", seed);
40
41     int64_t starttime;
42     int64_t endtime;
43     starttime = get_server_clock();
44     // per-partition malloc
45     printf("Initializing stats... ");
46     fflush(stdout);
47     stats.init(g_total_client_thread_cnt);
48     printf("Done\n");
49     printf("Initializing transport manager... ");
50     fflush(stdout);
51     tport_man.init();
52     printf("Done\n");
53     printf("Initializing client manager... ");
54     Workload *m_wl = new YCSBWorkload;
55     m_wl->Workload::init();
```

```
system > C++ client_thread.cpp > ...
79
80 RC ClientThread::run()
81 {
82     tsetup();
83     printf("Running ClientThread %ld\n", _thd_id);
84
85     while (true)
86     {
87         keyMTX.lock();
88         if (keyAvail)
89         {
90             keyMTX.unlock();
91             break;
92         }
93         keyMTX.unlock();
94     }
95
96     BaseQuery *m_query;
97     uint64_t iters = 0;
98     uint32_t num_txns_sent = 0;
99     int txns_sent[g_node_cnt];
100     for (uint32_t i = 0; i < g_node_cnt; ++i)
101         txns_sent[i] = 0;
102
103     run_starttime = get_sys_clock();
104
```


PBFT: Practical Byzantine Fault Tolerance

Process Messages

- Transport/message.cpp
- System/worker_thread.cpp
- System/worker_thread_pbft.cpp
- Worker Thread: Run function
- Worker Thread: Process function

```
C++ worker_thread.cpp ×
system > C++ worker_thread.cpp > WorkerThread::run()
626 /**
627  * Starting point for each worker thread.
628  *
629  * Each worker-thread created in the main() starts here. Each worker-thread is alive
630  * till the time simulation is not done, and continuously perform a set of tasks.
631  * These tasks involve, dequeuing a message from its queue and then processing it
632  * through call to the relevant function.
633  */
634 RC WorkerThread::run()
635 {
636     tsetup();
637     printf("Running WorkerThread %ld\n", _thd_id);
638
639     uint64_t agCount = 0, ready_starttime, idle_starttime = 0;
640
641     // Setting batch (only relevant for batching threads).
642     next_set = 0;
643
644     while (!simulation->is_done())
645     {
646         txn_man = NULL;
647         heartbeat();
648         progress_stats();
649
650         #if VIEW_CHANGES
651             // Thread 0 continuously monitors the timer for each batch.
652             if (get_thd_id() == 0)
653             {
654                 check_for_timeout();
655             }
656
657             if (g_node_id != get_current_view(get_thd_id()))
658             {
659                 check_switch_view();
660             }
661         #endif
662
663         // Dequeue a message from its work_queue.
664         Message *msg = work_queue.dequeue(get_thd_id());
```

```
C++ worker_thread.cpp ×
system > C++ worker_thread.cpp > WorkerThread::process(Message *msg)
87 void WorkerThread::process(Message *msg)
88 {
89     RC rc __attribute__((unused));
90
91     switch (msg->get_rtype())
92     {
93     case KEYEX:
94         rc = process_key_exchange(msg);
95         break;
96     case CL_BATCH:
97         rc = process_client_batch(msg);
98         break;
99     case BATCH_REQ:
100         rc = process_batch(msg);
101         break;
102     case PBFT_CHKPT_MSG:
103         rc = process_pbft_chkpt_msg(msg);
104         break;
105     case EXECUTE_MSG:
106         rc = process_execute_msg(msg);
107         break;
108     #if VIEW_CHANGES
109     case VIEW_CHANGE:
110         rc = process_view_change_msg(msg);
111         break;
112     case NEW_VIEW:
113         rc = process_new_view_msg(msg);
114         break;
115     #endif
116     case PBFT_PREP_MSG:
117         rc = process_pbft_prep_msg(msg);
118         break;
119     case PBFT_COMMIT_MSG:
120         rc = process_pbft_commit_msg(msg);
121         break;
122     default:
123         printf("Msg: %d\n", msg->get_rtype());
124         fflush(stdout);
125         assert(false);
126         break;
127     }
128 }
```

PBFT: Practical Byzantine Fault Tolerance

Process Client Message

- System/worker_thread_pbft.cpp
 - process_client_batch Function
 - Create and Send Batch Request
 - create_and_send_batchreq Function
 - Create Transactions
 - Create Digest
- BatchRequest Class
 - Pre-Prepare Message

```
C++ worker_thread_pbft.cpp ×
system > C++ worker_thread_pbft.cpp > ...
18  /**
19   * Processes an incoming client batch and sends a Pre-prepare message to all
20   *
21   * This function assumes that a client sends a batch of transactions and
22   * for each transaction in the batch, a separate transaction manager is created.
23   * Next, this batch is forwarded to all the replicas as a BatchRequests Message
24   * which corresponds to the Pre-Prepare stage in the PBFT protocol.
25   *
26   * @param msg Batch of Transactions of type ClientQueryBatch from the client.
27   * @return RC
28   */
29 RC WorkerThread::process_client_batch(Message *msg)
30 {
31     //printf("ClientQueryBatch: %ld, THD: %ld :: CL: %ld :: RQ: %ld\n", msg->txn_id,
32     //fflush(stdout);
33
34     ClientQueryBatch *clbtch = (ClientQueryBatch *)msg;
35
36     // Authenticate the client signature.
37     validate_msg(clbtch);
38
39 #if VIEW_CHANGES
40     // If message forwarded to the non-primary.
41     if (g_node_id != get_current_view(get_thd_id()))
42     {
43         client_query_check(clbtch);
44         return RCOK;
45     }
46
47     // Partial failure of Primary 0.
48     fail_primary(msg, 9);
49 #endif
50
51     // Initialize all transaction managers and
52     create_and_send_batchreq(clbtch, clbtch->txn_id);
53
54     return RCOK;
55 }
```

```
C++ worker_thread.cpp ×
system > C++ worker_thread.cpp > WorkerThread::create_and_send_batchreq(ClientQueryBatch *, uint64_t)
1123 * This function is used by the primary replicas to create and set
1124 * transaction managers for each transaction part of the ClientQueryBatch message
1125 * by the client. Further, to ensure integrity a hash of the complete batch is
1126 * generated, which is also used in future communication.
1127 *
1128 * @param msg Batch of transactions as a ClientQueryBatch message.
1129 * @param tid Identifier for the first transaction of the batch.
1130 */
1131 void WorkerThread::create_and_send_batchreq(ClientQueryBatch *msg, uint64_t tid)
1132 {
1133     // Creating a new BatchRequests Message.
1134     Message *bmsg = Message::create_message(BATCH_REQ);
1135     BatchRequests *breq = (BatchRequests *)bmsg;
1136     breq->init(get_thd_id());
1137
1138     // Starting index for this batch of transactions.
1139     next_set = tid;
1140
1141     // String of transactions in a batch to generate hash.
1142     string batchStr;
1143
1144     // Allocate transaction manager for all the requests in batch.
1145     for (uint64_t i = 0; i < get_batch_size(); i++)
1146     {
1147         uint64_t txn_id = get_next_txn_id() + i;
1148
1149         //cout << "Txn: " << txn_id << " :: Thd: " << get_thd_id() << "\n";
1150         //fflush(stdout);
1151         txn_man = get_transaction_manager(txn_id, 0);
1152
1153         // Unset this txn man so that no other thread can concurrently use.
1154         while (true)
1155         {
1156             bool ready = txn_man->unset_ready();
1157             if (!ready)
1158             {
1159                 continue;
1160             }
1161             else
1162             {
1163                 break;
1164             }
1165         }
1166     }
```

PBFT: Practical Byzantine Fault Tolerance

Process Batch Request (Prepare)

- System/worker_thread_pbft.cpp
- process_batch Function
- Create and Send Prepare Message
 - Create Transactions
 - Save Digest
- PBFTPrepare Class
 - Prepare Message

```
C++ worker_thread_pbft.cpp ×
system > C++ worker_thread_pbft.cpp > WorkerThread::process_batch(Message *)

57  /**
58   * Process incoming BatchRequests message from the Primary.
59   *
60   * This function is used by the non-primary or backup replicas to process an incoming
61   * BatchRequests message sent by the primary replica. This processing would require
62   * sending messages of type PBFTPrepMessage, which correspond to the Prepare phase of
63   * the PBFT protocol. Due to network delays, it is possible that a replica may have
64   * received some messages of type PBFTPrepMessage and PBFTCommitMessage, prior to
65   * receiving this BatchRequests message.
66   *
67   * @param msg Batch of Transactions of type BatchRequests from the primary.
68   * @return RC
69   */
70  RC WorkerThread::process_batch(Message *msg)
71  {
72      uint64_t cntime = get_sys_clock();
73
74      BatchRequests *breq = (BatchRequests *)msg;
75
76      //printf("BatchRequests: TID:%ld : VIEW: %ld : THD: %ld\n",breq->txn_id, breq->view, get_
77      //fflush(stdout);
78
79      // Assert that only a non-primary replica has received this message.
80      assert(g_node_id != get_current_view(get_thd_id()));
81
82      // Check if the message is valid.
83      validate_msg(breq);
84  }
```

PBFT: Practical Byzantine Fault Tolerance

Process Prepare and Commit Messages(Prepare)

- System/worker_thread_pbft.cpp
- process_pbft_prepare Function
 - Count Prepare Messages
 - Create and Send commit Message
 - PBFTCommit Message
- process_pbft_commit Function
 - Count commit messages
 - Create and Send execute Message
 - ExecuteMessage Class

```
C++ worker_thread_pbft.cpp X
system > C++ worker_thread_pbft.cpp > ...

186  /**
187   * Processes incoming Prepare message.
188   *
189   * This functions precessing incoming messages of type PBFTPrepMessage. If
190   * received 2f identical Prepare messages from distinct replicas, then it c
191   * and sends a PBFTCommitMessage to all the other replicas.
192   *
193   * @param msg Prepare message of type PBFTPrepMessage from a replica.
194   * @return RC
195   */
196  RC WorkerThread::process_pbft_prep_msg(Message *msg)
197  {
198      //cout << "PBFTPrepMessage: TID: " << msg->txn_id << " FROM: " << msg->
199      //fflush(stdout);
200
201      // Start the counter for prepare phase.
202      if (txn_man->prep_rsp_cnt == 2 * g_min_invalid_nodes)
203      {
204          txn_man->txn_stats.time_start_prepare = get_sys_clock();
205      }
206
207      // Check if the incoming message is valid.
208      PBFTPrepMessage *pmsg = (PBFTPrepMessage *)msg;
209      validate_msg(pmsg);
210
211      // Check if sufficient number of Prepare messages have arrived.
212      if (prepared(pmsg))
213      {
214          // Send Commit messages.
215          txn_man->send_pbft_commit_msgs();
216
217          // End the prepare counter.
218          INC_STATS(get_thd_id(), time_prepare, get_sys_clock() - txn_man->tx
219      }
220
221      return RCOK;
222  }
```

```
C++ worker_thread_pbft.cpp X
system > C++ worker_thread_pbft.cpp > WorkerThread::process_pbft_commit_msg(Message *)

275  /**
276   * Processes incoming Commit message.
277   *
278   * This functions precessing incoming messages of type PBFTCommitMessage.
279   * received 2f+1 identical Commit messages from distinct replicas, then
280   * execute-thread to execute all the transactions in this batch.
281   *
282   * @param msg Commit message of type PBFTCommitMessage from a replica.
283   * @return RC
284   */
285  RC WorkerThread::process_pbft_commit_msg(Message *msg)
286  {
287      //cout << "PBFTCommitMessage: TID " << msg->txn_id << " FROM: " << m
288      //fflush(stdout);
289
290      if (txn_man->commit_rsp_cnt == 2 * g_min_invalid_nodes + 1)
291      {
292          txn_man->txn_stats.time_start_commit = get_sys_clock();
293      }
294
295      // Check if message is valid.
296      PBFTCommitMessage *pcmsg = (PBFTCommitMessage *)msg;
297      validate_msg(pcmsg);
298
299      txn_man->add_commit_msg(pcmsg);
300
301      // Check if sufficient number of Commit messages have arrived.
302      if (committed_local(pcmsg))
303      {
304          #if TIMER_ON
305              // End the timer for this client batch.
306              server_timer->endTimer(txn_man->hash);
307          #endif
308
309          // Add this message to execute thread's queue.
310          send_execute_msg();
311
312          INC_STATS(get_thd_id(), time_commit, get_sys_clock() - txn_man->tx
313      }
```


PBFT: Practical Byzantine Fault Tolerance

Process Execute Message

- System/worker_thread.cpp
- Internal Message
- process_execute Function
- Execute the Transactions in batch in order
- Create and send Client Response
- ClientResponse Class

```
worker_thread.cpp
system > C++ worker_thread.cpp > WorkerThread::process_execute_msg(Message *)
796 /**
797  * Execute transactions and send client response.
798  *
799  * This function is only accessed by the execute-thread, which executes the transactions
800  * in a batch, in order. Note that the execute-thread has several queues, and at any
801  * point of time, the execute-thread is aware of which is the next transaction to
802  * execute. Hence, it only loops on one specific queue.
803  *
804  * @param msg Execute message that notifies execution of a batch.
805  * @ret RC
806  */
807 RC WorkerThread::process_execute_msg(Message *msg)
808 {
809     //cout << "EXECUTE " << msg->txn_id << " :: " << get_thd_id() << "\n";
810     //fflush(stdout);
811
812     uint64_t ctime = get_sys_clock();
813
814     // This message uses txn man of index calling process_execute.
815     Message *rsp = Message::create_message(CL_RSP);
816     ClientResponseMessage *crsp = (ClientResponseMessage *)rsp;
817     crsp->init();
818
819     ExecuteMessage *emsg = (ExecuteMessage *)msg;
820
821     // Execute transactions in a shot
822     uint64_t i;
823     for (i = emsg->index; i < emsg->end_index - 4; i++)
824     {
825         //cout << "i: " << i << " :: next index: " << g_next_index << "\n";
826         //fflush(stdout);
827
828         TxnManager *tman = get_transaction_manager(i, 0);
829
830         inc_next_index();
831
832         // Execute the transaction
833         tman->run_txn();
834
835         // Commit the results.
836         tman->commit();
837
838         crsp->copy_from_txn(tman);
```

PBFT: Practical Byzantine Fault Tolerance

Work Queue

- Lock Free queues
- All the messages are being stored in these queues
- System/work_queue.cpp
- Multiple queues for different Threads
- Dequeue and Enqueue Interfaces
- Enqueue in IOThread
- Dequeue in Worker Thread

```
work_queue.cpp ×
system > C++ work_queue.cpp > ...
44 void QWorkQueue::enqueue(uint64_t thd_id, Message *msg, bool busy)
45 {
46     uint64_t starttime = get_sys_clock();
47     assert(msg);
48     DEBUG_M("QWorkQueue::enqueue work_queue_entry alloc\n");
49     work_queue_entry *entry = (work_queue_entry *)mem_allocator.align_alloc(sizeof(work_queue_entry));
50     entry->msg = msg;
51     entry->rtype = msg->rtype;
52     entry->txn_id = msg->txn_id;
53     entry->batch_id = msg->batch_id;
54     entry->starttime = get_sys_clock();
55     assert(ISSERVER || ISREPLICA);
56     DEBUG("Work Enqueue (%ld,%ld) %d\n", entry->txn_id, entry->batch_id, entry->rtype);
57
58     if (msg->rtype == CL_QRY || msg->rtype == CL_BATCH)
59     {
60         if (g_node_id == get_current_view(thd_id))
61         {
62             //cout << "Placing \n";
63             while (!new_txn_queue->push(entry) && !simulation->is_done())
64             {
65             }
66         }
67         else
68         {
69             assert(entry->rtype < 100);
70             while (!work_queue[0]->push(entry) && !simulation->is_done())
71             {
72             }
73         }
74     }
```

PBFT: Practical Byzantine Fault Tolerance

IO Thread and Transport Layer

- Multiple Input Threads
- Multiple Output Threads
- System/io_thread.cpp
- Transport Layer: TCP Sockets
- Nano Message Library
- Transport/transport.cpp

```
io_thread.cpp ×
system > C++ io_thread.cpp > ...

299 RC InputThread::server_rcv_loop()
300 {
301
302     myrand rdm;
303     rdm.init(get_thd_id());
304     RC rc = RCOK;
305     assert(rc == RCOK);
306     uint64_t starttime = 0;
307     uint64_t idle_starttime = 0;
308     std::vector<Message *> *msgs;
309     while (!simulation->is_done())
310     {
311         heartbeat();
312
313         #if VIEW_CHANGES
314         if (g_node_id != get_current_view(get_thd_id()))
315         {
316             uint64_t tid = get_thd_id() - 1;
317             uint32_t nchange = get_newView(tid);
318
319             if (nchange)
320             {
321                 set_current_view(get_thd_id(), get_current_view(get_thd_id()) + 1);
322                 set_newView(tid, false);
323             }
324         }
325         #endif
326
327         msgs = tport_man.rcv_msg(get_thd_id());
328     }
```

Configuration Parameters to Play

- **NODE_CNT** Total number of replicas, minimum 4, that is, $f=1$.
- **THREAD_CNT** Total number of threads at primary (at least 5)
- **CLIENT_NODE_CNT** Total number of clients (at least 1).
- **MAX_TXN_IN_FLIGHT** Multiple of Batch Size
- **DONE_TIMER** Amount of time to run the system.
- **BATCH_THREADS** Number of threads at primary to batch client transactions.
- **BATCH_SIZE** Number of transactions in a batch (at least 10)
- **TXN_PER_CHKPT** Frequency at which garbage collection is done.
- **USE_CRYPTO** To switch on and off cryptographic signing of messages.
- **CRYPTO_METHOD_ED25519** To use ED25519 based digital signatures.
- **CRYPTO_METHOD_CMAC_AES** To use CMAC + AES combination for authentication



Thank You