

ResilientDB: Global Scale Resilient Blockchain Fabric



Suyash Gupta



Sajjad Rahnama



Jelle Hellings



Mohammad Sadoghi

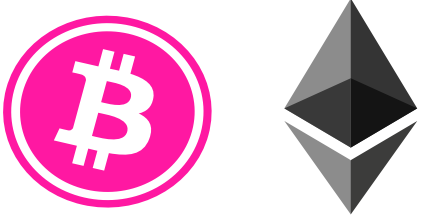
**Exploratory Systems Lab
University of California Davis**

Blockchain Applications



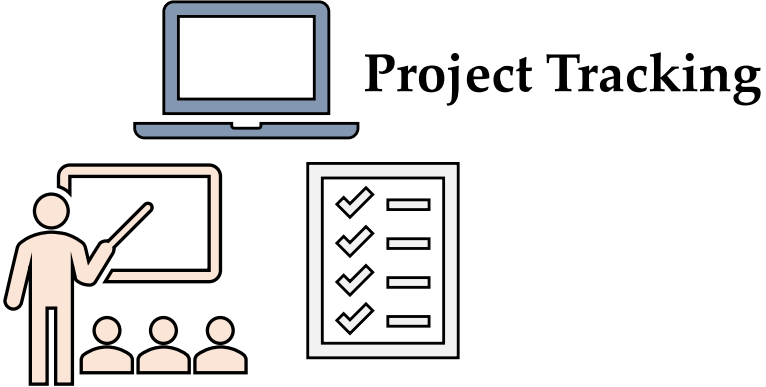
Waste Management

The Waste Management section features three icons: a yellow recycling symbol, a blue globe, and a green hand holding a plant. The text 'Waste Management' is centered below these icons.



Cryptocurrencies

The Cryptocurrencies section features two icons: a pink Bitcoin symbol and a grey Ethereum symbol. The text 'Cryptocurrencies' is centered below these icons.



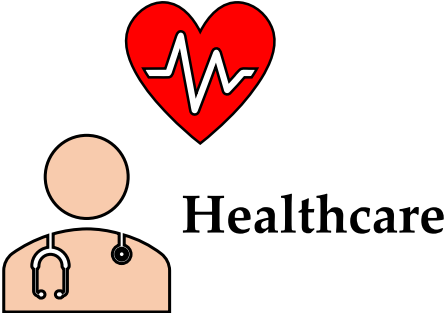
Project Tracking

The Project Tracking section features three icons: a laptop, a person presenting to an audience, and a checklist. The text 'Project Tracking' is centered to the right of these icons.



Energy Trading and Management

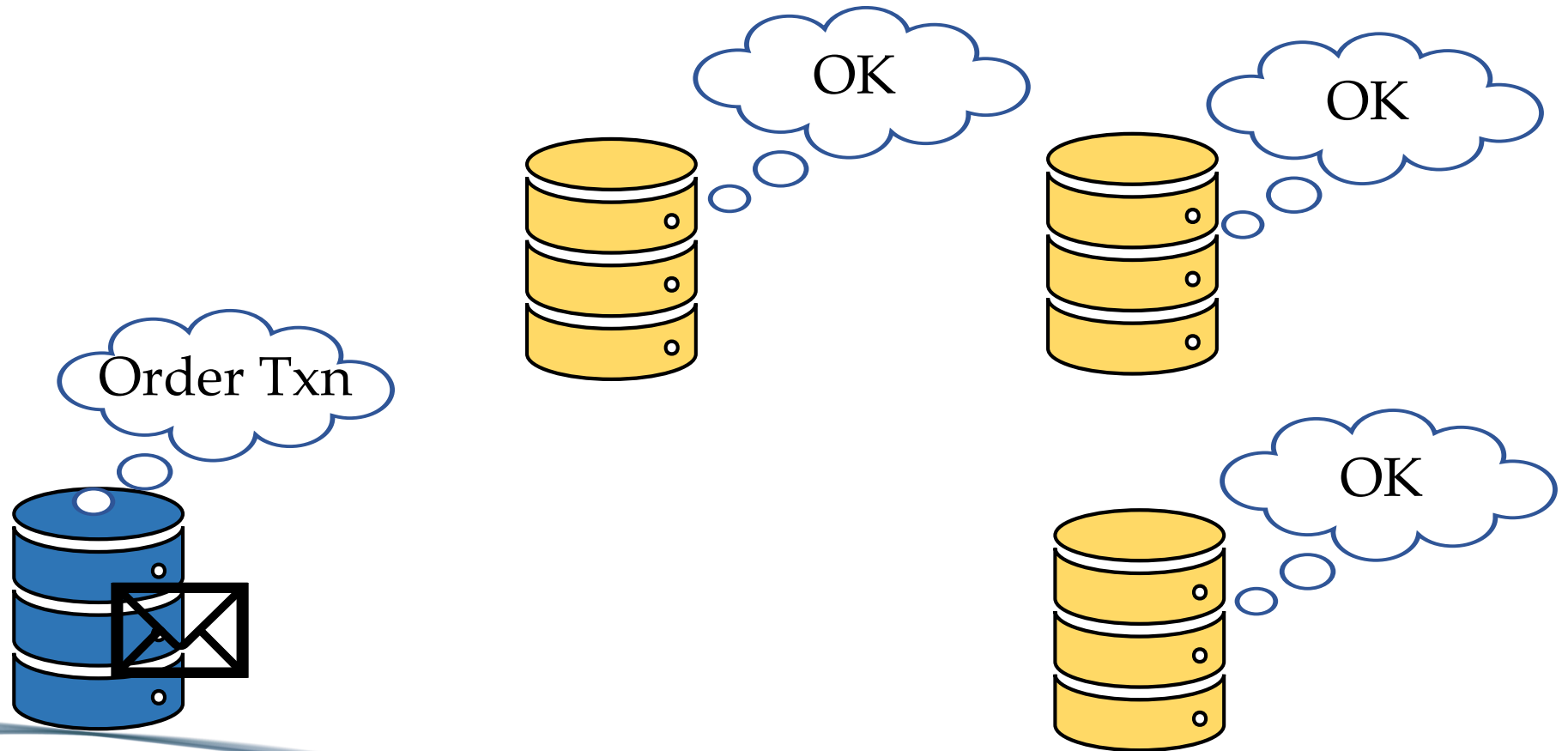
The Energy Trading and Management section features two icons: a yellow lightning bolt and a brown line graph showing an upward trend. The text 'Energy Trading and Management' is centered below these icons.



Healthcare

The Healthcare section features two icons: a red heart with a white ECG line and a brown doctor icon with a stethoscope. The text 'Healthcare' is centered to the right of these icons.

At the core of *any* Blockchain application is a Byzantine Fault-Tolerant (BFT) consensus protocol.



Challenges For Geo-Scale Blockchains

	<i>Ping round-trip times (ms)</i>						<i>Bandwidth (Mbit/s)</i>					
	<i>O</i>	<i>I</i>	<i>M</i>	<i>B</i>	<i>T</i>	<i>S</i>	<i>O</i>	<i>I</i>	<i>M</i>	<i>B</i>	<i>T</i>	<i>S</i>
Oregon (<i>O</i>)	≤ 1	38	65	136	118	161	7998	669	371	194	188	136
Iowa (<i>I</i>)		≤ 1	33	98	153	172		10004	752	243	144	120
Montreal (<i>M</i>)			< 1	82	186	202			7977	283	111	102
Belgium (<i>B</i>)				≤ 1	252	270				9728	79	66
Taiwan (<i>T</i>)					≤ 1	187					7998	100
Sydney (<i>S</i>)												

High!

Low!

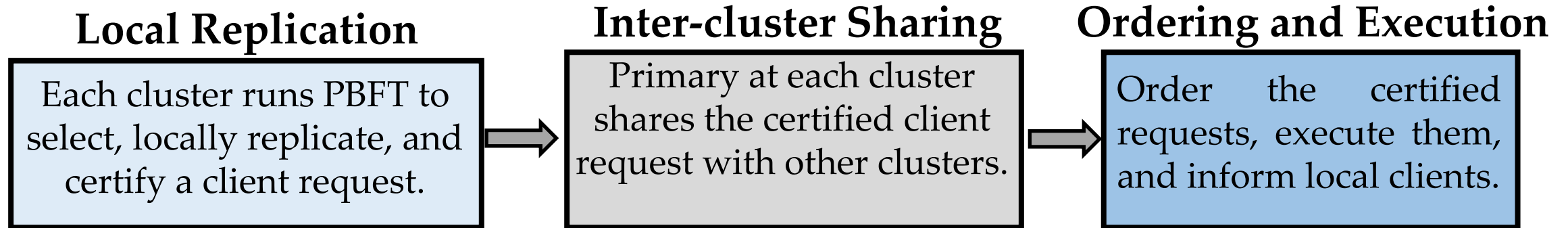
Blockchain expects decentralization → Replicas maybe geo-distributed

Limitations of Existing Consensus Protocols

Protocol	Decisions	Communication		Centralized
		(Local)	(Global)	
GEOBFT (our paper)	z	$\mathcal{O}(2zn^2)$	$\mathcal{O}(fz^2)$	No
↳ <i>single decision</i>	1	$\mathcal{O}(4n^2)$	$\mathcal{O}(fz)$	No
STEWARD	1	$\mathcal{O}(2zn^2)$	$\mathcal{O}(z^2)$	Yes
ZYZZYVA	1	$\mathcal{O}(zn)$		Yes
PBFT	1	$\mathcal{O}(2(zn)^2)$		Yes
PoE	1	$\mathcal{O}((zn)^2)$		Yes
HOTSTUFF	1	$\mathcal{O}(8(zn))$		Partly

Normal-case metrics for a system with z clusters, each with n replicas of which at most f , $n > 3f$, are Byzantine

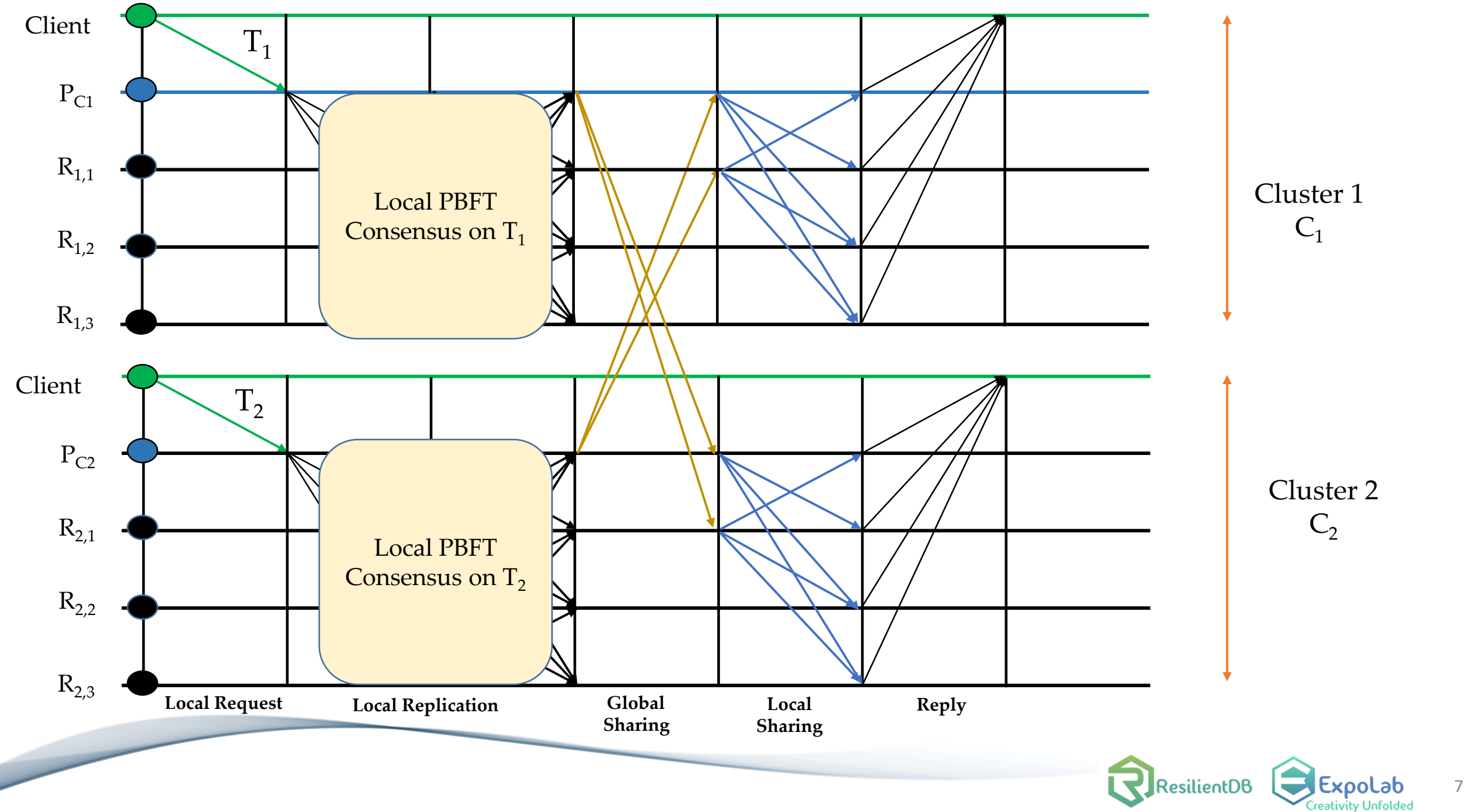
GeoBFT Protocol



Topology-aware protocol.

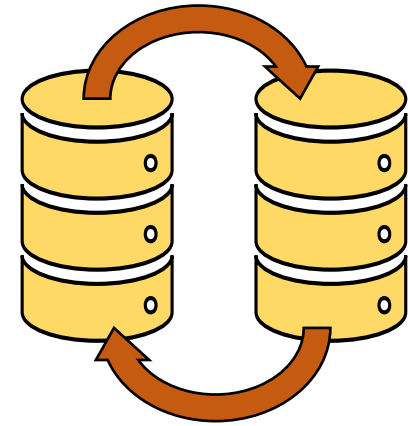
Parallel consensus at each cluster.

Low inter-cluster communication overheads.

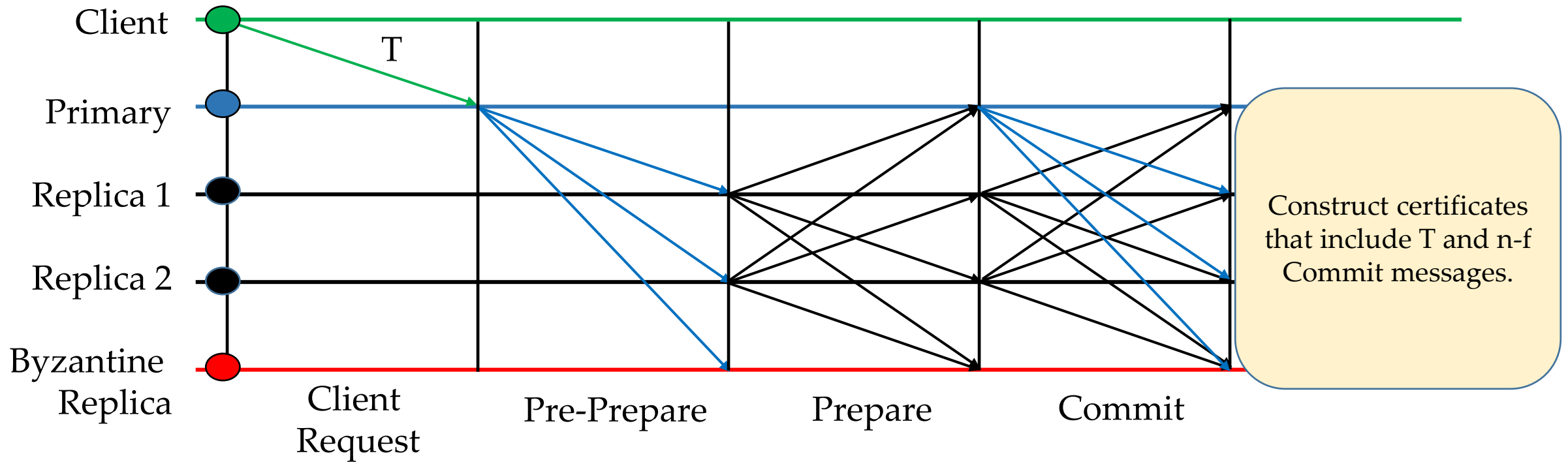


Local Replication

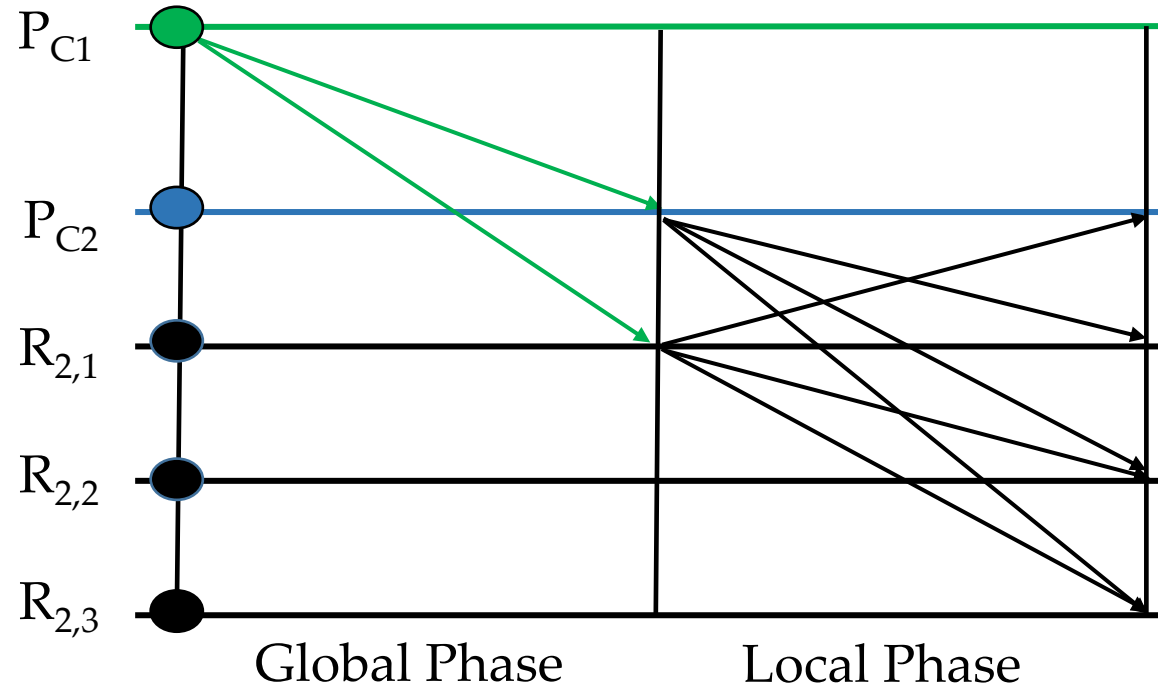
- GeoBFT employs PBFT for local replication.
- Tolerates up to f failure out of $3f+1$ replicas
- Three phases of which two require quadratic communication complexity.
- **Safety** is **always** guaranteed and **Liveness** is guaranteed in **periods of partial synchrony**.
- **View-Change** protocol for replacing malicious primary



PBFT Civil Execution



Inter-Cluster Sharing



The Primary P_{C1} sends a certificate that includes the client request and commit messages from $n-f$ replicas of Cluster C_1 .

Ordering and Execution

GeoBFT orders requests deterministically.

For $i < j$, requests of Cluster C_i are executed before requests of cluster C_j .

For example: requests of C_1 are executed before C_2 .

Implementation on ResilientDB



ResilientDB associates a multi-threaded deep-pipelined architecture with each replica.

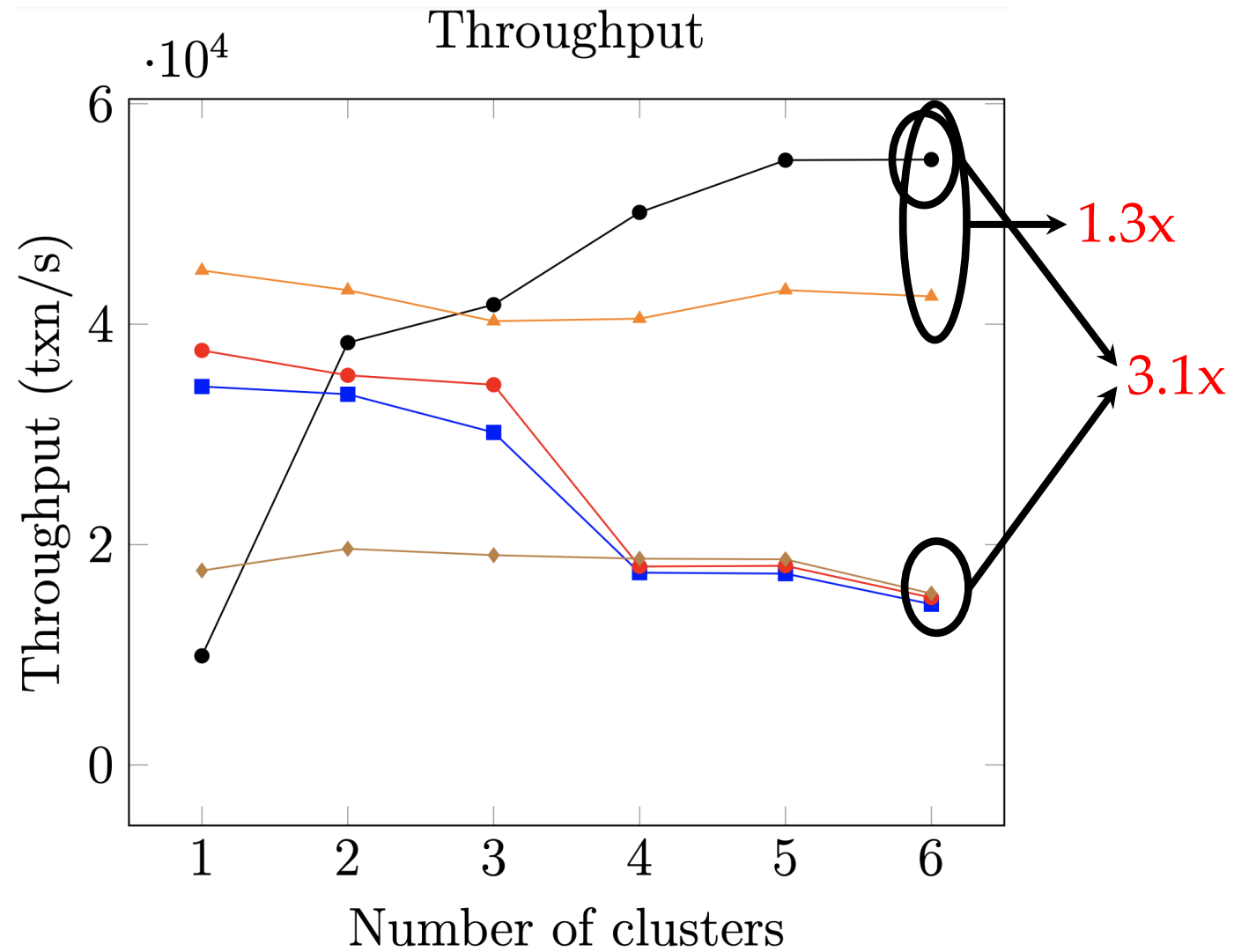
Ledger (Blockchain) Management

- i^{th} block in the ledger contains the i^{th} executed request.
- In each round of GeoBFT, each replica executes z requests, each belonging to a different cluster C_i , $1 \leq i \leq z$.
- Hence, in each round, each replica creates z blocks.
- To ensure immutability, each block includes both client requests and exchanged certificates.

Evaluation on ResilientDB

- Google cloud used for deploying replicas and clients.
- Each replica used 8-core Intel Skylake CPUs and had access to 16 GB memory.
- Total 160K clients deployed on eight 4-core machines.
- Workload provided by Yahoo Cloud Serving Benchmark (YCSB).
- Replicas deployed across six different regions: Oregon, Iowa, Montreal, Belgium, Taiwan and Sydney.
- Primaries for centralized protocol placed at Oregon (highest bandwidth).

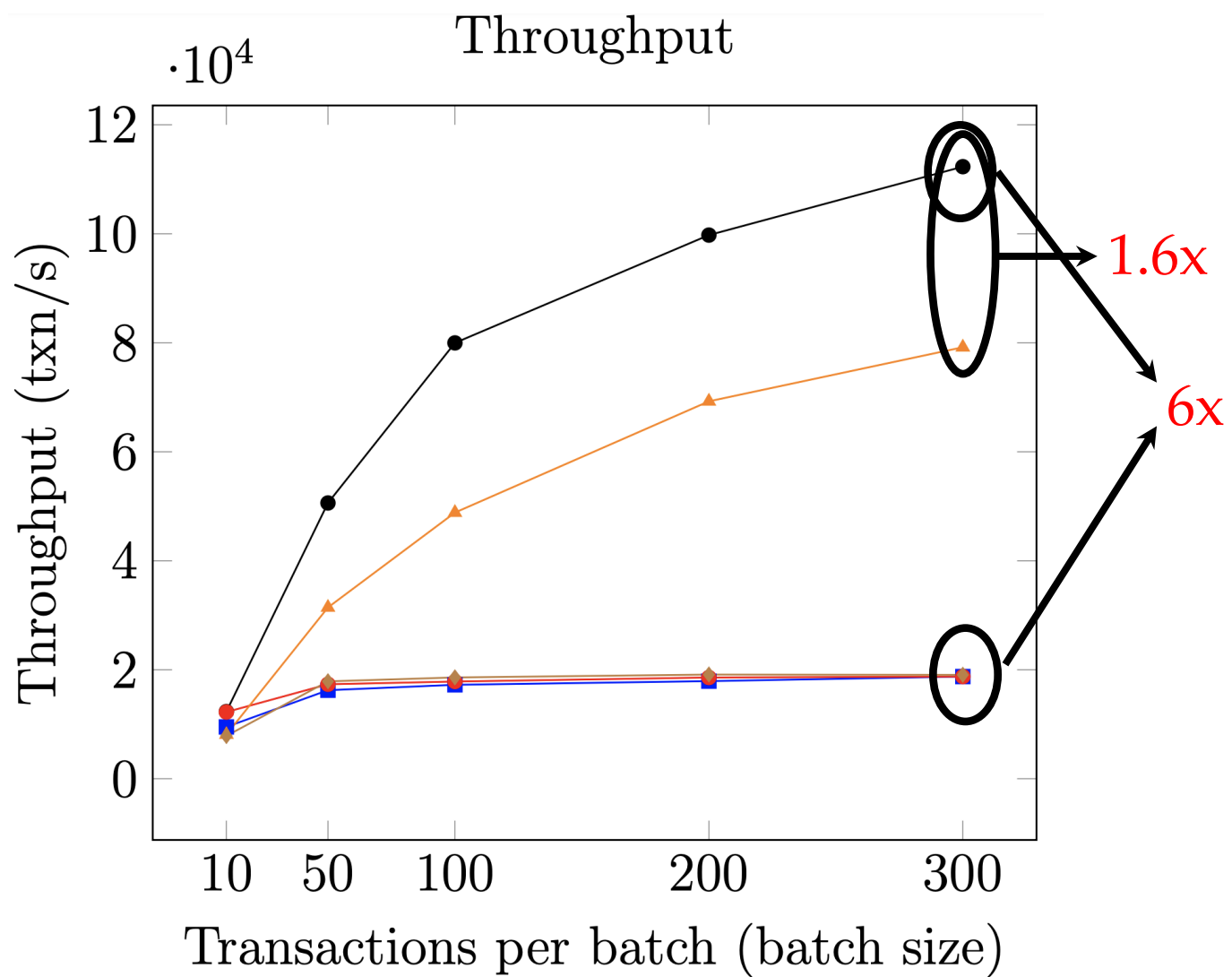
Scalability



—●— GEOBFT —■— PBFT —●— ZYZZYVA —▲— HOTSTUFF —◆— STEWARD

Total replicas = 60

Batching



● GEOBFT ■ PBFT ● ZYZZYVA ▲ HOTSTUFF ◆ STEWARD

4 clusters and 7 replicas per cluster

Conclusions and Final Remarks

- For achieving faster local replication, GeoBFT can employ other efficient BFT protocols, such as PoE.
- Modern cryptographic techniques such as Threshold signatures can be used in place of sending $n-f$ Commit messages.
- If a cluster does not have a request, it can send “no-op” messages.
- GeoBFT optimizes consensus by reducing global communication costs.
- Parallel local replication helps to increase system throughput.

References

1. S. Gupta, S. Rahnema, J. Hellings, and M. Sadoghi. ResilientDB: Global Scale Resilient Blockchain Fabric. Proc. VLDB Endow., 13(6):868–883, Feb. 2020.
2. Miguel Castro and Barbara Liskov. Practical byzantine fault tolerance. In Proceedings of the Third Symposium on Operating Systems Design and Implementation, pages 173–186. USENIX Association, 1999.
3. Miguel Castro and Barbara Liskov. Practical byzantine fault tolerance and proactive recovery. ACM Transactions on Computer Systems, 20(4):398–461, 2002. doi:10.1145/571637.571640.
4. M. Yin, D. Malkhi, M. K. Reiter, G. G. Gueta, and I. Abraham. HotStuff: BFT consensus with linearity and responsiveness. In Proceedings of the 2019 ACM Symposium on Principles of Distributed Computing, PODC, pages 347–356. ACM, 2019.
5. R. Kotla, L. Alvisi, M. Dahlin, A. Clement, and E. Wong. Zyzzyva: Speculative byzantine fault tolerance. In Proceedings of Twenty-first ACM SIGOPS Symposium on Operating Systems Principles, SOSP, pages 45–58. ACM, 2007.
6. Yair Amir, Claudiu Danilov, Danny Dolev, Jonathan Kirsch, John Lane, Cristina Nita-Rotaru, Josh Olsen, and David Zage. Steward: Scaling byzantine fault-tolerant replication to wide area networks. IEEE Transactions on Dependable and Secure Computing, 7(1):80–93, 2010. doi:10.1109/TDSC.2008.53.
7. S. Gupta, J. Hellings, S. Rahnema, and M. Sadoghi, Proof-of-Execution: Reaching Consensus through Fault-Tolerant Speculation, CoRR, vol. abs/1911.00838, 2019.