



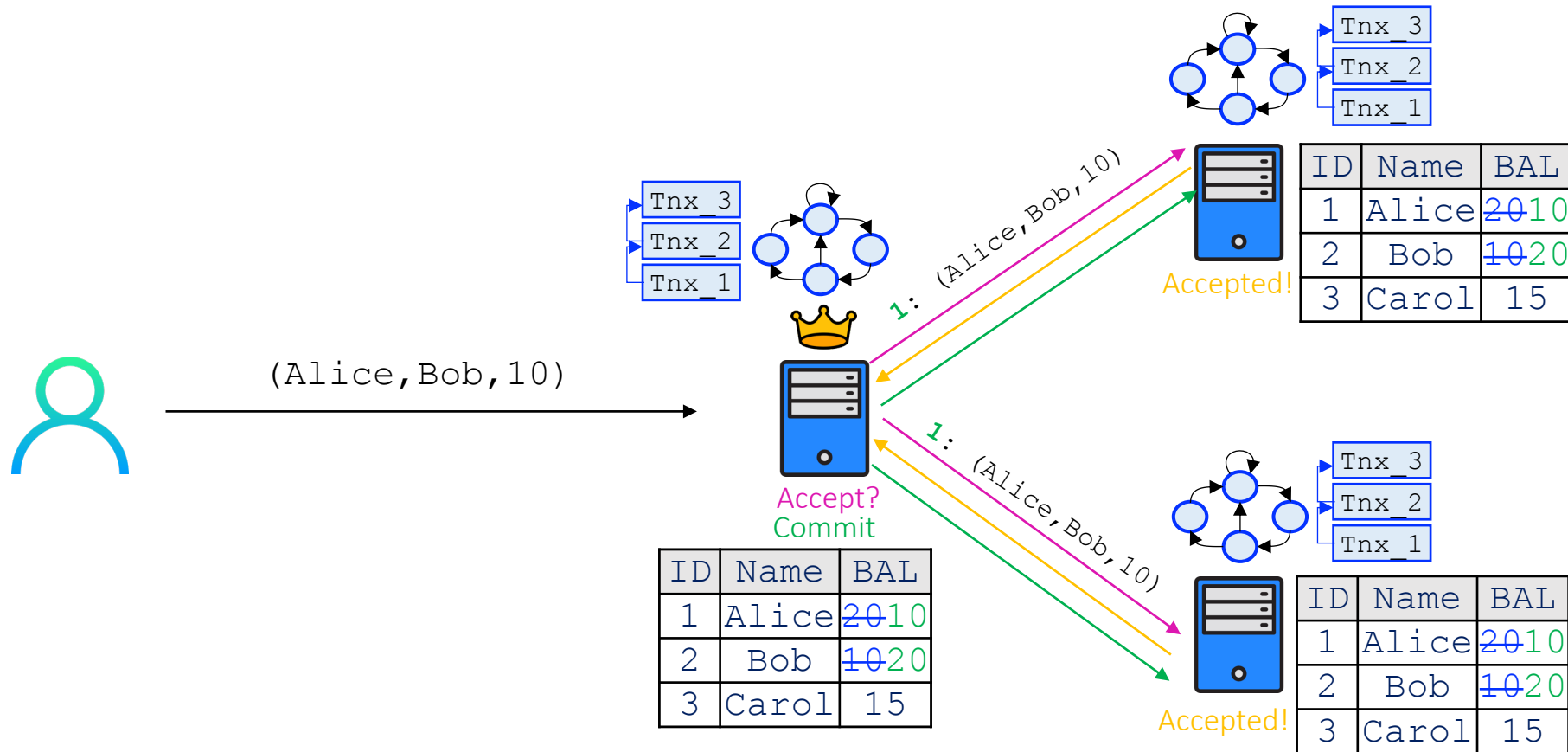
The Bedrock of Byzantine Fault Tolerance

A Unified Platform for BFT Protocols Analysis, Implementation, and Experimentation

Mohammad Javad Amiri¹, Chenyuan Wu², Divyakant Agrawal³, Amr El Abbadi³,
Boon Thau Loo², Mohammad Sadoghi⁴

¹Stony Brook University, ²UPenn, ³UC Santa Barbara, ⁴UC Davis

Distributed transaction processing

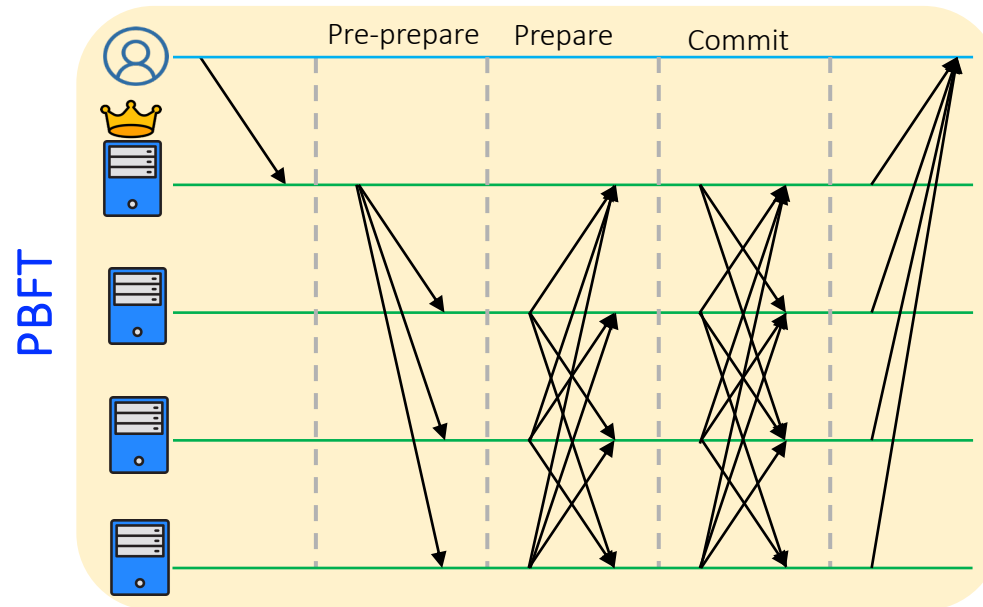


State Machine Replication: a replicated service whose state is mirrored across different deterministic replicas

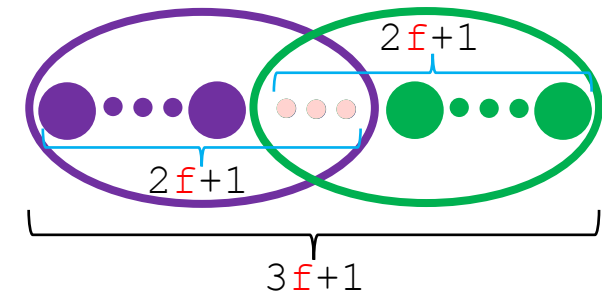
- Assign **order** to each client request in the global service history and execute it in that order

Byzantine fault-tolerant protocol: PBFT

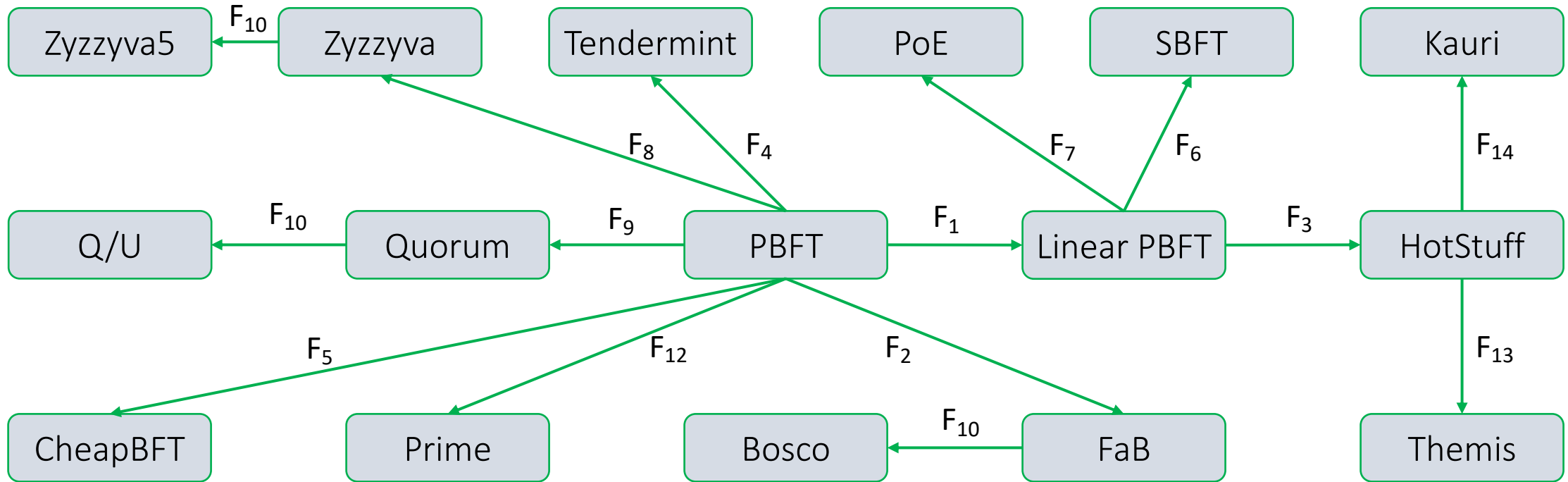
Nodes can fail arbitrarily, including deviating from the protocol



At Most f
Byzantine Failures



BFT protocols landscape

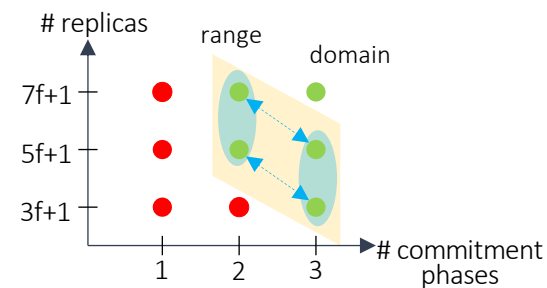


What protocol best fits our needs?

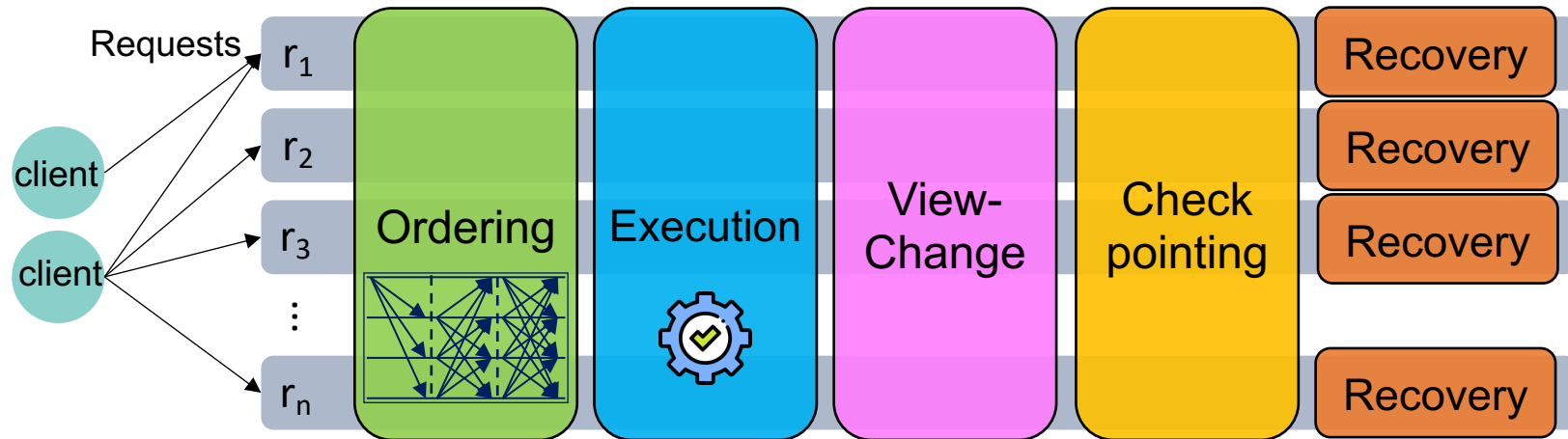
Analysis Implementation Experimentation

BFT protocols design space and design dimensions

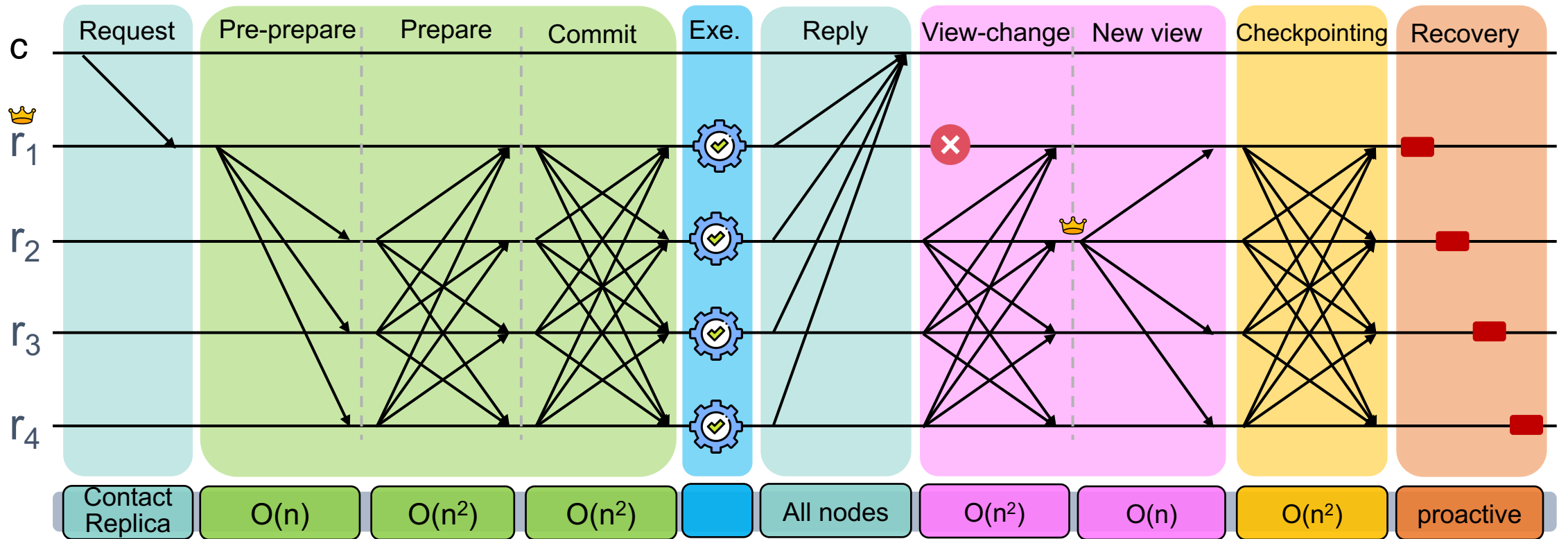
- Design space
 - A set of dimensions to analyze BFT protocols
- Design choices
 - Trade-offs between dimensions
 - A set of one-to-one functions, each maps protocols in its domain to protocols in its range
- Focus on partially synchronous BFT protocols



Different stages of replicas in a BFT protocol



PBFT



Design space of BFT protocols

Protocol structure

- P1. Commitment strategy
- P2. Number of commitment phases
- P3. View-change
- P4. Checkpointing
- P5. Recovery
- P6. Types of clients

Environmental Settings

- E1. Number of replicas
- E2. Communication topology
- E3. Authentication
- E4. Responsiveness, synchronization, and timers

Quality of Service

- Q1. Order-fairness
- Q2. Load balancing

Performance Optimization

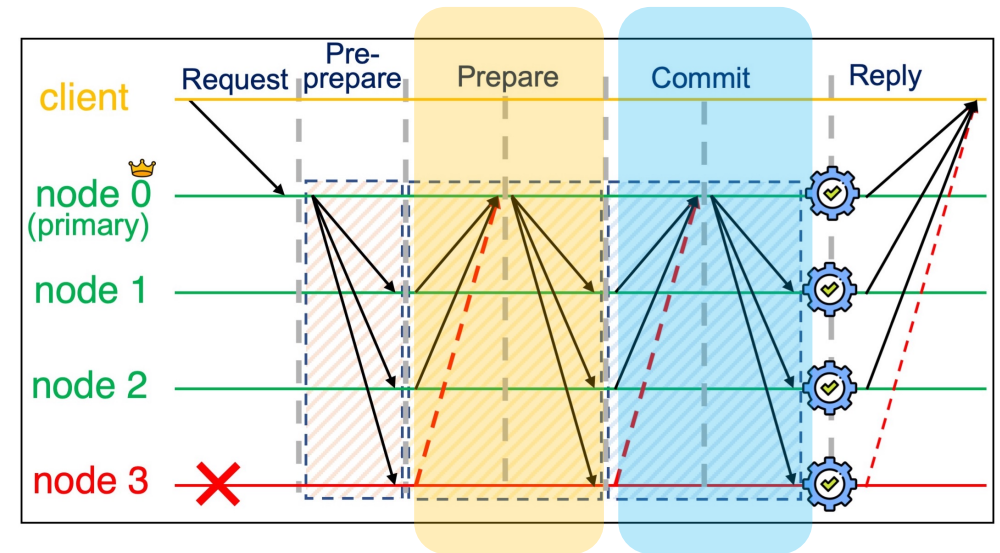
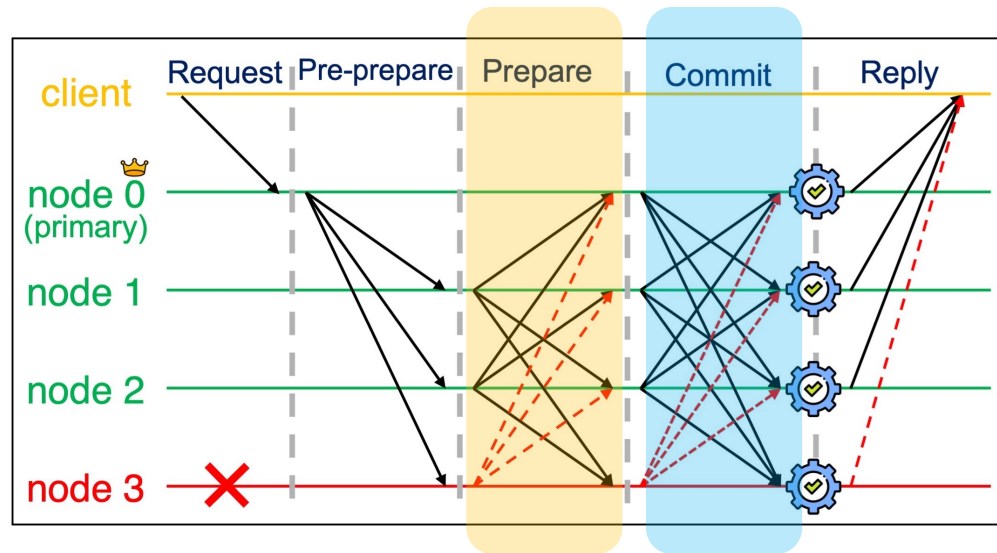
- O1. Out-of-order processing
- O2. Request pipelining
- O3. Parallel ordering
- O4. Parallel execution
- O5. Read-only requests processing
- O6. Separating ordering and execution
- O7. Trusted hardware
- O8. Request/reply dissemination

Design choices

1. Linearization
2. Phase reduction through redundancy
3. Leader rotation
4. Non-responsive leader rotation
5. Optimistic replica reduction
6. Optimistic phase reduction
7. Speculative phase reduction
8. Speculative execution
9. Optimistic conflict-free
10. Resilience
11. Authentication
12. Robust
13. Fair
14. Tree-based LoadBalancer

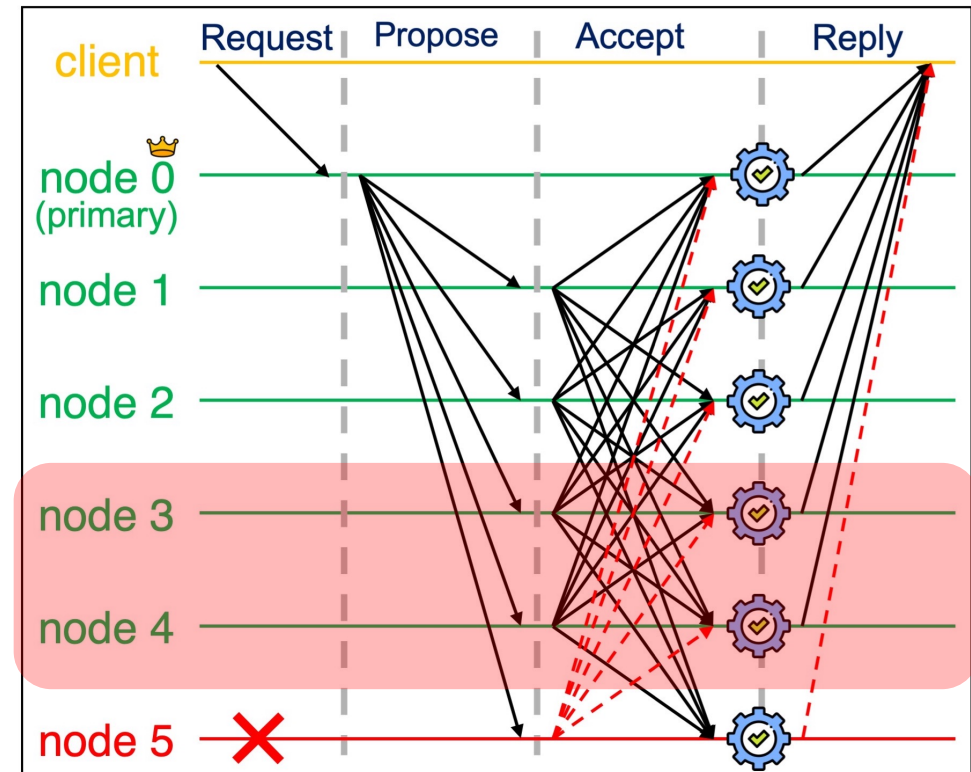
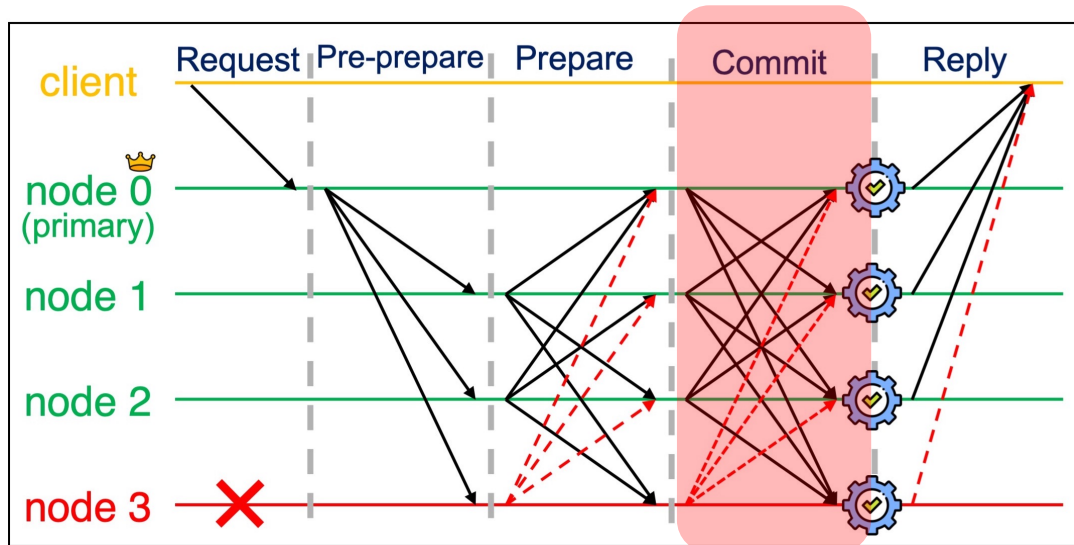
Design choice 1: Linearization

- Trade-off between communication topology and communication phases.
 - Linear PBFT
 - The collector needs to send a certificate of having received the required signatures.



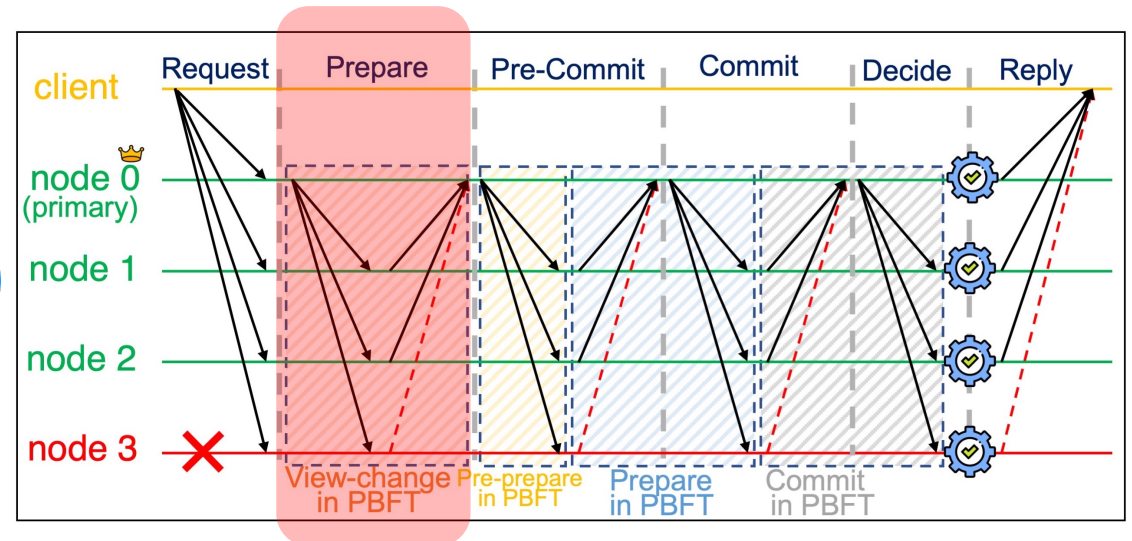
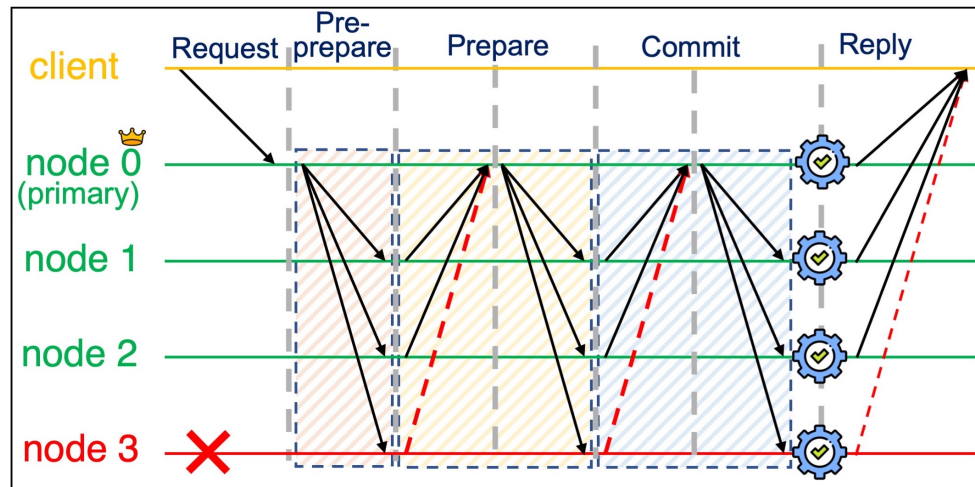
Design choice 2: Phase reduction through redundancy

- Trade-off between the number of ordering phases and the number of replicas
 - FaB



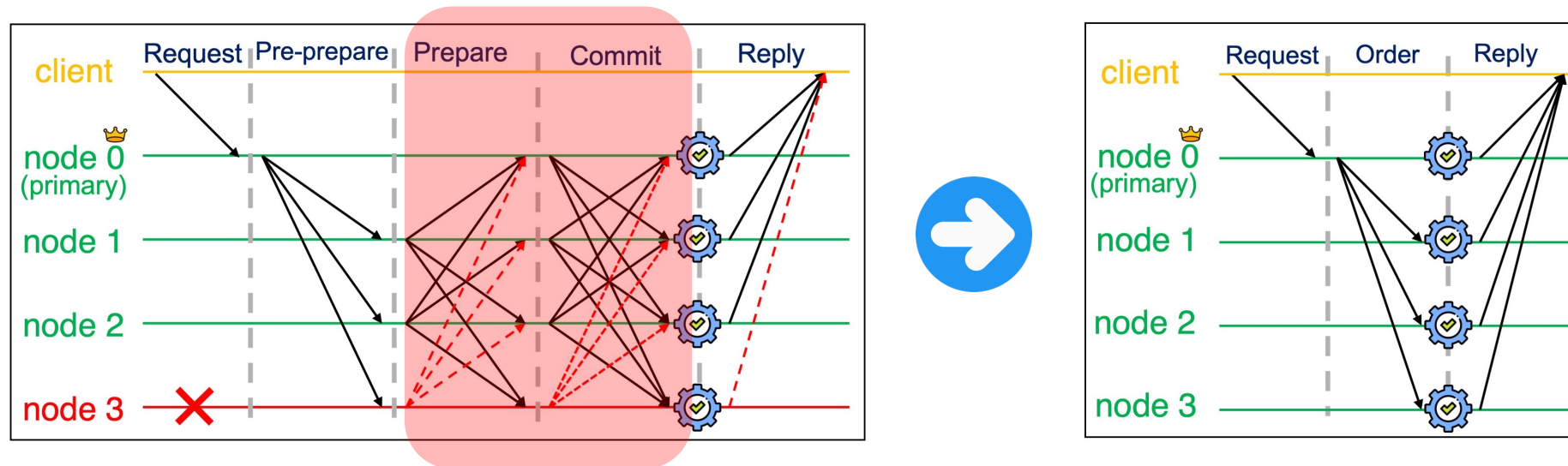
Design choice 3: Leader rotation

- Replace the stable leader with the rotating leader mechanism by adding one phase
 - HotStuff



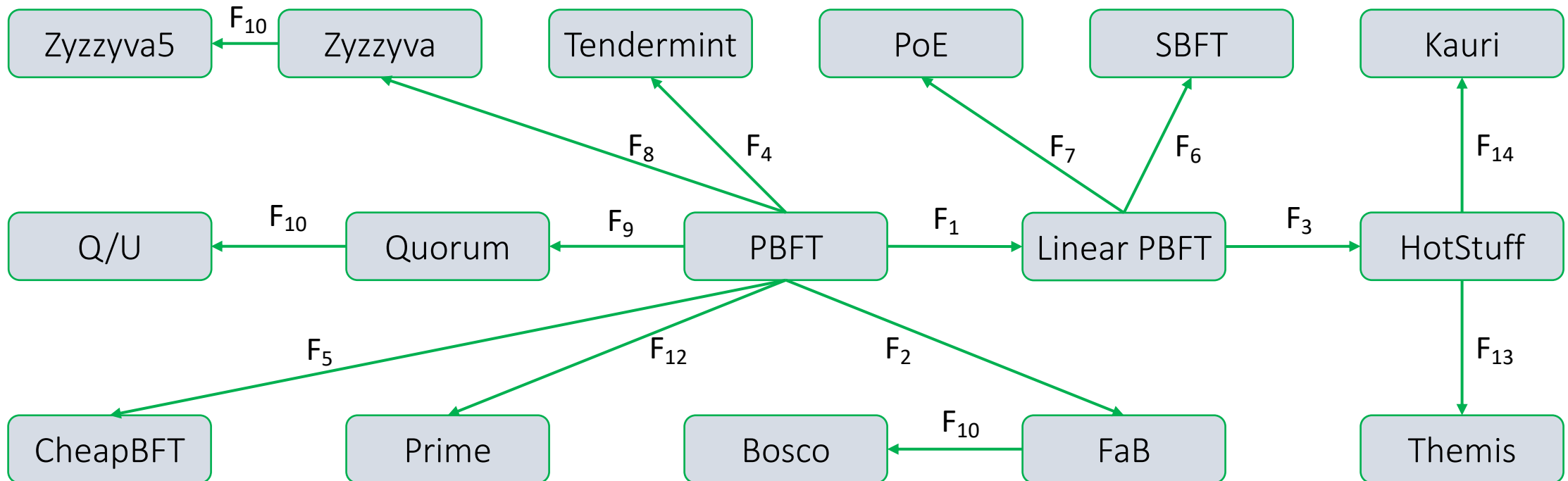
Design choice 8: Speculative execution

- Eliminate the prepare and commit phases while optimistically assuming that all replicas are non-faulty
 - Zyzzzyva



Derivation of protocols from PBFT using design choices

1. Linearization
2. Phase reduction through redundancy
3. Leader rotation
4. Non-responsive leader rotation
5. Optimistic replica reduction
6. Optimistic phase reduction
7. Speculative phase reduction
8. Speculative execution
9. Optimistic conflict-free
10. Resilience
11. Authentication
12. Robust
13. Fair
14. Tree-based LoadBalancer



Implementation

- The core unit

- Defines entities, e.g., clients and nodes, and maintains the application logic and data
- Defines workloads and benchmarks

The state manager

- Enables the core unit to track the states and transitions of each entity according to the protocol
- Defines a [domain-specific language \(DSL\)](#) to rapidly prototype BFT protocols

The plugin manager

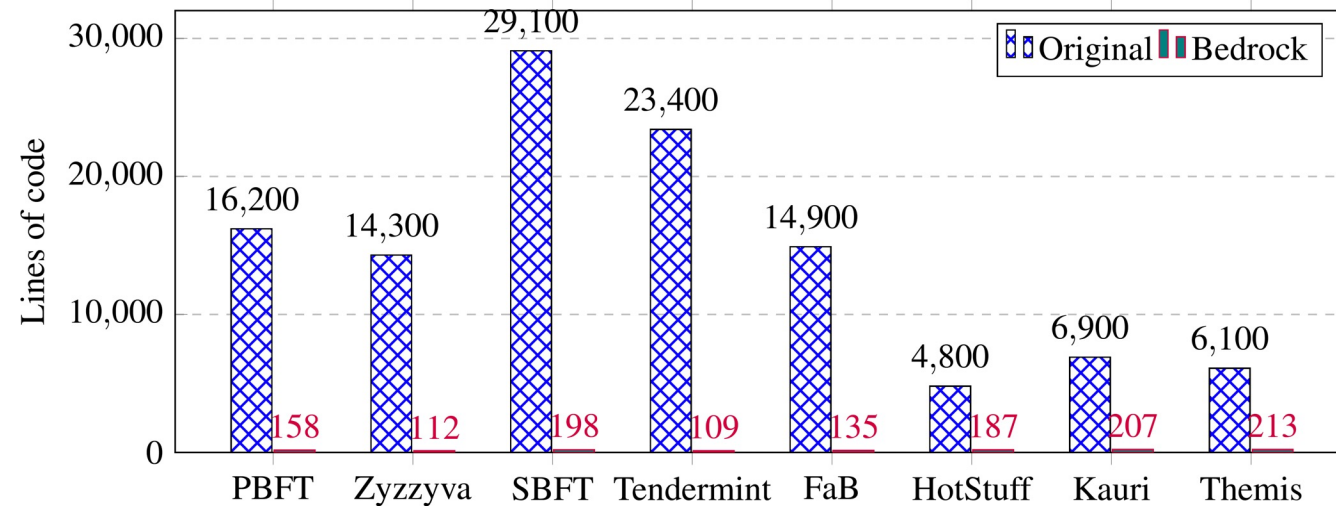
- Implements [protocol-specific behaviors](#) that cannot be handled by the protocol config
- Enables users to [define their own dimensions/values](#) or to update existing dimensions without requiring changes to the platform code or rebuilding the platform binaries

The run-time unit

- Manages the run-time execution
- E.g., manages benchmarks, setups all entities, enables plugins to run, reports results

DSL code

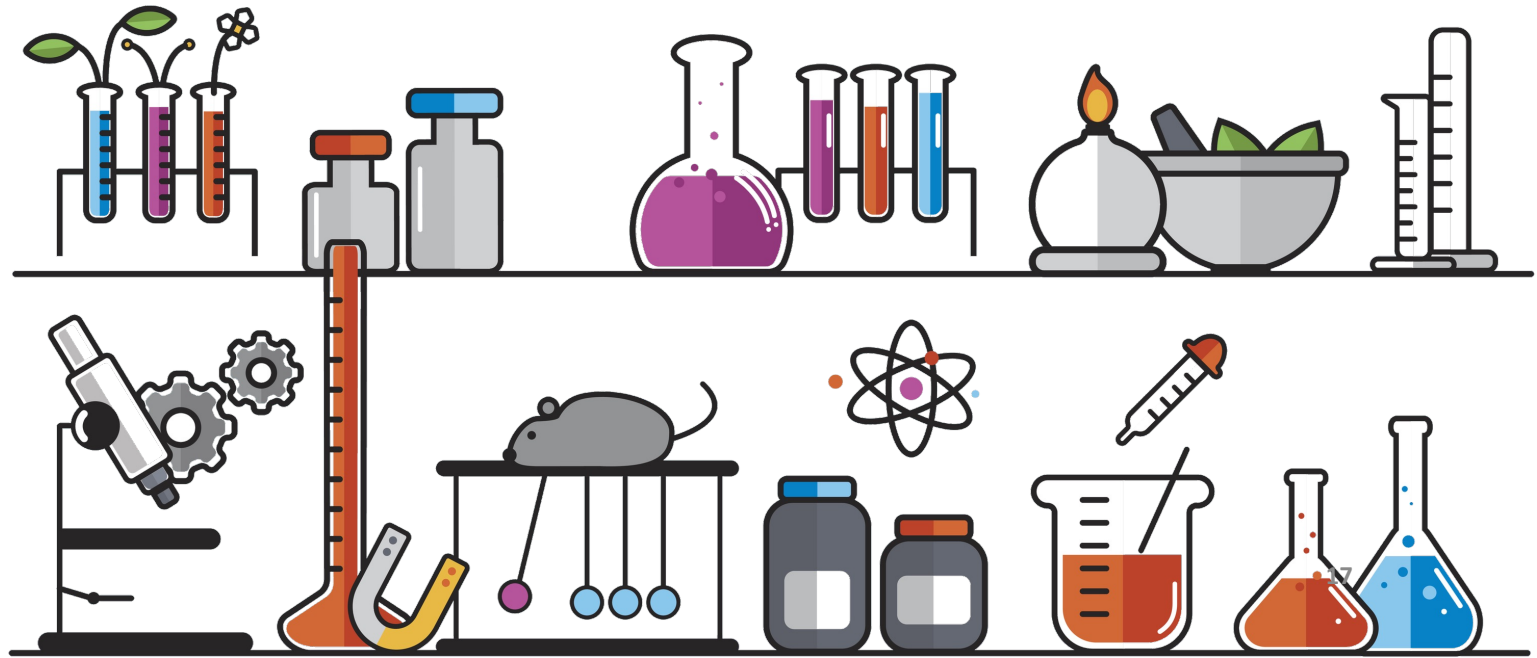
- Written in the protocol config
- Defines different dimensions and their chosen values, the list of roles, phases, states, messages, quorum conditions, and plugins
- Reduces the effort needed to write a BFT protocol



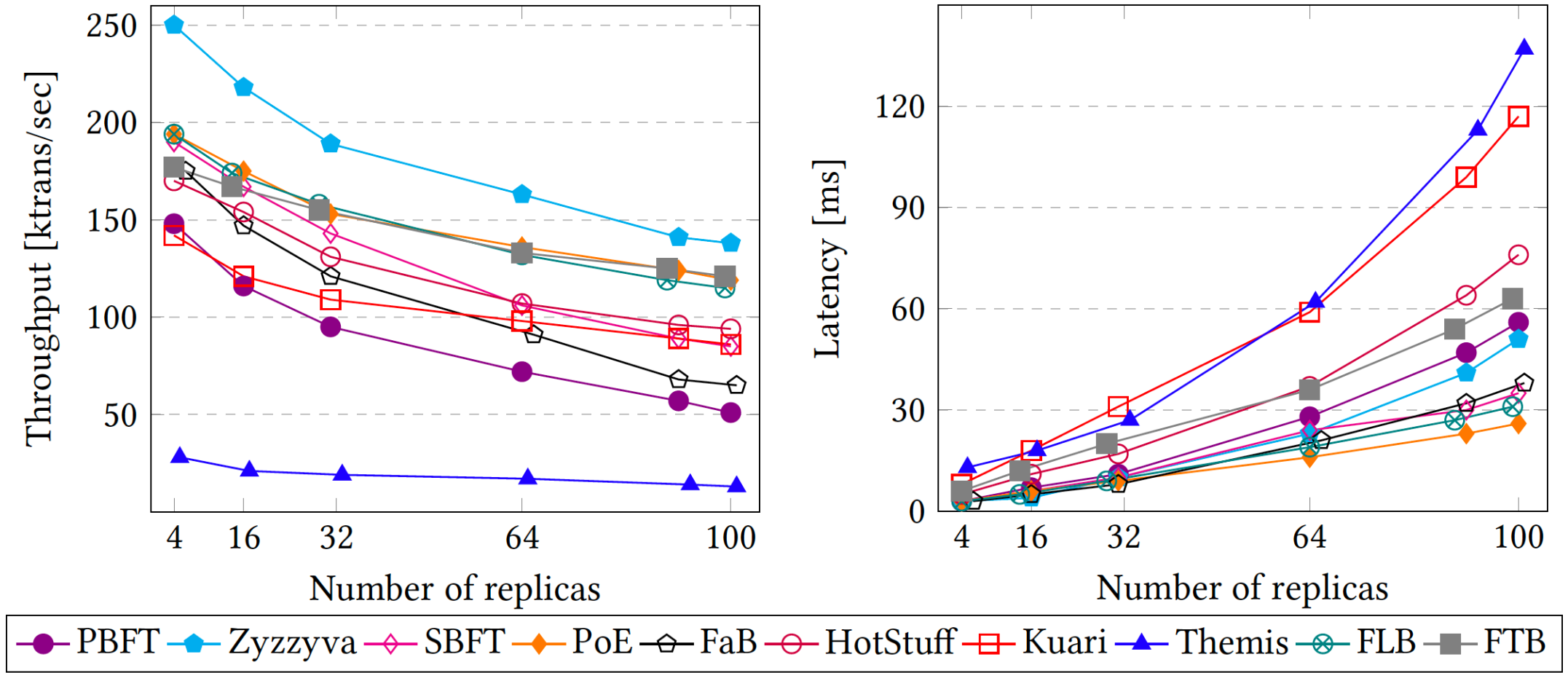
```
11 protocol:
12   general:
13     leader: stable
14     requestTarget: primary
15
16   roles:
17     - primary
18     - nodes
19     - client
20
21   phases:
22     - name: normal
23       states:
24         - idle
25         - wait_prepare
26         - wait_commit
27         - executed
28       messages:
29         - name: request
30           requestBlock: true
31         - name: reply
32           requestBlock: true
33         - name: preprepare
34           requestBlock: true
35         - prepare
36         - commit
37     - name: view_change
38       states:
39         - wait_view_change
40         - wait_new_view
41       messages:
42         - view_change
43         - new_view
44     - name: checkpoint
45       messages:
46         - checkpoint
47
48   transitions:
49     from:
50       - role: client
51         state: idle
52     to:
53       - state: executed
54         update: sequence
55         condition:
56           type: msg
57           message: reply
58           quorum: 2f + 1
```


Experimental settings

- Platform: [Amazon EC2](#)
- PBFT, Zyzzyva, SBFT, FaB, PoE, (Chained-)HotStuff, Kauri, Themis, FLB and FTB.
- Evaluate the impact of the impact of design choices 1, 2, 3, 6, 7, 8, 10, 11, 13, and 14
- Workload with client request/reply payload sizes of 128/128 byte.
- Measuring performance
 - Throughput
 - Latency

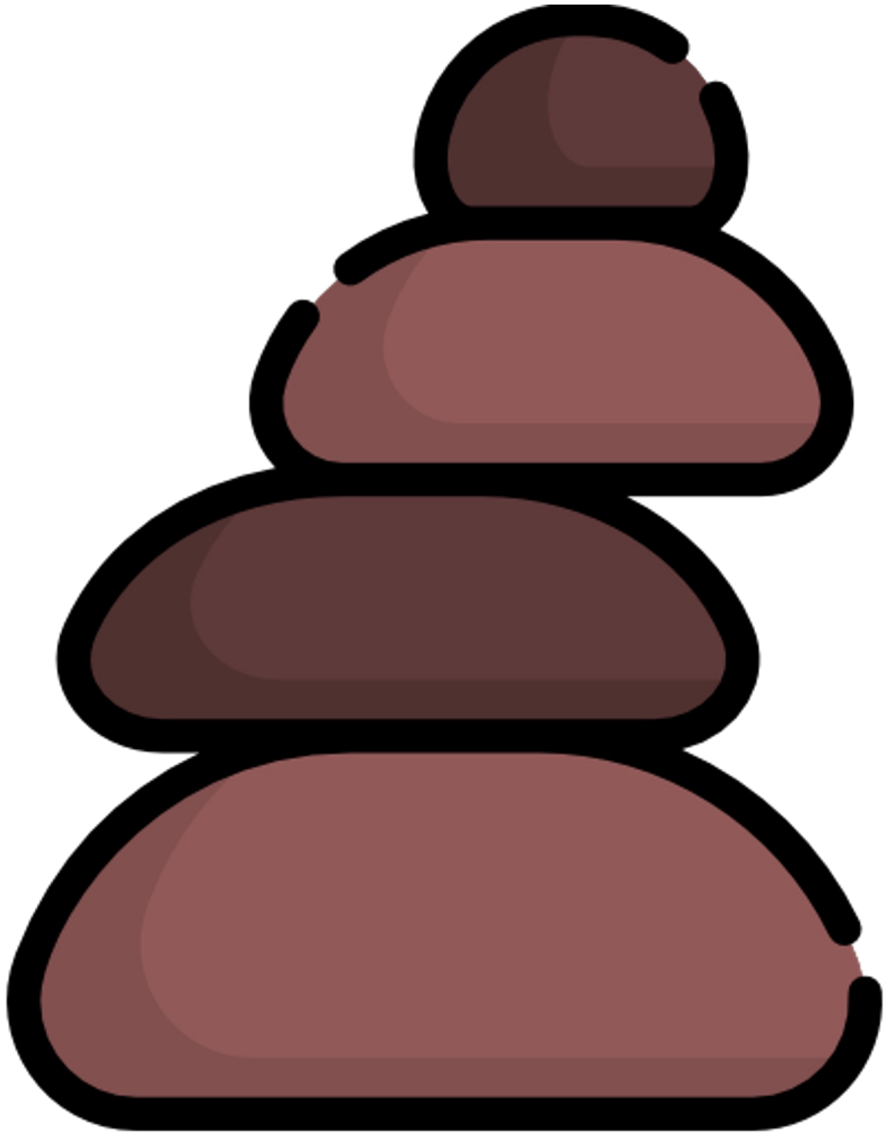


Performance with different number of replicas



Future work

- Enabling users to check the correctness of their written protocols
 - Transforming the DSL code written in Bedrock to the language used by verification tools
- Diversifying replica implementation using n-version programming
 - To ensure the independent failure of replicas
- Extending the supported protocols
 - E.g., adding synchronous and fully asynchronous protocols
- Enabling scalable transaction processing
 - Running different instances of consensus protocols in parallel
- Incorporating automatic selection strategies in Bedrock
 - Using machine learning to select the appropriate BFT protocol, or switch protocols at runtime



Questions?