



UCDAVIS

THE JOURNEY OF BUILDING GLOBAL-SCALE RESILIENTDB BLOCKCHAIN FABRIC

Mohammad Sadoghi

URCS Seminar Series
University of Rochester
November 5, 2021



Mohammad Sadoghi
Exploratory Systems Lab
Department of Computer Science

UCDAVIS
UNIVERSITY OF CALIFORNIA





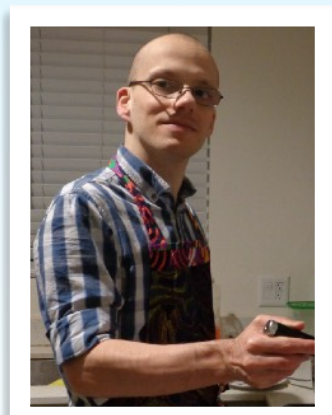
ExpoLab Team



Mohammad Sadoghi
(Principal Investigator)



Thamir Qadah, PhD
(Distributed & Coordination-free Concurrency)



Jelle Hellings, PostDoc
(Fault-tolerant Complexity Analysis)



Suyash Gupta, PhD
(Scalable Consensus Meta-Protocols)



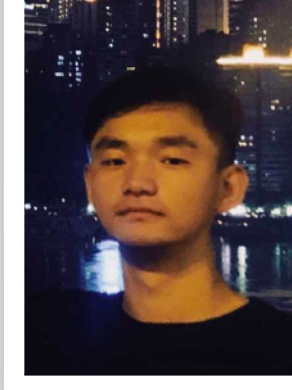
Sajjad Rahnema, PhD
(Global Scale Consensus)



Haojun (Howard) Zhu, MSc
(Re-Configurable Consensus Protocols)



Alejandro Armas, BSc
(Re-engineering ResilientDB Toolkits)



Dakai Kang, BSc
(View-change-less Protocols)



Dhruv Krishnan, MSc
(Scaling Fabric via Sharding)



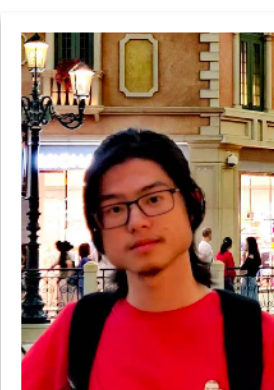
Priya Holani, MSc
(Scaling Fabric via Sharding)



Shubham Pandey, MSc
(Scaling Fabric via RDMA)

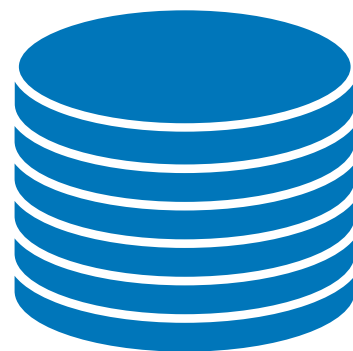


Rohan Sogani, MSc
(Scaling Fabric via Sharding)

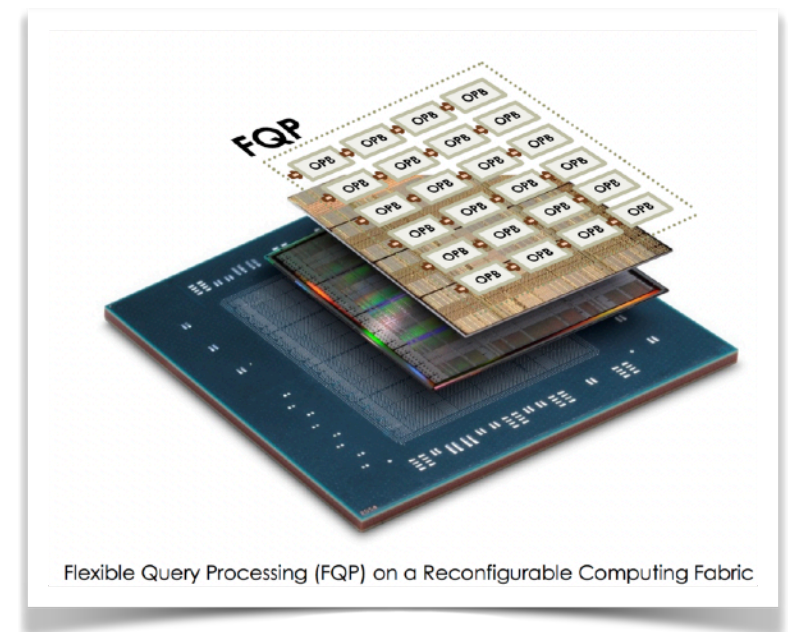


Xinyuan Sun, MSc
(Scaling Fabric via RDMA)

Resilient Journey...



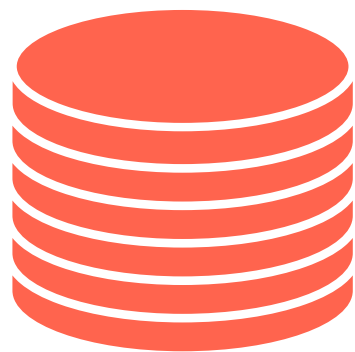
**SQL
Analytics**



FPGA Acceleration: FQP (Flexible Query Processor)

[VLDB'10, ICDE'12, VLDB'13, ICDE'15, SIGMOD Record'15, ICDE'16, USENIX ATC'16, ICDCS'17, ICDE'18, TKDE'19]

Resilient Journey...



**SQL
Transactions**

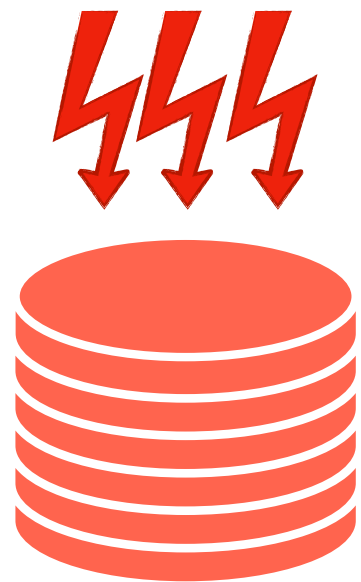


**SQL
Analytics**

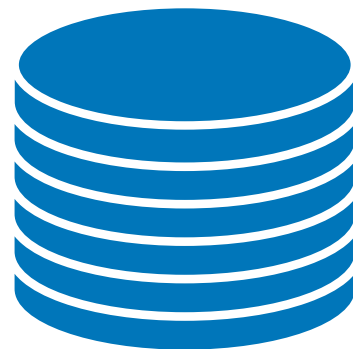
High-dimensional Indexing: (e.g., BE-Tree, BE-topK)
[SIGMOD'11, ICDE'12, TODS'13, ICDCS'13, ICDE'14, ICDCS'17, Middleware'17]

Resilient Journey...

Concurrency Protocols



**SQL
Transactions**

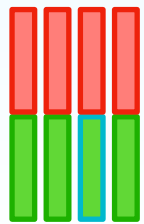


**SQL
Analytics**

Concurrency Control Protocols: (e.g., 2VCC, QueCC - Best Paper Award)
[VLDB'13, VLDB'14, VLDBJ'16, Middleware'16, TDKE'15, SIGMOD'15, ICDE'16, Middleware'18]

Resilient Journey...

QueCC: Queue-Oriented Planning and Execution Architecture

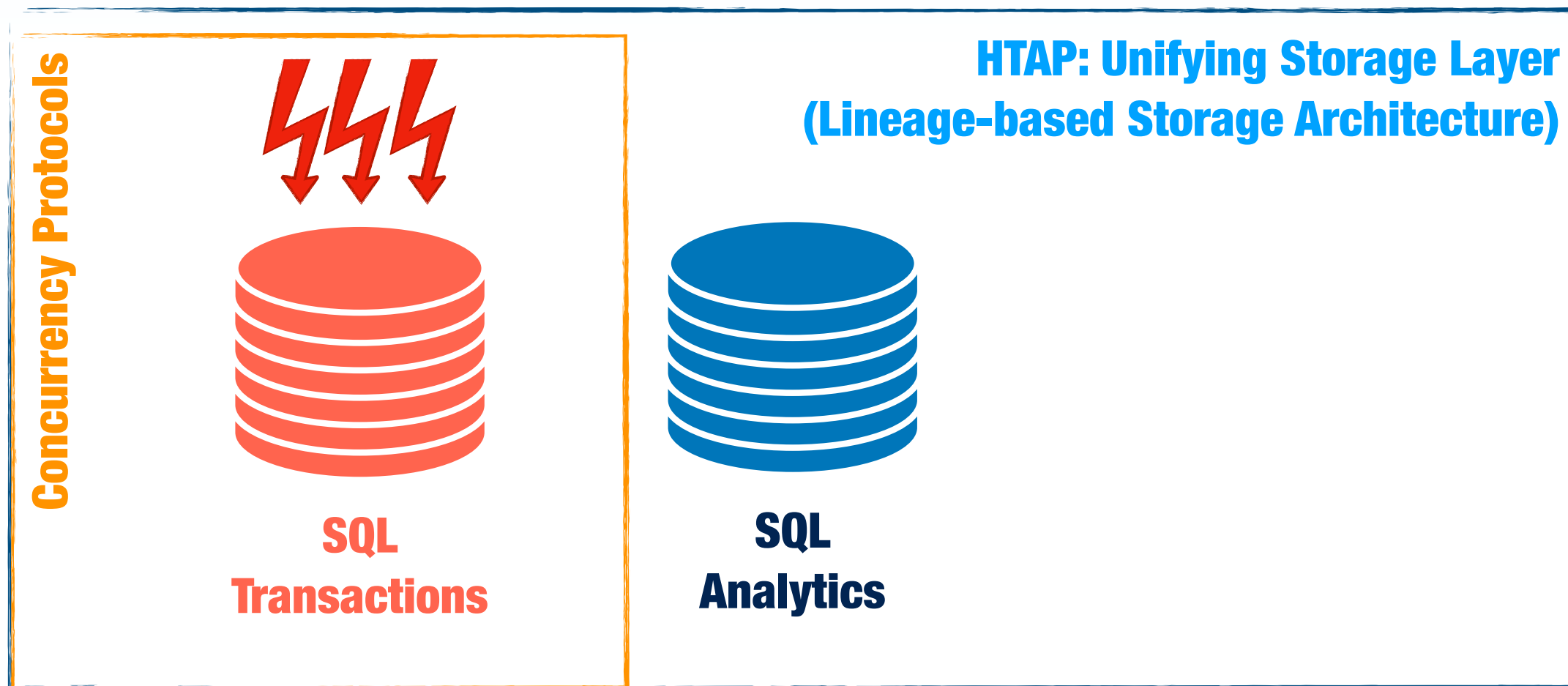


Concurrency Protocols



Concurrency Control Protocols: (e.g., 2VCC, QueCC - Best Paper Award)
[VLDB'13, VLDB'14, VLDBJ'16, Middleware'16, TDKE'15, SIGMOD'15, ICDE'16, Middleware'18]

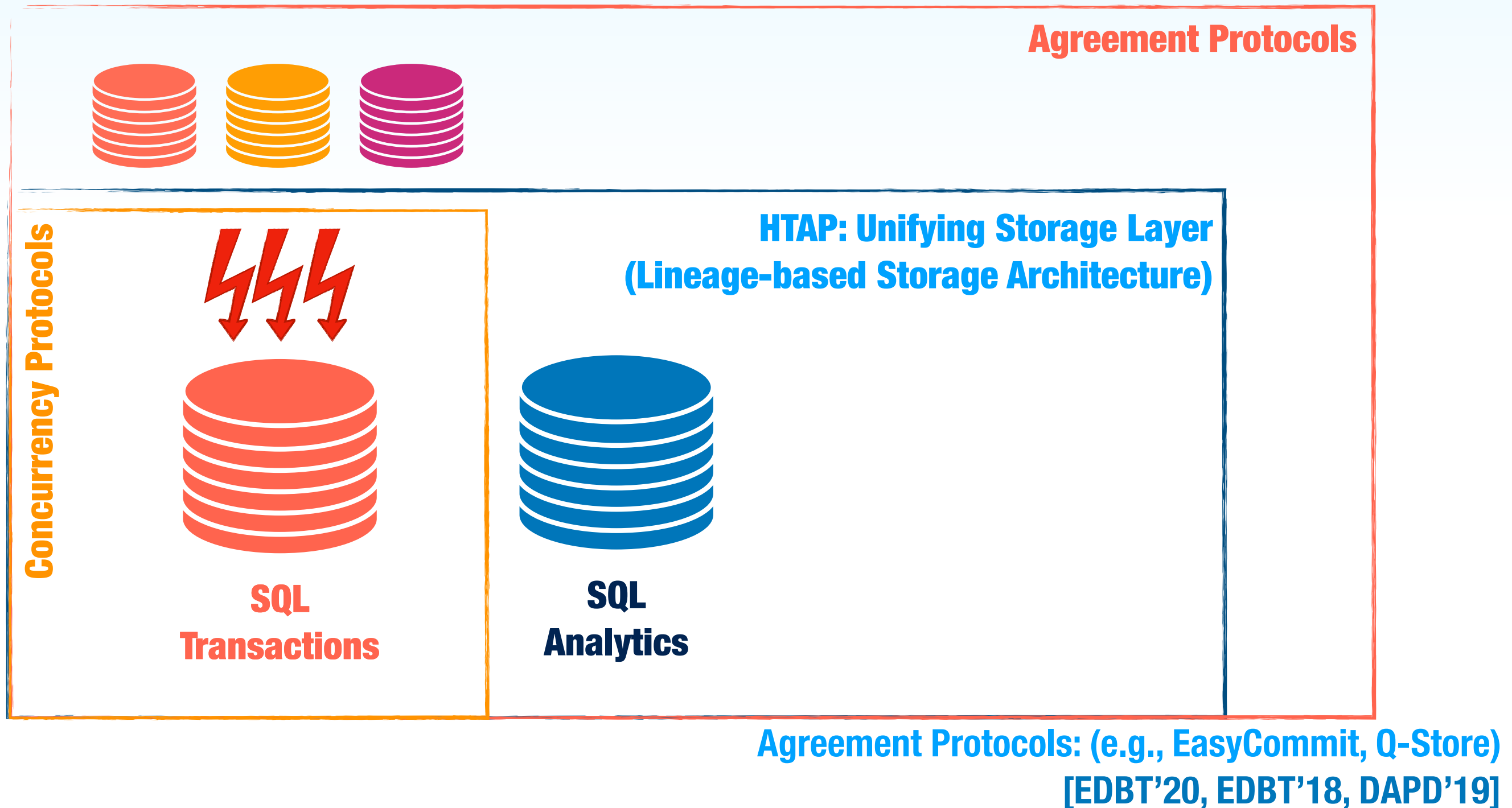
Resilient Journey...



HTAP Column-store: L-Store (Lineage-based Data Store)
[VLDB'12, ICDE'14, ICDCS'16, EDBT'18, TKDE'20 (2x) 34 filed US patents]

Graphs on SQL: (e.g., GRFusion) [SIGMOD'18, EDBT'18] 7

Resilient Journey...



Resilient Journey...

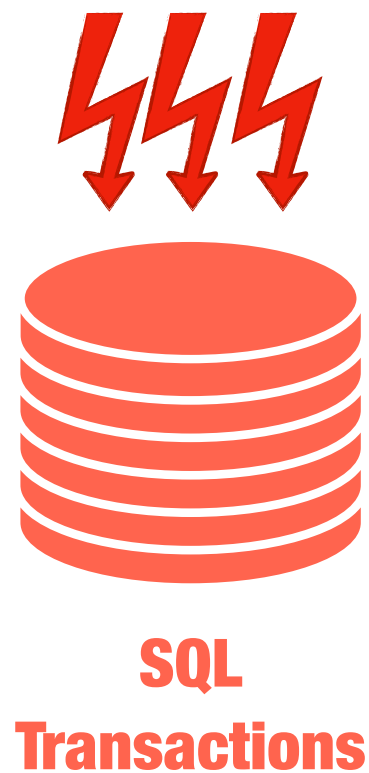
Resilient Consensus Protocols



Agreement Protocols

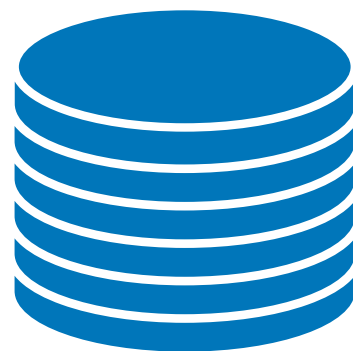


Concurrency Protocols



**SQL
Transactions**

HTAP: Unifying Storage Layer (Lineage-based Storage Architecture)



**SQL
Analytics**

Consensus Protocols: (e.g., GeoBFT, PoE, RCC, ByShard, RingBFT, Delayed Replication, CSP, Blockplane)
[VLDB'21, ICDE'21, EDTB'21, VLDB'20, ICDCS'20, ICDT'20, DISC'19 (2x), SC'19, ICDE'19, arXiv'19 (8x)]

Layer 1 (e.g., Proof-of-Work)

Layer 2 (e.g., PBFT, Po*)

Chain Management (off-chain, on-chain)

Database Stack

Query Optimization & Evaluation

Concurrency Control Protocols

Relational Operators

Files and Access Methods

Buffer Management

Disk Space Management

Storage



Log



Analytics
(Read-only)

Resilient Replication

Sharding (Isolation Semantics, Consistency Levels)

Cross-chain Network

Global Distribution

Reconfigurable Network

Recovery (View-change)

Identity Management

Applications: DeFi, Smart Contracts, IoT, Serverless

Waif-free BFT [DISC'19]

Resilient Concurrent Consensus [ICDE'21]

AHL [SIGMOD'19]

Cerberus [arXiv'20]

SharPer [SIGMOD'21]

Cluster Sending Primitive [DISC'19]

Delayed Replication [ICDT'20]

Proof-of-Execution [EDBT'21]

ByShard [VLDB'21]

RingBFT [arXiv'21]

Atomic Commitment [VLDB'20]

Cross-chain Deals [VLDB'20]

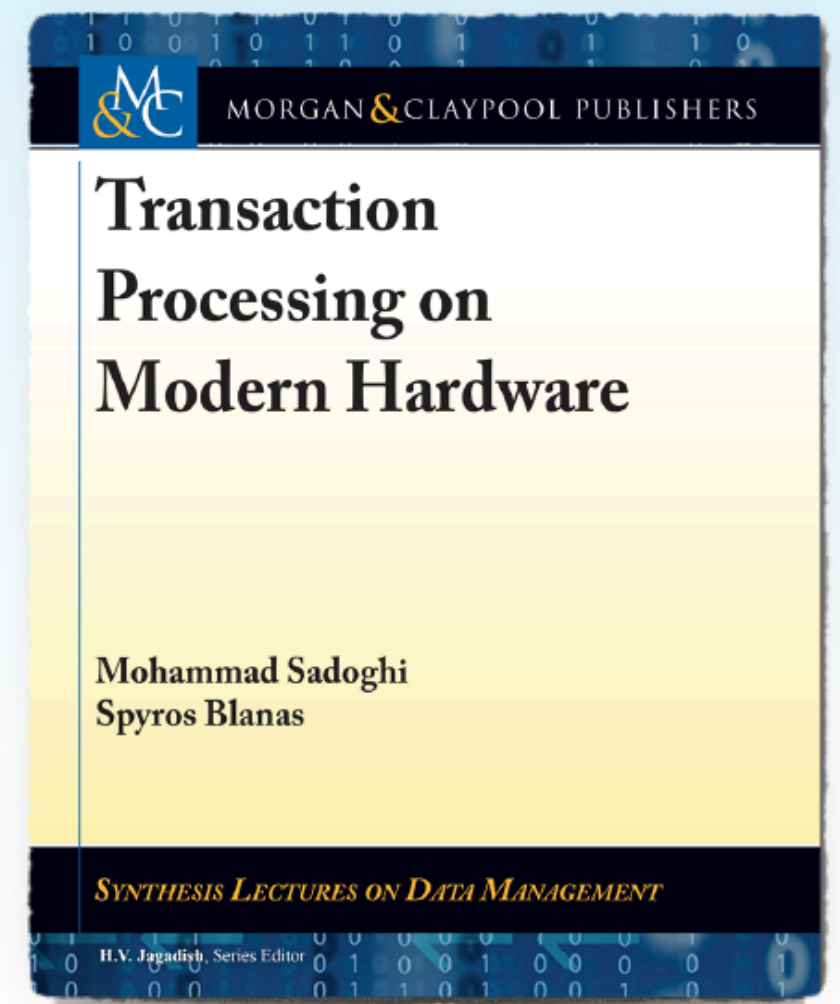
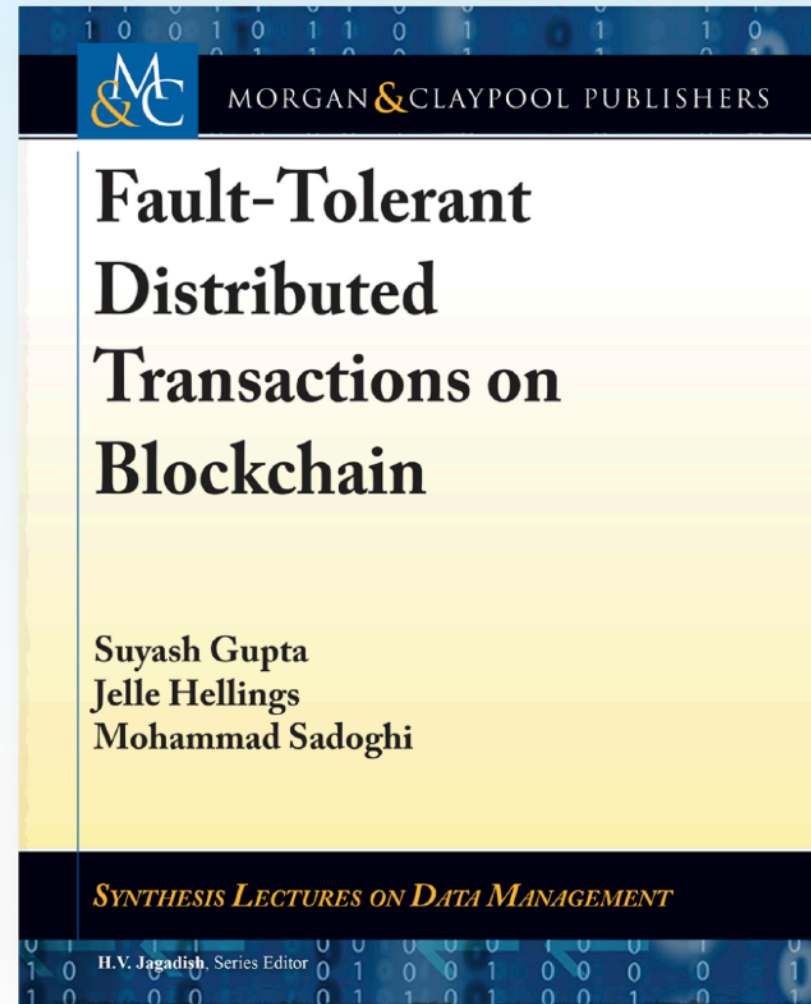
GeoBFT [VLDB'20]

Permissioned

Permissionless

BlockBench [SIGMOD'17]

Blockplane [ICDE'19]



Books

Transaction Processing on Modern Hardware.

Synthesis Lectures on Data Management, Morgan & Claypool Publishers 2019

Fault-Tolerant Distributed Transactions on Blockchain.

Synthesis Lectures on Data Management, Morgan & Claypool Publishers 2021



Press

Advancements TV With Ted Danson - CNBC, CityAM, Medium, Yahoo! Finance, Market Insider, CoinDesk, Crypto Media, Davis Enterprise, Times Union, WBOC TV/Radio

Books

Transaction Processing on Modern Hardware.

Synthesis Lectures on Data Management, Morgan & Claypool Publishers 2019

Fault-Tolerant Distributed Transactions on Blockchain.

Synthesis Lectures on Data Management, Morgan & Claypool Publishers 2021

Quantifiable Resiliency

(Graduate Student Experiments)

Aloha Lake, Desolation Wilderness
15 Miles Long
2,500 Feet Elevation Gain
(8,700 Feet at Summit)



Tomales Point Trail, Point Reyes National Seashore

9.4 Miles Long

1,579 Feet Elevation Gain



Non-Quantifiable Resiliency

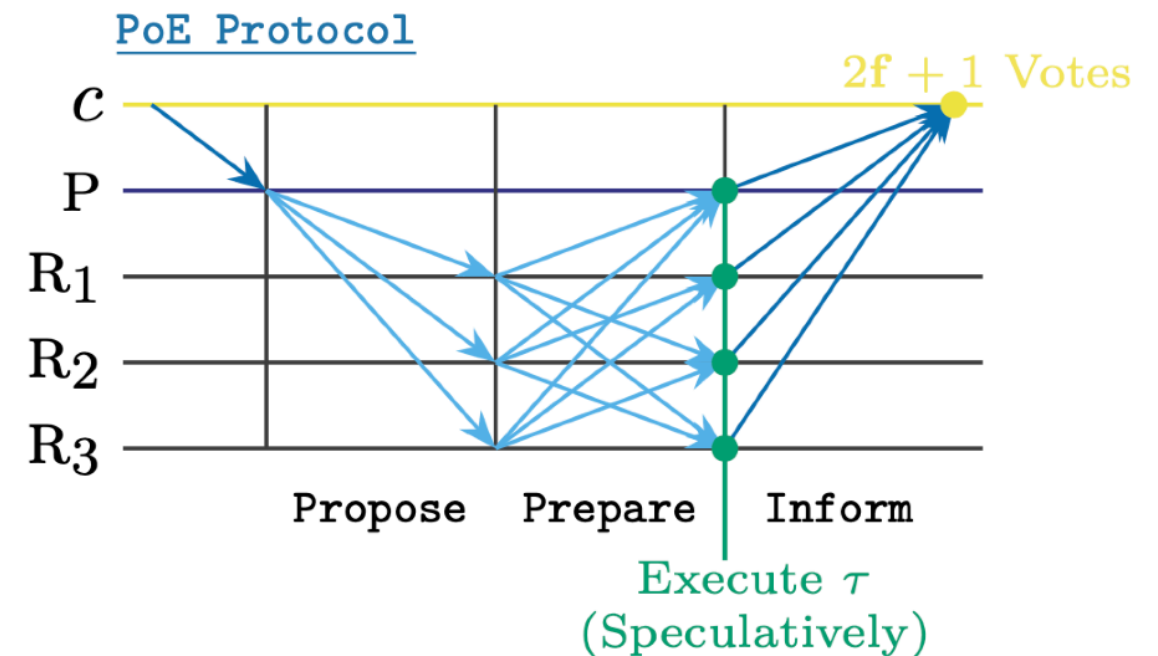
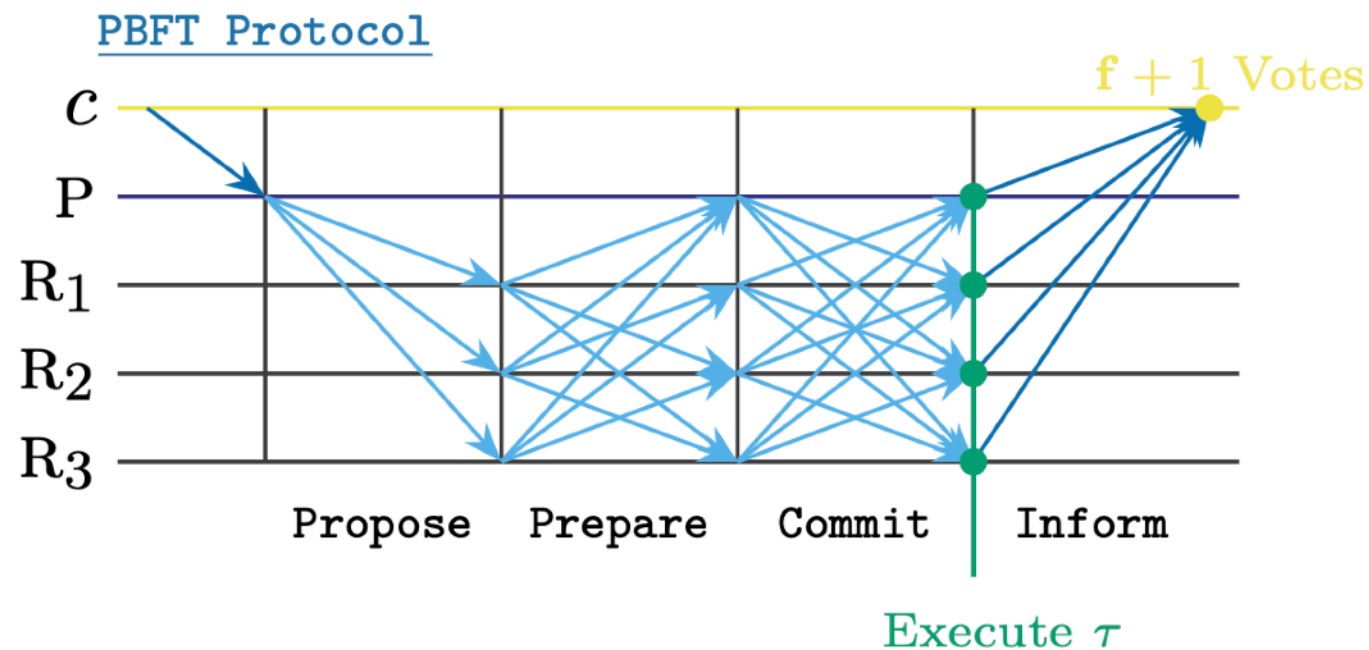
Proof-of-Execution: Reaching Consensus Through Fault-Tolerant Speculation [EDBT'21]

Out-of-Order message processing to reduce replica idleness

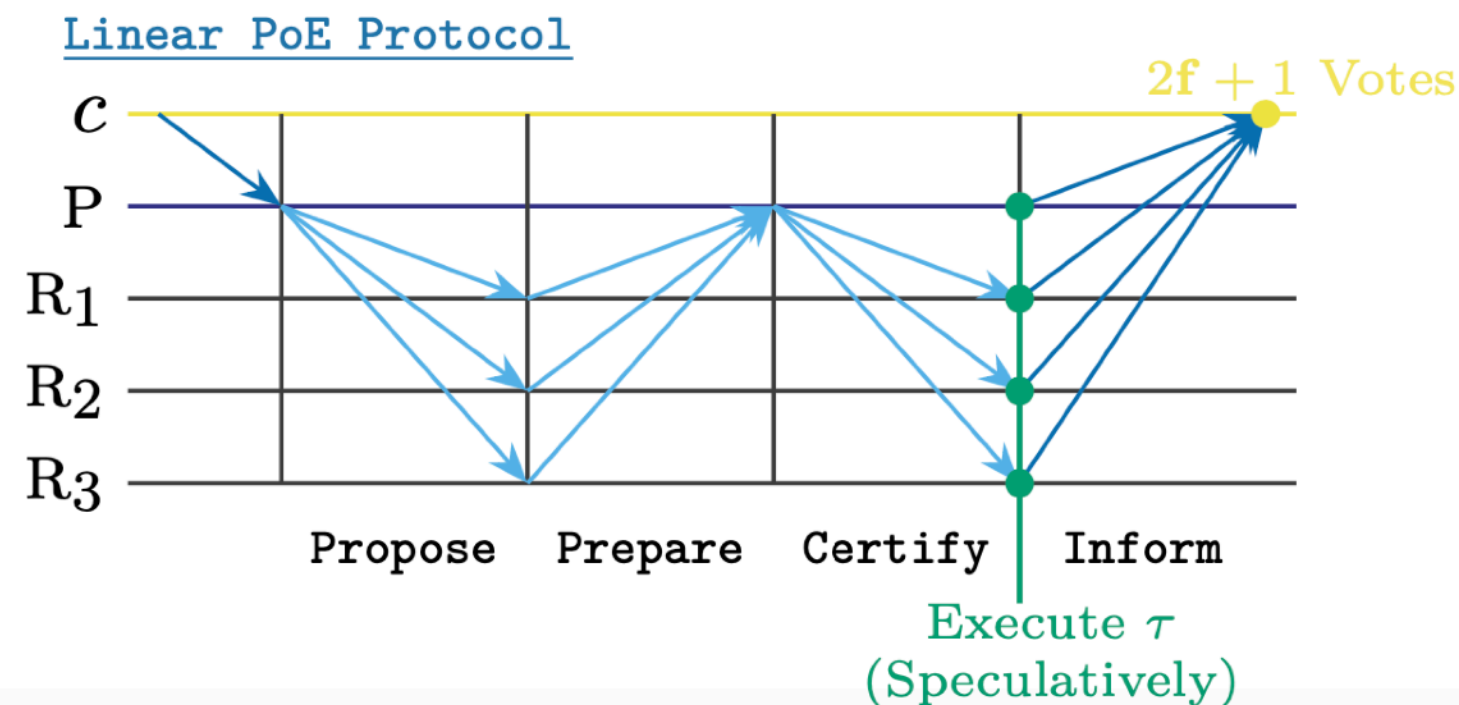
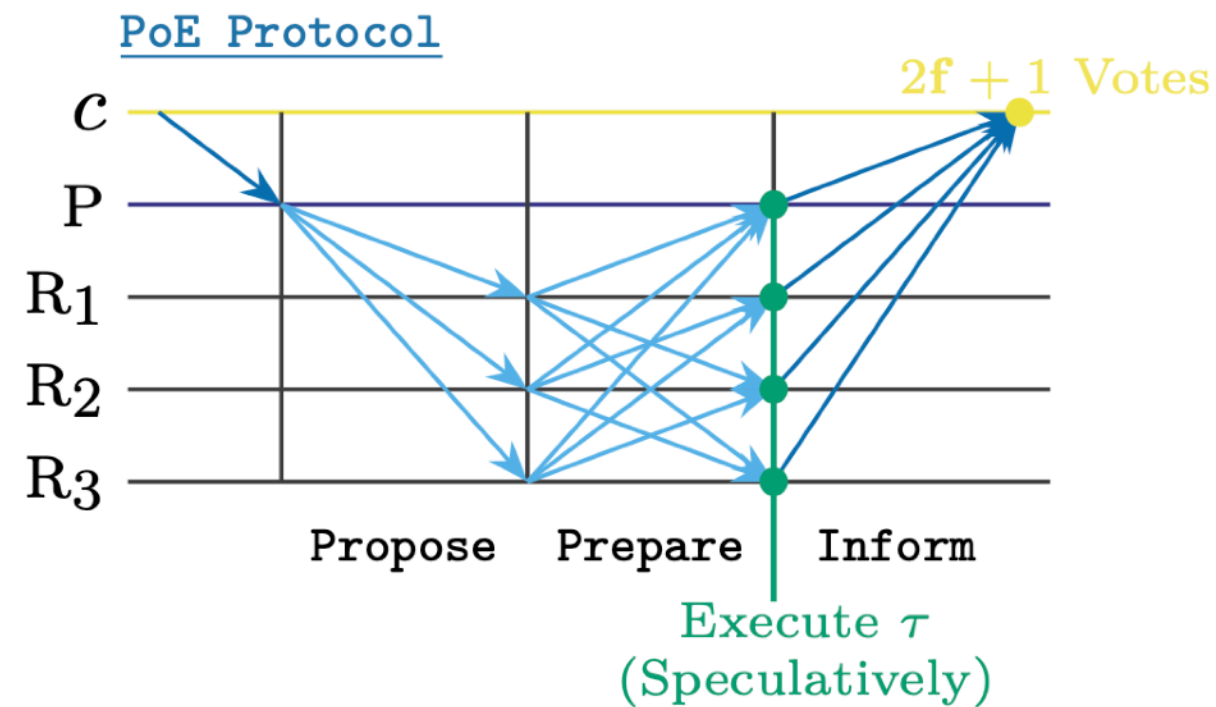
Speculative Execution with revertible/divergent replicas &

eager/irrevertible client commit

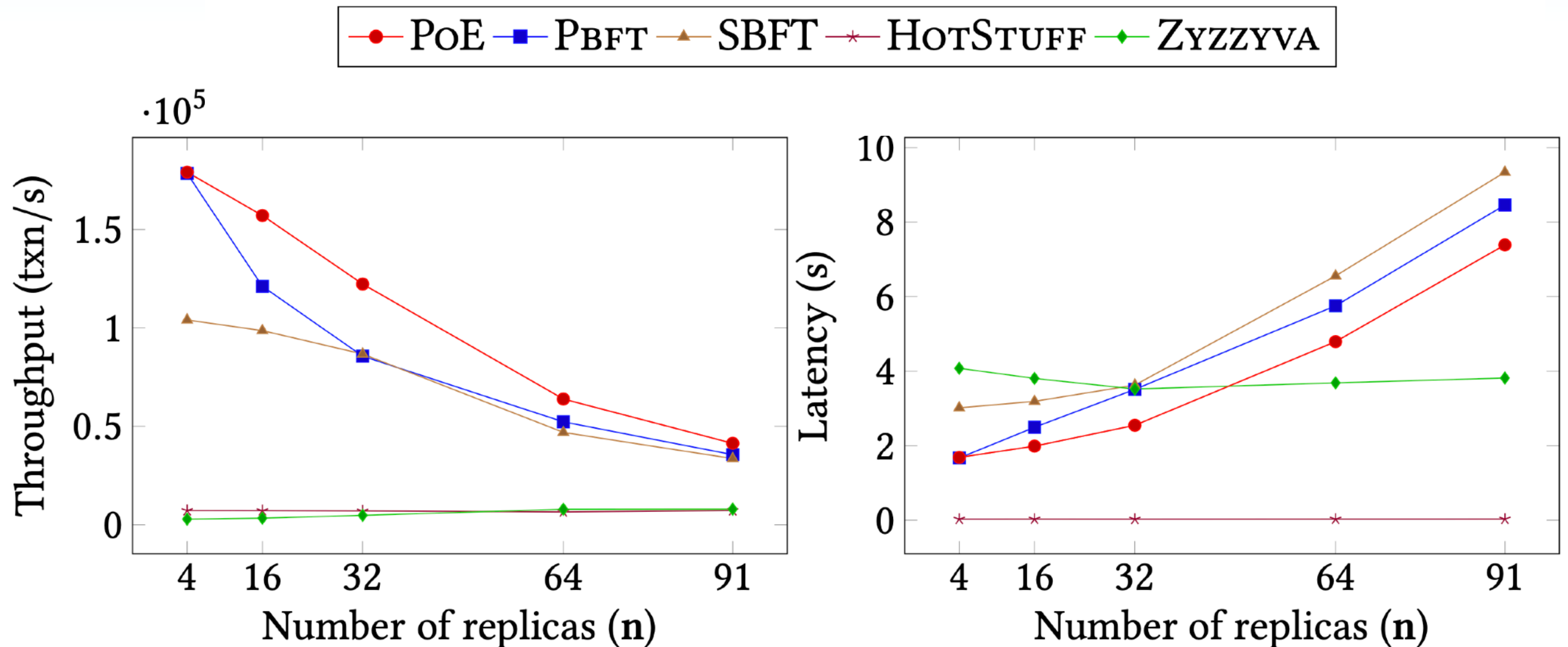
introducing linear message complexity



Proof-of-Execution: Reaching Consensus Through Fault-Tolerant Speculation [EDBT'21]



Proof-of-Execution: Reaching Consensus Through Fault-Tolerant Speculation [EDBT'21]



PoE scales beyond 91 replicas, in presence of failures, outperforms PBFT up to 43%

RCC: Resilient Concurrent Consensus Paradigm [ICDE'21]

A wait-free meta-protocol...

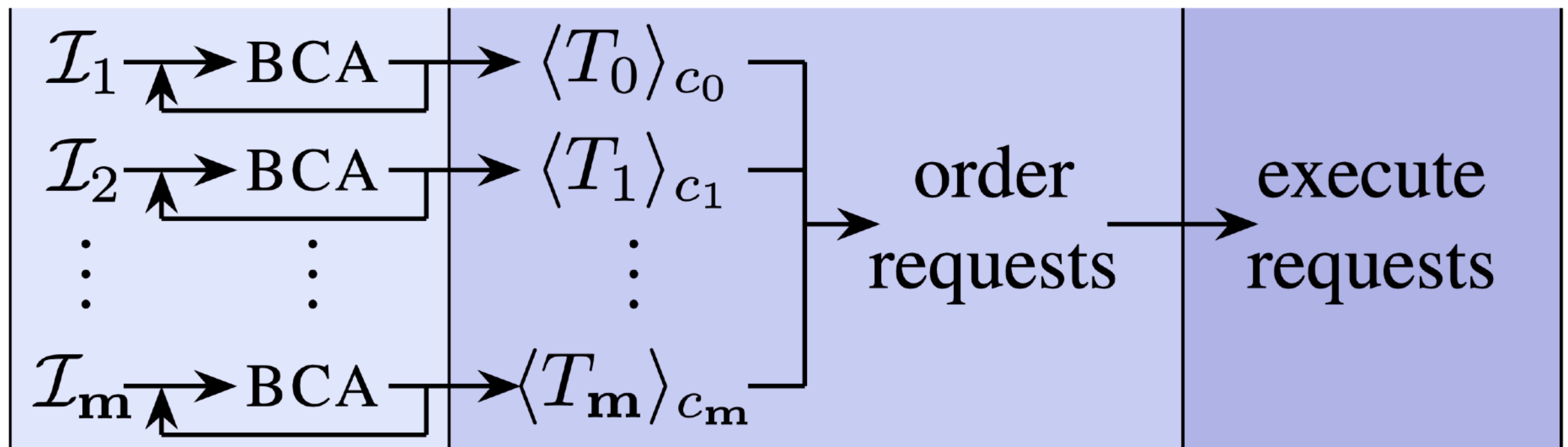
Designate multiple replicas as primaries!

Run multiple parallel consensus on each replica independently

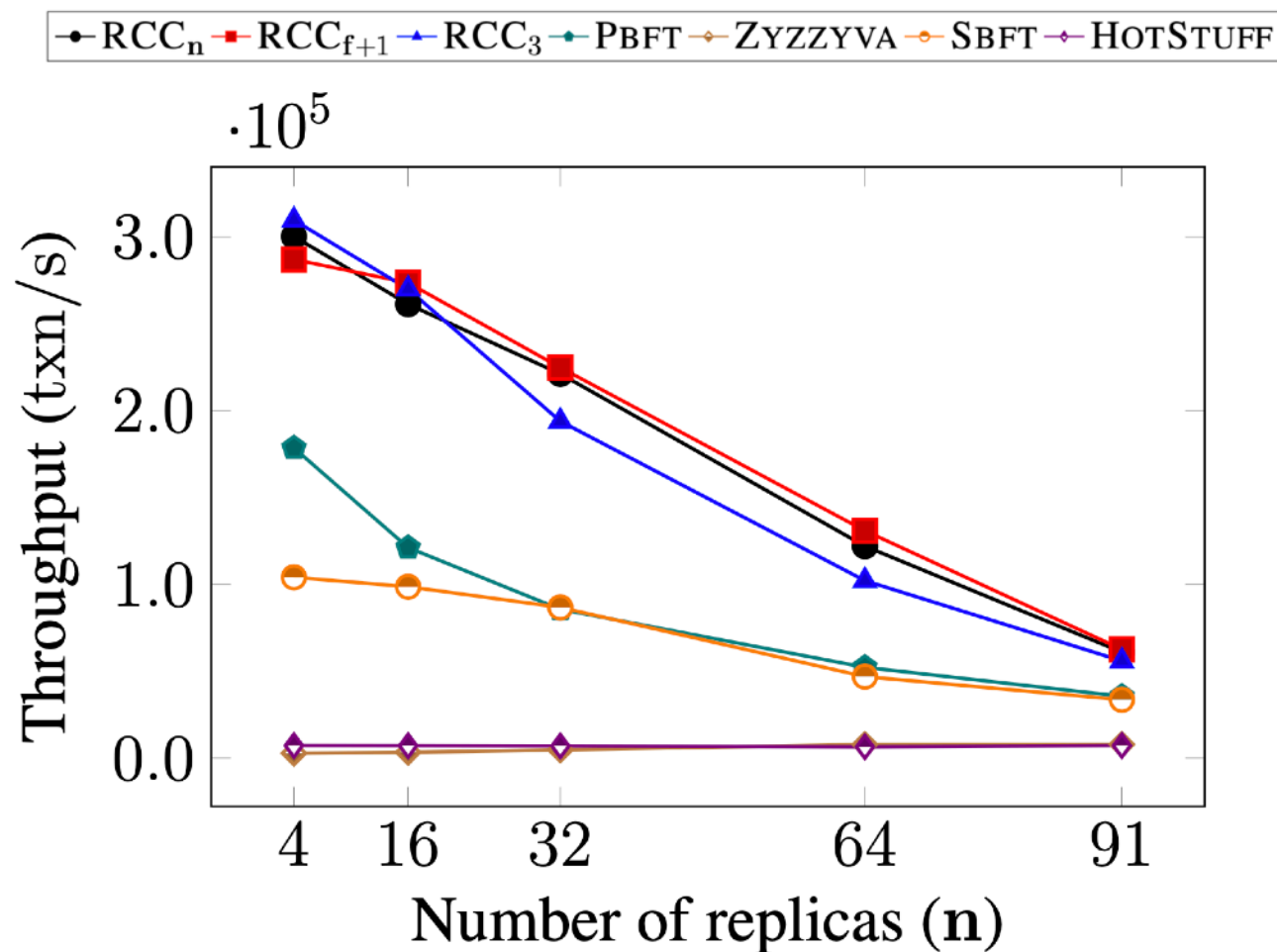
Concurrent BCA

Ordering

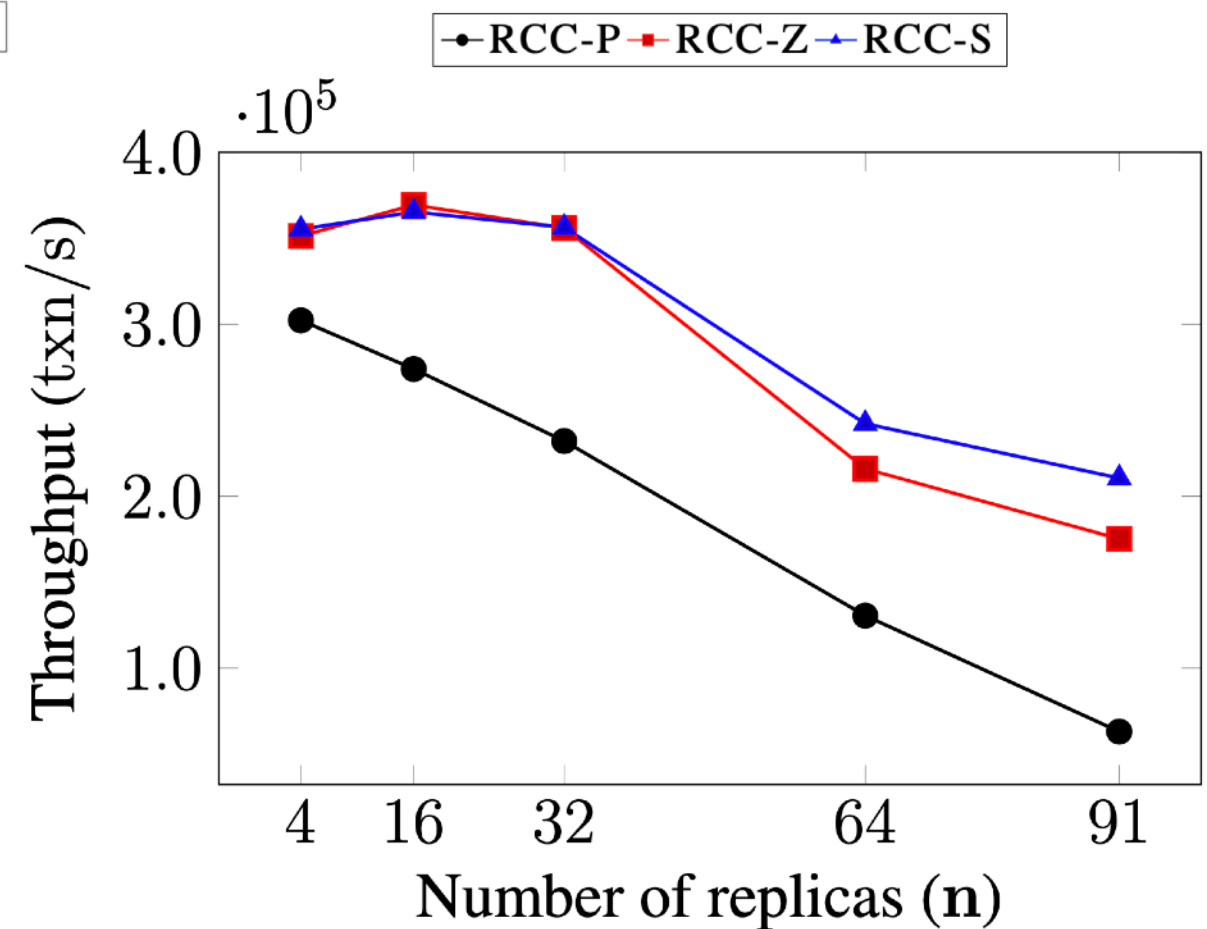
Execution



RCC: Resilient Concurrent Consensus Paradigm [ICDE'21]



Throughput up to 300,000 txns/s
(with failures)



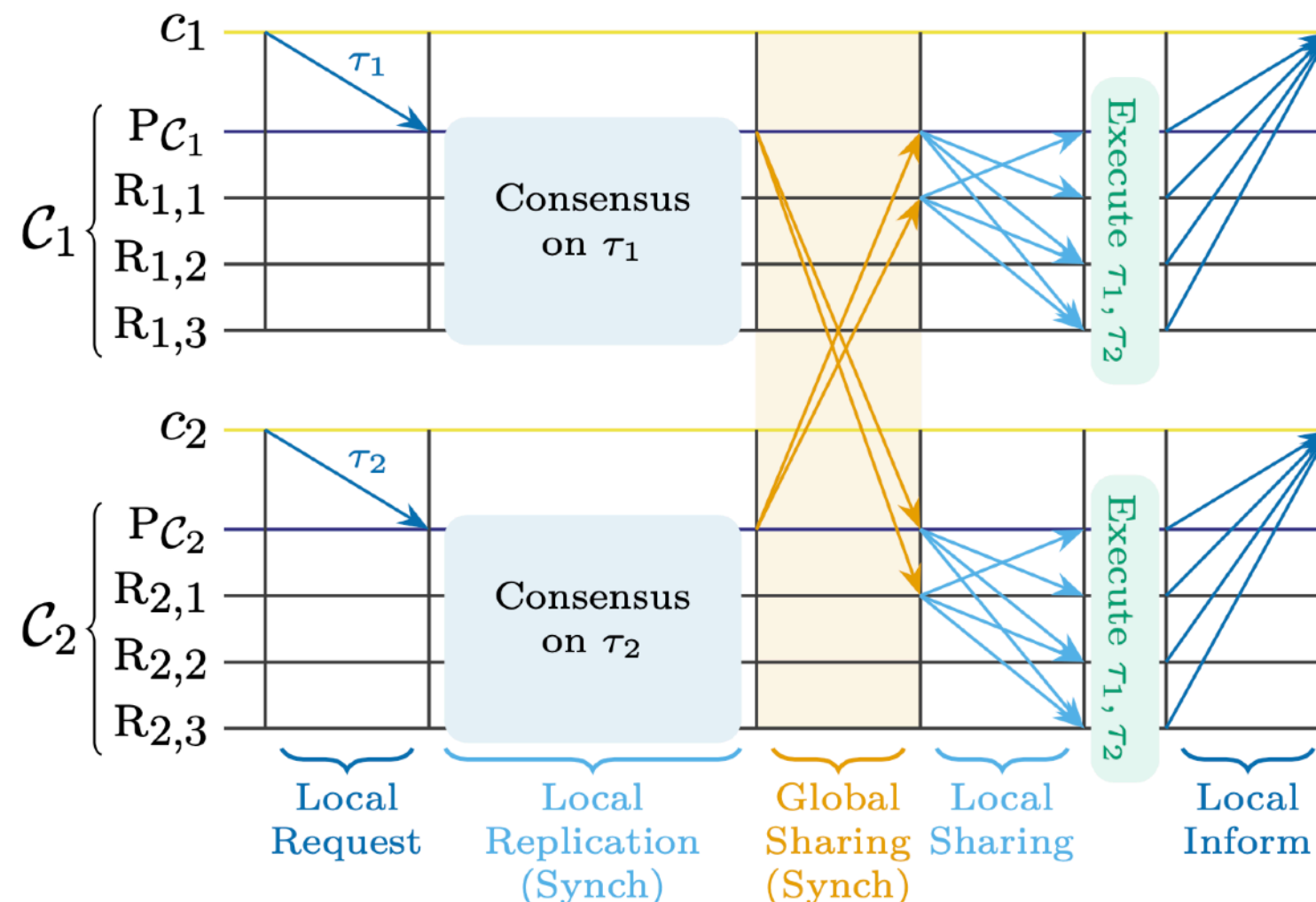
Throughput up to 400,000 txns/s
(without failures)

GeoBFT: Global Scale Resilient Consensus [VLDB'20]

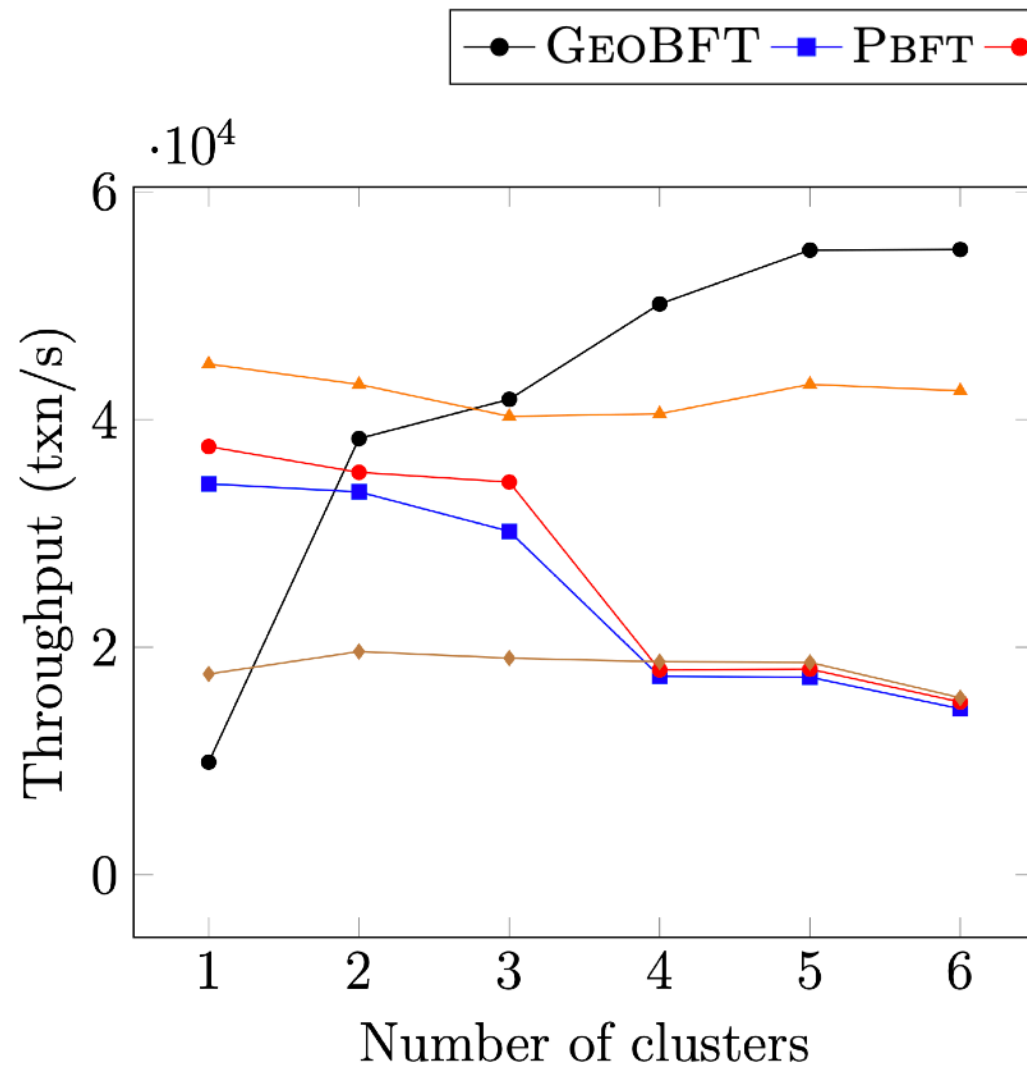
A meta-protocol, locally running any BFT in parallel and independently

Global ordering provably requires only linear communication

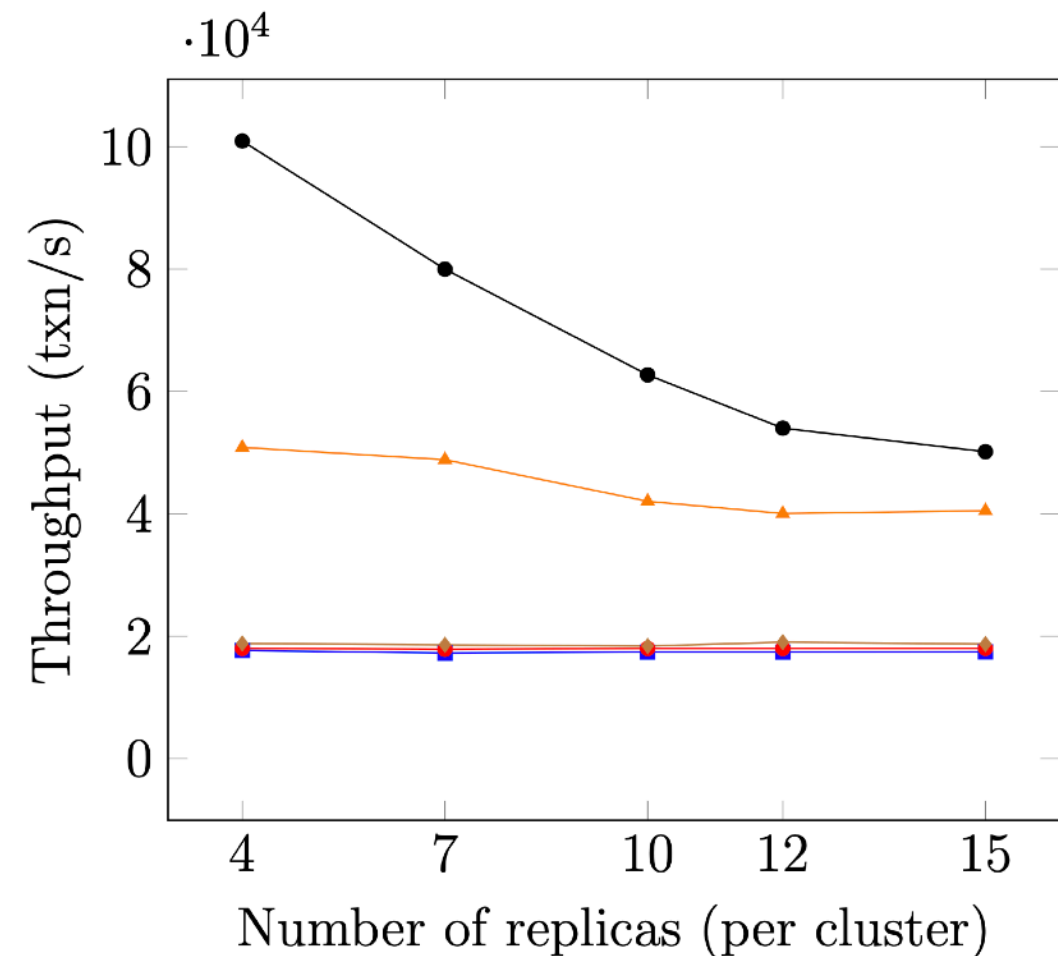
Provably sufficient for primary to send a certificate to at most $f+1$ replicas,
malicious primary is detectable and replaceable



GeoBFT: Global Scale Resilient Consensus [VLDB'20]



GeoBFT easily scales across 6 countries in 4 continents due to GeoBFT protocol.

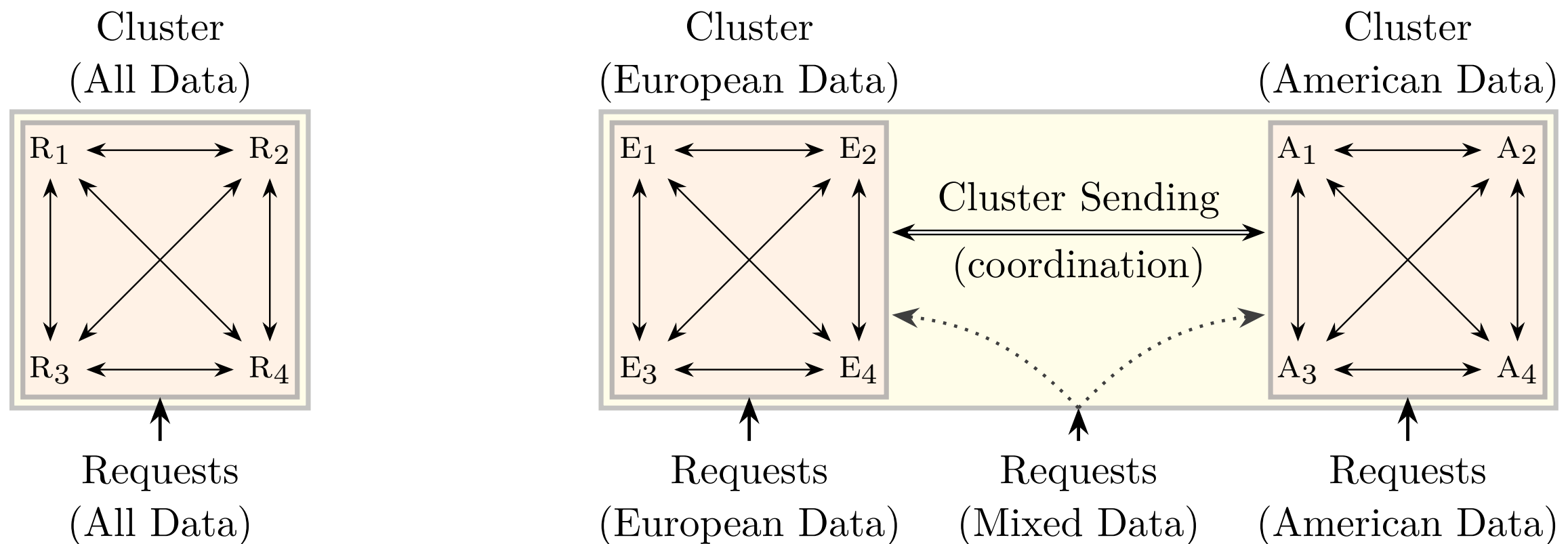


GeoBFT scales a permissioned blockchain up to 60 replicas globally.

The Fault-Tolerant Cluster-Sending Problem [DISC'19]

formalizing the problem of sending a message from one Byzantine cluster to another Byzantine cluster in a reliable manner,

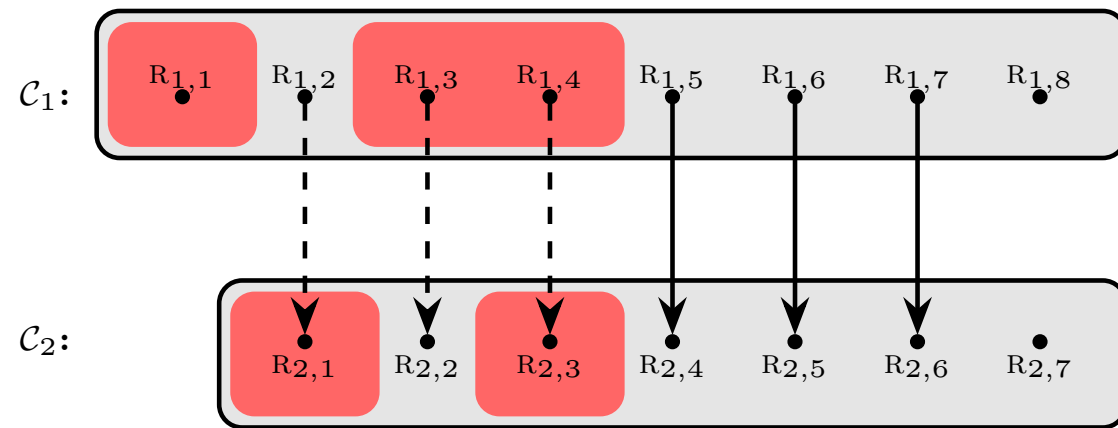
establishing lower bounds on the complexity of this problem under crash failures and Byzantine failures (linear in the size of clusters)



The Fault-Tolerant Cluster-Sending Problem [DISC'19]

formalizing the problem of sending a message from one Byzantine cluster to another Byzantine cluster in a reliable manner,

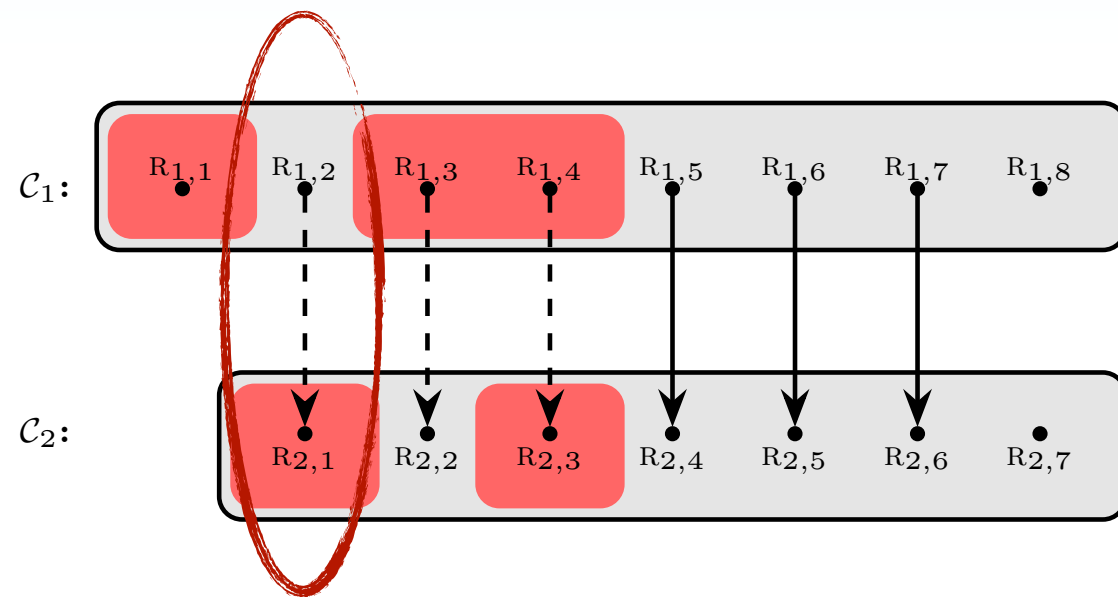
establishing lower bounds on the complexity of this problem under crash failures and Byzantine failures
(linear in the size of clusters)



	Protocol	System	Robustness	Messages	Message size
non-linear	RB-bcs	Omit	$nc_1 > 2fc_1, nc_2 > fc_2$	$(fc_1 + 1) \cdot (fc_2 + 1)$	$\mathcal{O}(\ v\)$
	RB-brs	Byzantine, RS	$nc_1 > 2fc_1, nc_2 > fc_2$	$(2fc_1 + 1) \cdot (fc_2 + 1)$	$\mathcal{O}(\ v\)$
	RB-bcs	Byzantine, RS	$nc_1 > 2fc_1, nc_2 > fc_2$	$(fc_1 + 1) \cdot (fc_2 + 1)$	$\mathcal{O}(\ v\ + fc_1)$
	RB-bcs	Byzantine, CS	$nc_1 > 2fc_1, nc_2 > fc_2$	$(fc_1 + 1) \cdot (fc_2 + 1)$	$\mathcal{O}(\ v\)$
linear	PBS-bcs	Omit	$nc_1 > 3fc_1, nc_2 > 3fc_2$	$\mathcal{O}(\max(nc_1, nc_2))$ (optimal)	$\mathcal{O}(\ v\)$
	PBS-brs	Byzantine, RS	$nc_1 > 4fc_1, nc_2 > 4fc_2$	$\mathcal{O}(\max(nc_1, nc_2))$ (optimal)	$\mathcal{O}(\ v\)$
	PBS-bcs	Byzantine, RS	$nc_1 > 3fc_1, nc_2 > 3fc_2$	$\mathcal{O}(\max(nc_1, nc_2))$	$\mathcal{O}(\ v\ + fc_1)$
	PBS-bcs	Byzantine, CS	$nc_1 > 3fc_1, nc_2 > 3fc_2$	$\mathcal{O}(\max(nc_1, nc_2))$ (optimal)	$\mathcal{O}(\ v\)$

Byzantine Cluster-Sending in Expected Constant Communication [arXiv'21]

formalizing the problem of probabilistically sending a message from one Byzantine cluster to another Byzantine cluster in a reliable manner,
 establishing lower bounds on the complexity of this problem under crash failures and Byzantine failures
 (expected constant message complexity)



	Protocol	Robustness	Message Steps		O. U.	
			(expected)	(worst)		
	PBS-CS [13]	$\min(\mathbf{n}_{C_1}, \mathbf{n}_{C_2}) > \mathbf{f}_{C_1} + \mathbf{f}_{C_2}$	$\mathbf{f}_{C_1} + \mathbf{f}_{C_2} + 1$		✓	✗
	PBS-CS [13]	$\mathbf{n}_{C_1} > 3\mathbf{f}_{C_1}, \mathbf{n}_{C_2} > 3\mathbf{f}_{C_2}$	$\max(\mathbf{n}_{C_1}, \mathbf{n}_{C_2})$		✓	✗
	GEOBFT [12]	$\mathbf{n}_{C_1} = \mathbf{n}_{C_2} > 3 \max(\mathbf{f}_{C_1}, \mathbf{f}_{C_2})$	$\mathbf{f}_{C_2} + 1^{\ddagger}$	$\Omega(\mathbf{f}_{C_1} \mathbf{n}_{C_2})$	✗	✓
This Paper	PPCS	$\mathbf{n}_{C_1} > 2\mathbf{f}_{C_1}, \mathbf{n}_{C_2} > 2\mathbf{f}_{C_2}$	4	$(\mathbf{f}_{C_1} + 1)(\mathbf{f}_{C_2} + 1)$	✗	✓
	PPCS	$\mathbf{n}_{C_1} > 3\mathbf{f}_{C_1}, \mathbf{n}_{C_2} > 3\mathbf{f}_{C_2}$	$2\frac{1}{4}$	$(\mathbf{f}_{C_1} + 1)(\mathbf{f}_{C_2} + 1)$	✗	✓
	PLCS	$\min(\mathbf{n}_{C_1}, \mathbf{n}_{C_2}) > \mathbf{f}_{C_1} + \mathbf{f}_{C_2}$	4	$\mathbf{f}_{C_1} + \mathbf{f}_{C_2} + 1$	✓	✓
	PLCS	$\min(\mathbf{n}_{C_1}, \mathbf{n}_{C_2}) > 2(\mathbf{f}_{C_1} + \mathbf{f}_{C_2})$	$2\frac{1}{4}$	$\mathbf{f}_{C_1} + \mathbf{f}_{C_2} + 1$	✓	✓
	PLCS	$\mathbf{n}_{C_1} > 3\mathbf{f}_{C_1}, \mathbf{n}_{C_2} > 3\mathbf{f}_{C_2}$	3	$\max(\mathbf{n}_{C_1}, \mathbf{n}_{C_2})$	✓	✓

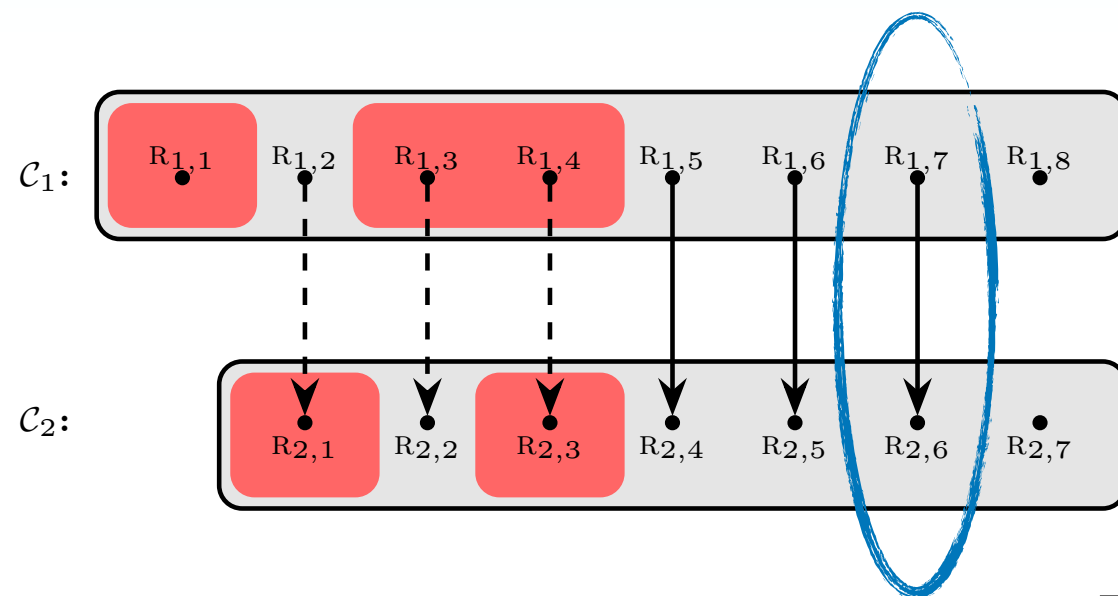
Byzantine Cluster-Sending in Expected Constant Communication [arXiv'21]

formalizing the problem of probabilistically sending a message from one

Byzantine cluster to another Byzantine cluster in a reliable manner,

establishing lower bounds on the complexity of this problem under crash failures and Byzantine failures

(expected constant message complexity)



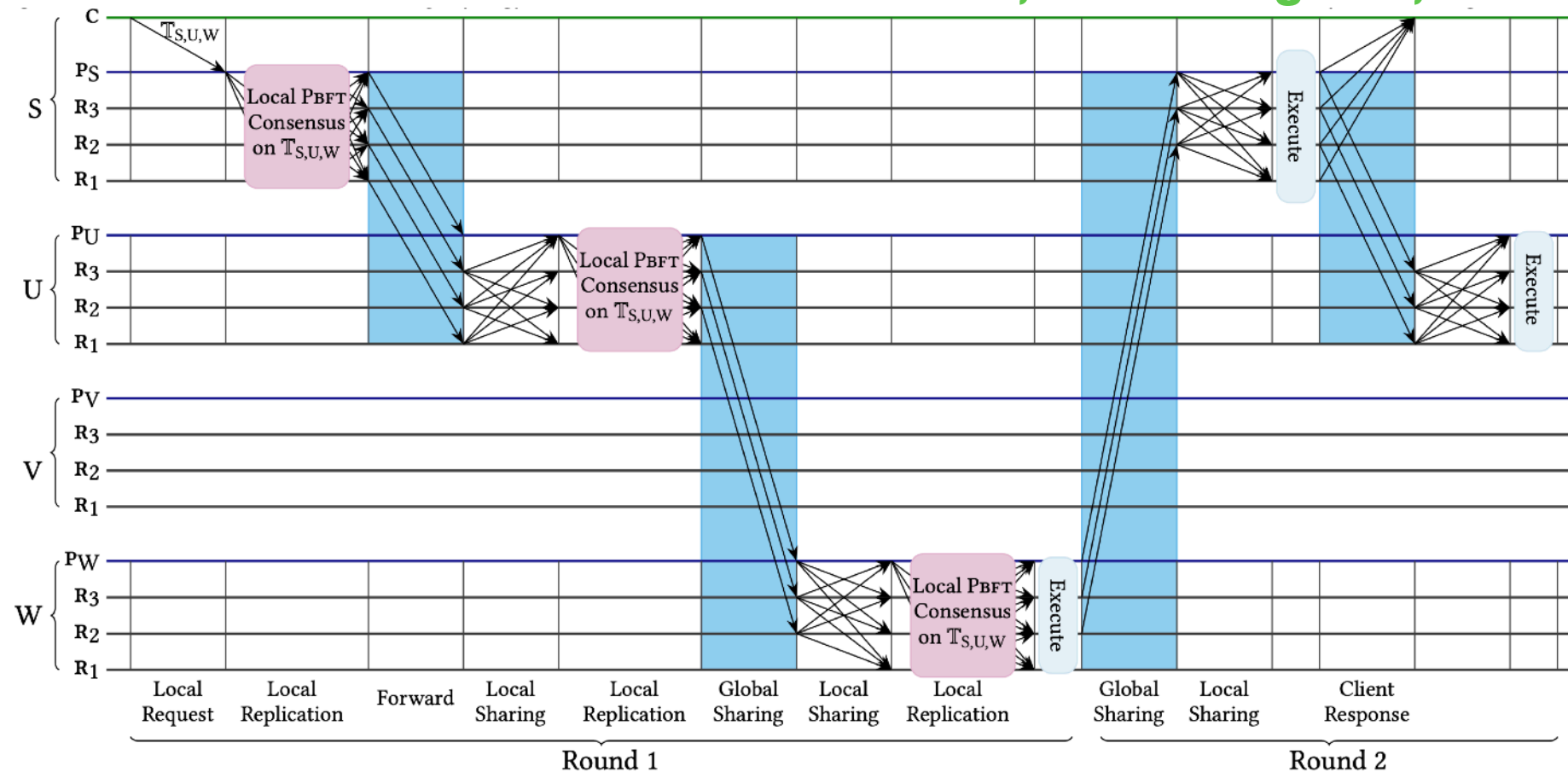
	Protocol	Robustness	Message Steps		O. U.	
			(expected)	(worst)		
	PBS-CS [13]	$\min(\mathbf{n}_{C_1}, \mathbf{n}_{C_2}) > \mathbf{f}_{C_1} + \mathbf{f}_{C_2}$	$\mathbf{f}_{C_1} + \mathbf{f}_{C_2} + 1$		✓	✗
	PBS-CS [13]	$\mathbf{n}_{C_1} > 3\mathbf{f}_{C_1}, \mathbf{n}_{C_2} > 3\mathbf{f}_{C_2}$	$\max(\mathbf{n}_{C_1}, \mathbf{n}_{C_2})$		✓	✗
	GEOBFT [12]	$\mathbf{n}_{C_1} = \mathbf{n}_{C_2} > 3 \max(\mathbf{f}_{C_1}, \mathbf{f}_{C_2})$	$\mathbf{f}_{C_2} + 1^{\ddagger}$	$\Omega(\mathbf{f}_{C_1} \mathbf{n}_{C_2})$	✗	✓
This Paper	PPCS	$\mathbf{n}_{C_1} > 2\mathbf{f}_{C_1}, \mathbf{n}_{C_2} > 2\mathbf{f}_{C_2}$	4	$(\mathbf{f}_{C_1} + 1)(\mathbf{f}_{C_2} + 1)$	✗	✓
	PPCS	$\mathbf{n}_{C_1} > 3\mathbf{f}_{C_1}, \mathbf{n}_{C_2} > 3\mathbf{f}_{C_2}$	$2\frac{1}{4}$	$(\mathbf{f}_{C_1} + 1)(\mathbf{f}_{C_2} + 1)$	✗	✓
	PLCS	$\min(\mathbf{n}_{C_1}, \mathbf{n}_{C_2}) > \mathbf{f}_{C_1} + \mathbf{f}_{C_2}$	4	$\mathbf{f}_{C_1} + \mathbf{f}_{C_2} + 1$	✓	✓
	PLCS	$\min(\mathbf{n}_{C_1}, \mathbf{n}_{C_2}) > 2(\mathbf{f}_{C_1} + \mathbf{f}_{C_2})$	$2\frac{1}{4}$	$\mathbf{f}_{C_1} + \mathbf{f}_{C_2} + 1$	✓	✓
	PLCS	$\mathbf{n}_{C_1} > 3\mathbf{f}_{C_1}, \mathbf{n}_{C_2} > 3\mathbf{f}_{C_2}$	3	$\max(\mathbf{n}_{C_1}, \mathbf{n}_{C_2})$	✓	✓

RingBFT: Resilient Consensus Over Sharded Ring Topology [arXiv'21]

A meta-protocol adhering to the ring order, and follow the principle of
process, forward, and re-transmit

Guarantees consensus for each cross-shard transaction in
at most two rotations around the ring

Sustaining over 1,200,000 transactions per second when deployed globally spanning ten
countries, fifteen regions, nearly 500 replicas.



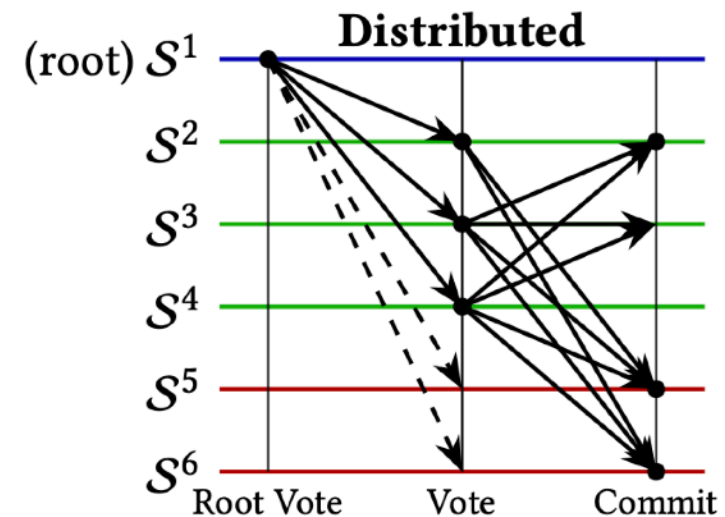
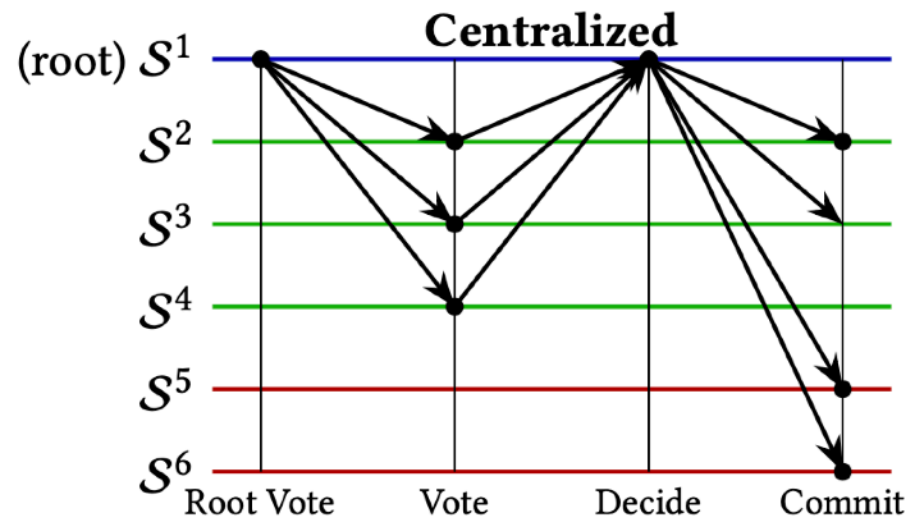
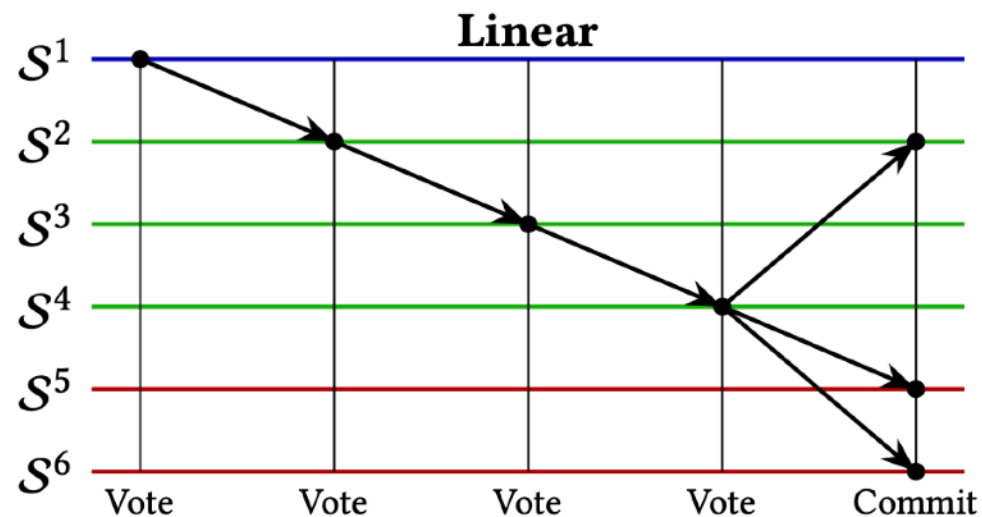
ByShard: Sharding in a Byzantine Environment [VLDB'21]

Processing multi-shard transaction via the orchestrate-execute model

Processing is broken down into three types of shard-steps: **vote**, **commit**, and **abort**

Each shard-step is performed **via one consensus step**

Steps are communicated via **cluster-sending**



Coordination-Free Byzantine Replication With Minimal Communication Costs [ICDT'20]

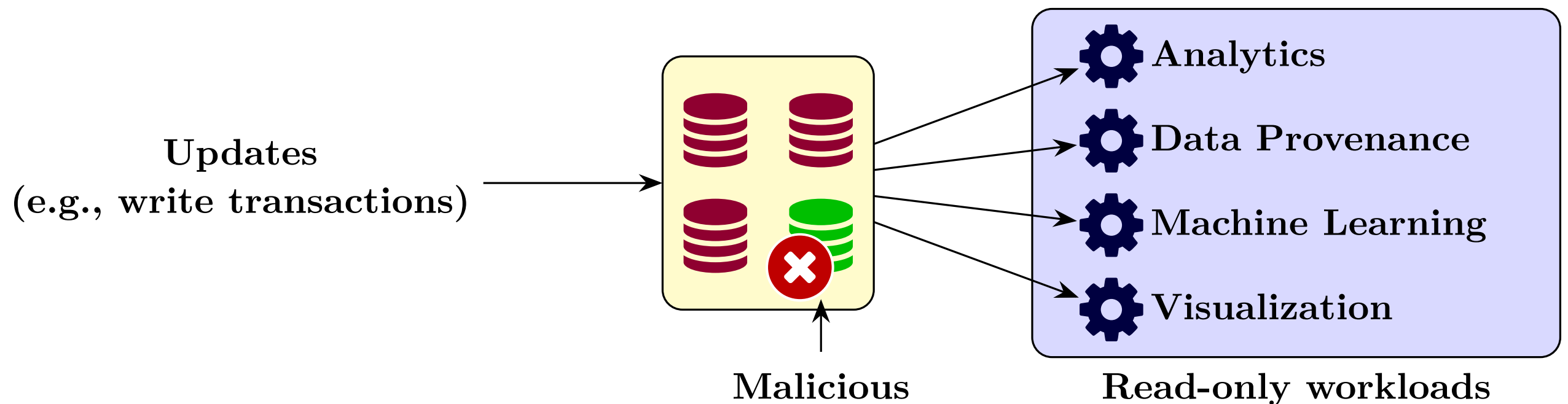
formalizing the Byzantine learner problem to support efficient

analytics for blockchain applications

introducing the delayed-replication algorithm,

utilizing information dispersal techniques,

giving rise to a coordination-free, push-based, minimal communication protocol



Coordination-Free Byzantine Replication With Minimal Communication Costs [ICDT'20]

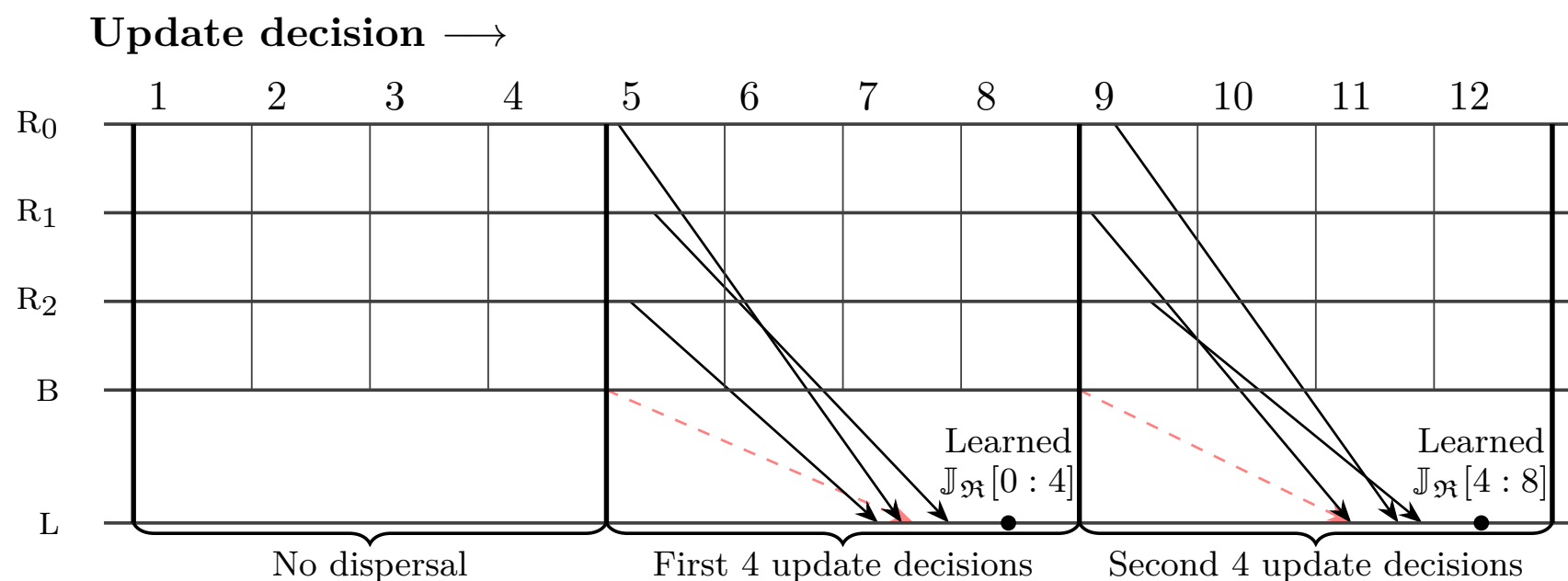
formalizing the Byzantine learner problem to support efficient

analytics for blockchain applications

introducing the delayed-replication algorithm,

utilizing information dispersal techniques,

giving rise to a coordination-free, push-based, minimal communication protocol



System	Checksum	Complexity for the learner		
		<i>Data sent per replica</i>	<i>Data received</i>	<i>Decode steps</i>
$\mathbf{b} = 0$	None	$\mathcal{O}(s/\mathbf{g})$	$\mathcal{O}(s(\mathbf{n}/\mathbf{g}))$	u/\mathbf{n}
$\mathbf{b} < \mathbf{g}$	Simple	$\mathcal{O}(s/\mathbf{g})$	$\mathcal{O}(s(\mathbf{n}/\mathbf{g}))$	$\binom{\mathbf{g}+\mathbf{b}}{\mathbf{g}}(u/\mathbf{n})$
$\mathbf{b} < \mathbf{g}$	Tree	$\mathcal{O}(s/\mathbf{g} + (u/\mathbf{n}) \log(\mathbf{n}))$	$\mathcal{O}(s(\mathbf{n}/\mathbf{g}) + u \log(\mathbf{n}))$	u/\mathbf{n}

Permissioned Blockchain Through the Looking Glass: Architectural and Implementation Lessons Learned [ICDCS'20]

Single-threaded Monolithic Design

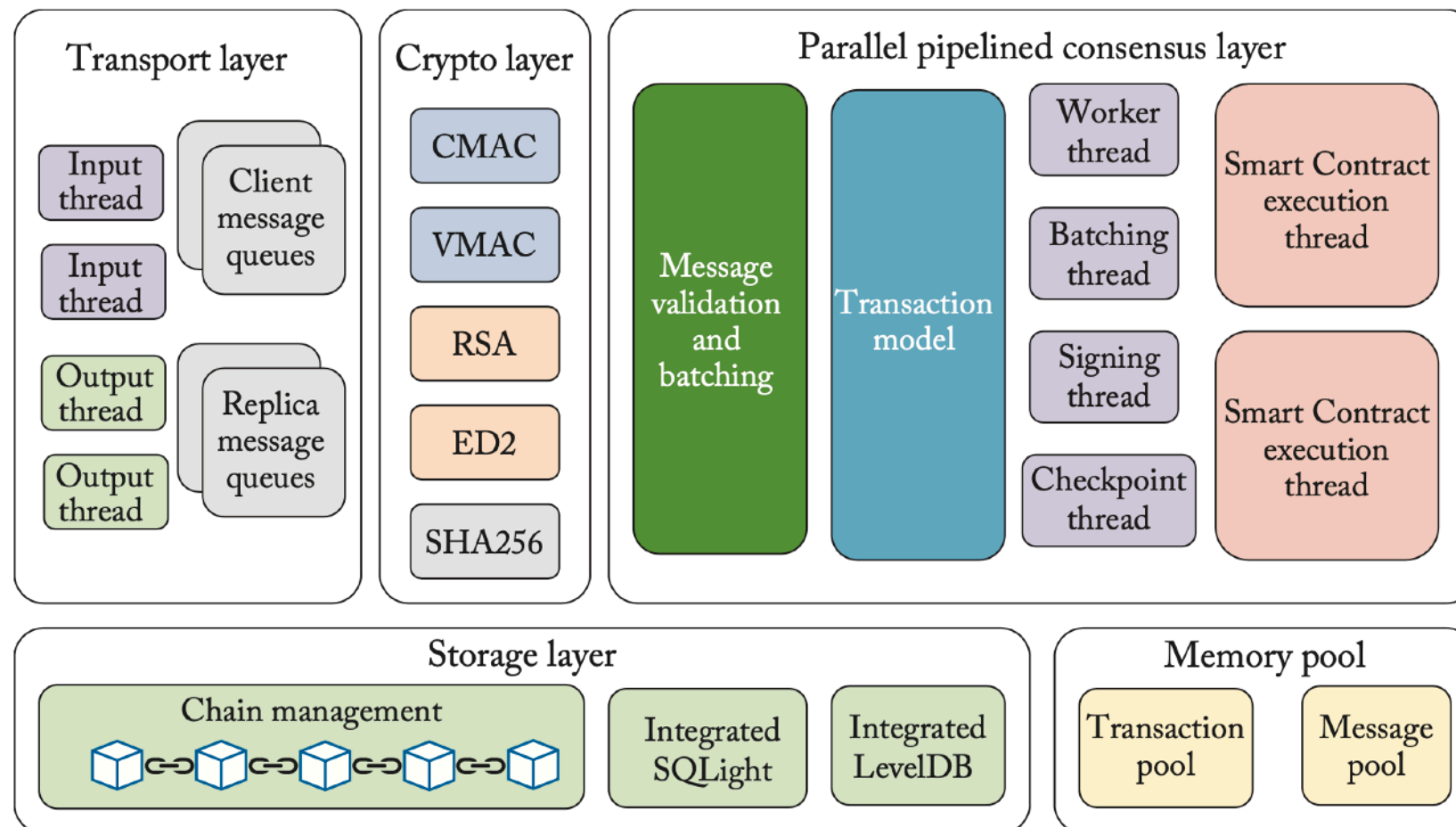
Out-of-ordering Consensus Communication

De-coupled Ordering and Execution

Off-Chain Memory Management

Expensive Cryptographic Practices (DS vs. MAC)

Smart Contracts Code Generation (Pre-compilation)



Permissioned Blockchain Through the Looking Glass: Architectural and Implementation Lessons Learned [ICDCS'20]

Single-threaded Monolithic Design

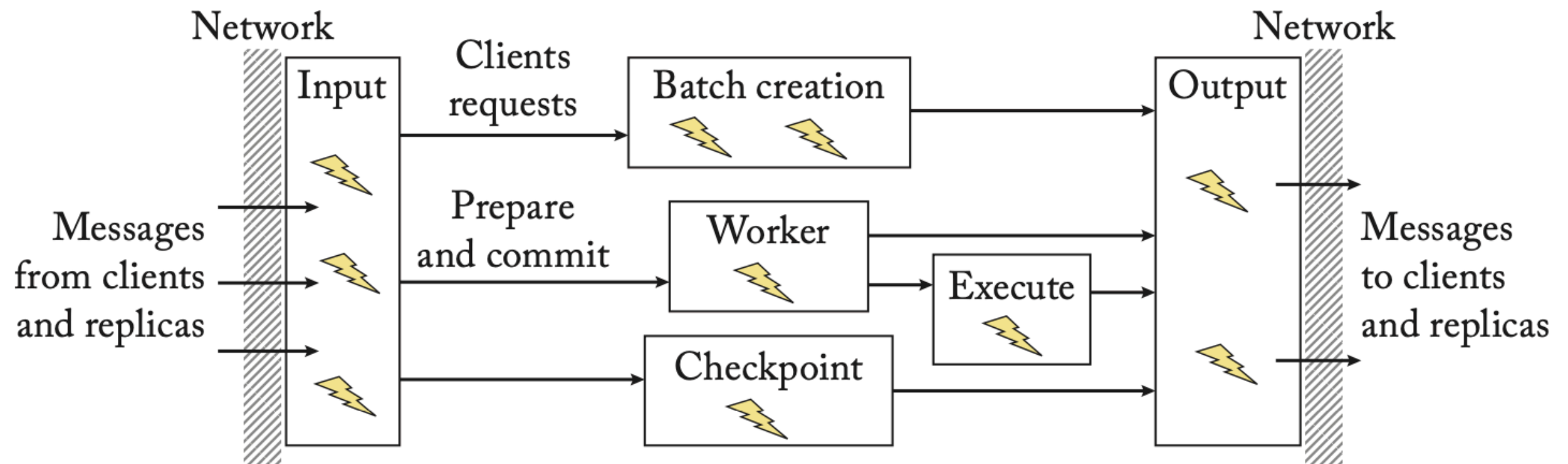
Out-of-ordering Consensus Communication

De-coupled Ordering and Execution

Off-Chain Memory Management

Expensive Cryptographic Practices (DS vs. MAC)

Smart Contracts Code Generation (Pre-compilation)



Multi-Threaded Deep Pipeline

Permissioned Blockchain Through the Looking Glass: Architectural and Implementation Lessons Learned [ICDCS'20]

Single-threaded Monolithic Design

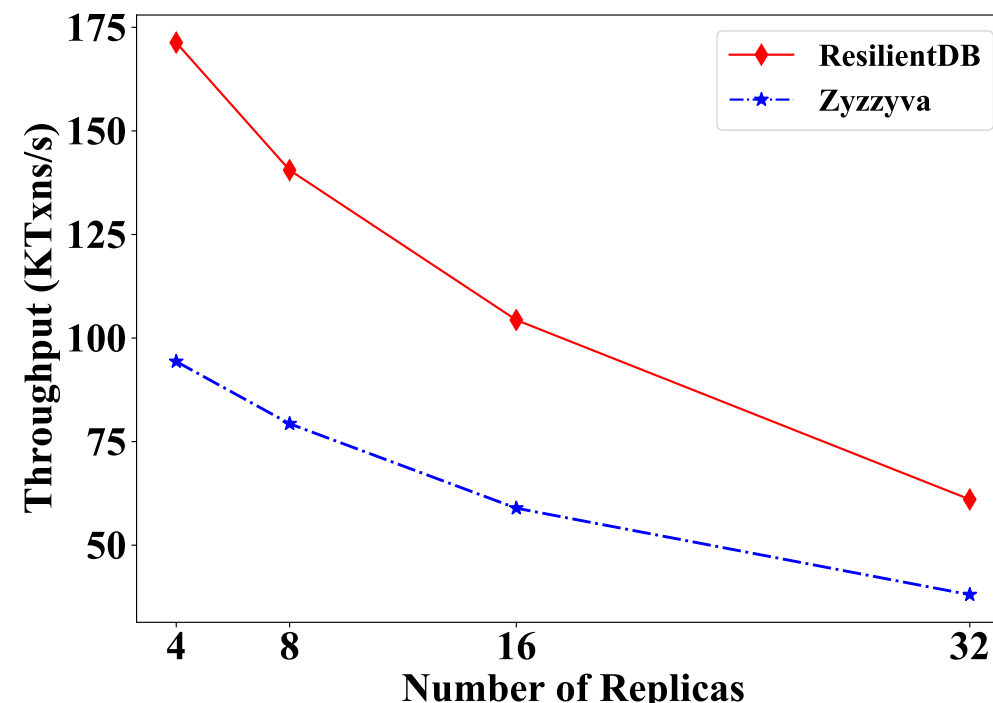
Out-of-ordering Consensus Communication

De-coupled Ordering and Execution

Off-Chain Memory Management

Expensive Cryptographic Practices (DS vs. MAC)

Smart Contracts Code Generation (Pre-compilation)



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

Revisit Resiliency

(Graduate Student Experiment Continues)

Mount Tallac, Lake Tahoe

12.1 Miles Long

3,931 Feet Elevation Gain

(9,738 Feet at Summit)



Fostering Resiliency

(Offering Stress Management and Well-Being Courses at UC Davis)



REDUCE STRESS

Spring 2020
ECS 298 (CRN 66553):
Graduate Survival Kit

Learn the foundation & working knowledge of stress reduction based on a unique heart-centered meditation practice referred to as Tamarkoz®.

The M.T.O. Tamarkoz® method is the art of self-knowledge through concentration and meditation.

**Release Tension
Increase Focus**

Days: Wednesdays
Time: 7:00 pm - 8:00 pm
Location: Zoom (Live Online Class)

INSTRUCTORS:
Mohammad Sadoghi, Ph.D.
Nasim Bahadorani, DrPH.

Computer Science Department
UC DAVIS
UNIVERSITY OF CALIFORNIA



Becoming an EXTRAORDINARY Human

Spring 2020

Days: Thursdays
Time: 7:00 pm - 8:00 pm

Location: Zoom

CRN: 57877

INSTRUCTORS:
Mohammad Sadoghi, Ph.D.
Nasim Bahadorani, DrPH.

No one wants to be ordinary. This course focuses on the personal development of the characteristics of human beings deemed extraordinary. Outcomes include enhanced concentration for higher-level cognition, increased capacity to handle stress, development of increased self-confidence, increased mastery of emotional and mental processes, development of physical awareness and control, and development of positive personal characteristics. Physical activities include movements and visualizations.

msadoghi@ucdavis.edu

Computer Science Department
UC DAVIS
UNIVERSITY OF CALIFORNIA



Reduce Stress

First-year Seminar (FYS): Undergraduate Survival Kit
Learn the foundation & working knowledge of stress reduction based on a unique heart-centered meditation method referred to as Tamarkoz®.

The M.T.O. Tamarkoz® method is the art of self-knowledge through concentration and meditation.

Spring 2020
Time: Tuesdays from 7:00pm-8:00pm
Location: UC DAVIS Zoom (CRN: 66553)

**Release Tension
Increase Focus**

INSTRUCTORS:
Mohammad Sadoghi, Ph.D.
Nasim Bahadorani, DrPH.

Dress Code: Loose comfortable clothing, sweat shirts and pants with socks.

UC DAVIS
UNIVERSITY OF CALIFORNIA

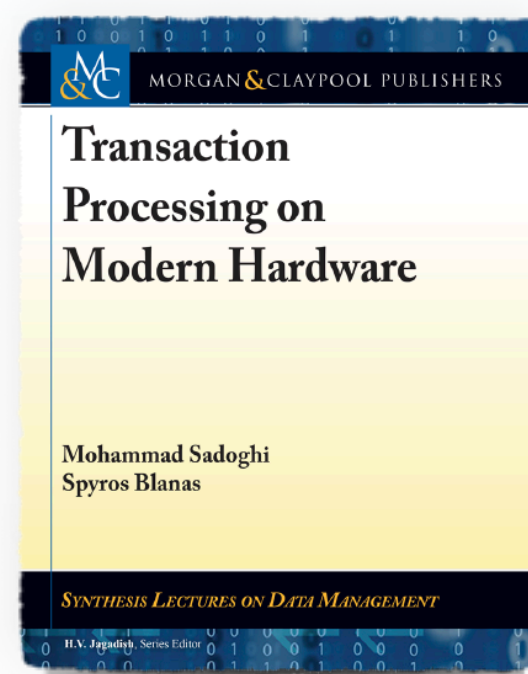
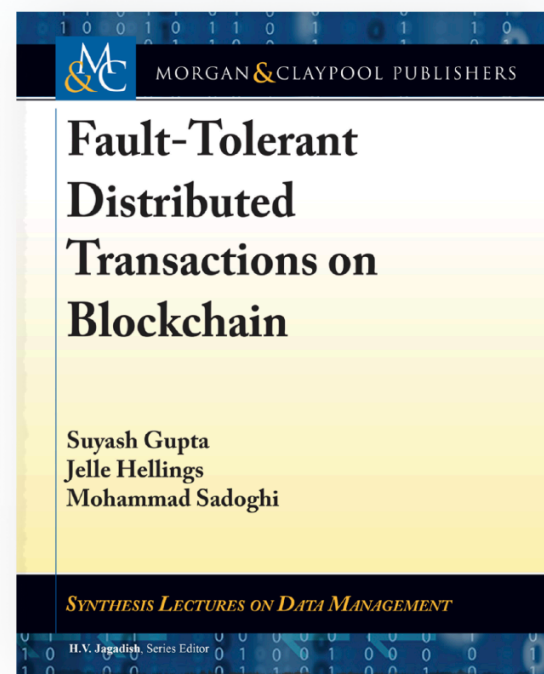


 **THE CALIFORNIA AGGIE**
Seminar spotlight: "Becoming an Extraordinary Human"
The California Aggie, April 6, 2020

 **Tamarkoz®**
BE BALANCED



THANK YOU



FOR COMPLETE REFERENCES

