

Algorithms and Distributed Systems 2023/2024 (Lab Six)

**MIEI - Integrated Master in Computer Science and
Informatics**

**MEI – Master in Computer Science and
Informatics**

Specialization block

Nuno Preguiça (nmp@fct.unl.pt)

Alex Davidson (a.davidson@fct.unl.pt)



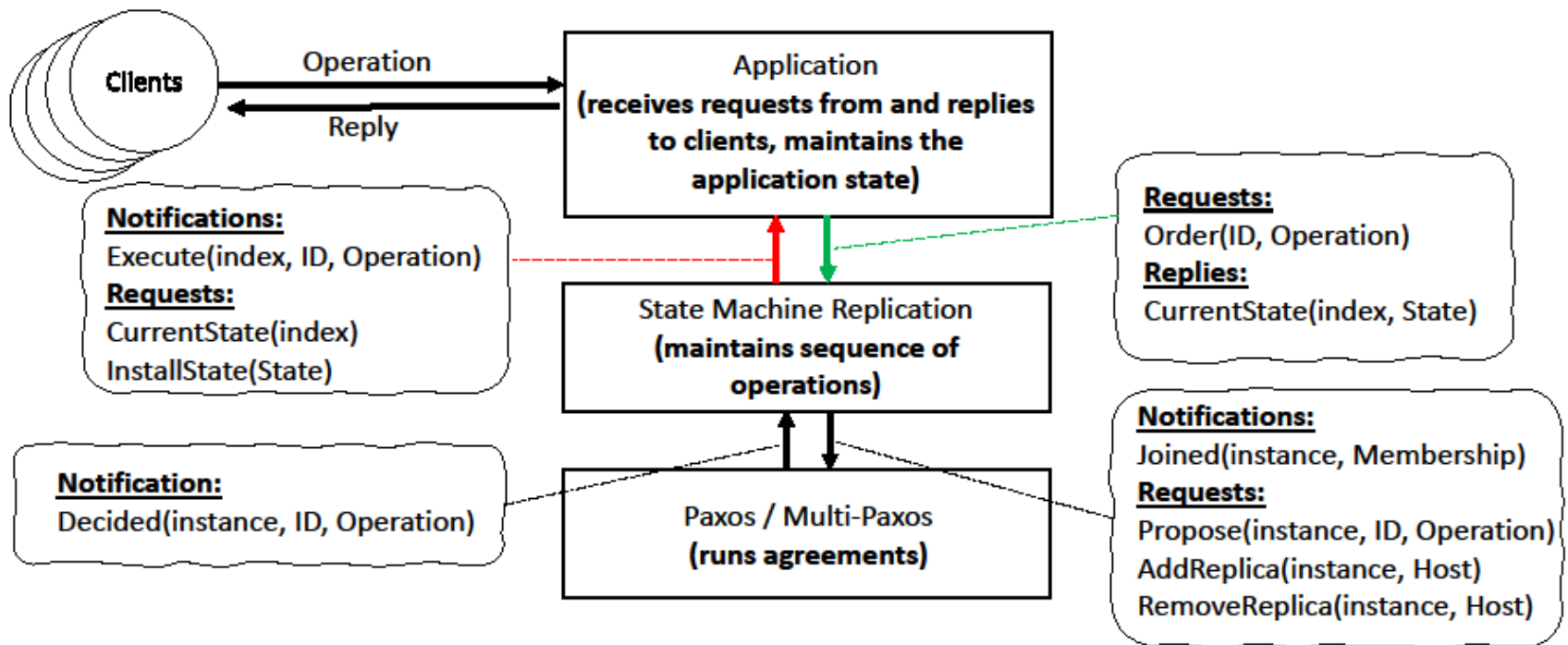
NOVA SCHOOL OF
SCIENCE & TECHNOLOGY

Based on slides from João Leitão

Second phase of the project

- **Deadline:** 29th November 2023 23:59:59
- **Link:** <https://classroom.github.com/a/PetG1iuo>
- **Submission:** Google Form link to be provided.
- **Appendix:**
 - Will be covered in future labs as theory classes catch-up
 - Extra recommendations for Paxos/Multi-Paxos implementations

Overview of Project



Paxos (Agreement protocols)

- Will be covered in theory classes from next week onwards
- The agreement protocol manages the ordering of arbitrary proposals and commands.
- You can focus on the SMR computations initially

State Machine Replication

- Maintains the sequence of commands being executed by the replicated system (the index of the command is used as the instance of the agreement protocol typically).
- Manages the membership of the system and maintains the communication channel:
 - Opens TCP connections (and attempts to reconnect on failure for a maximum number of times).
 - Investigates possible failed replicas.
 - Issues commands to remove failed replicas.
 - If joining a system already executing, acts as a client and requests to join the system directly to another process (by contacting the state machine replication protocol there).

State Machine Replication

- It receives the initial membership of the system as a parameter (if not in the membership it should request to be added).
- The reply to that request should indicate the instance in which the new replica joins the system and the application state in that instance (that must be installed by the new replica).

State Machine Replication

- While the implementation can be mostly agnostic to the underlying protocol (Paxos or Multi-Paxos), you might need to have an additional parameter to inform it and adjust the behavior:
 - Paxos has no leader, all state machine replication protocols propose commands to Paxos instances to be ordered.
 - Paxos will most likely require that when you initialize a new instance you are provided with the current membership of the system (in Multi-Paxos it might make more sense to have the membership “replicated” in the agreement protocol).
 - Multi-Paxos has a leader, only the leader proposed commands to Multi-Paxos to be ordered. Non-leader replicas must forward client operations received by them to the leader (add a notification to Multi-Paxos).
 - (be careful, in Multi-Paxos, if the leader changes you might need to resend pending client operations to the new leader).

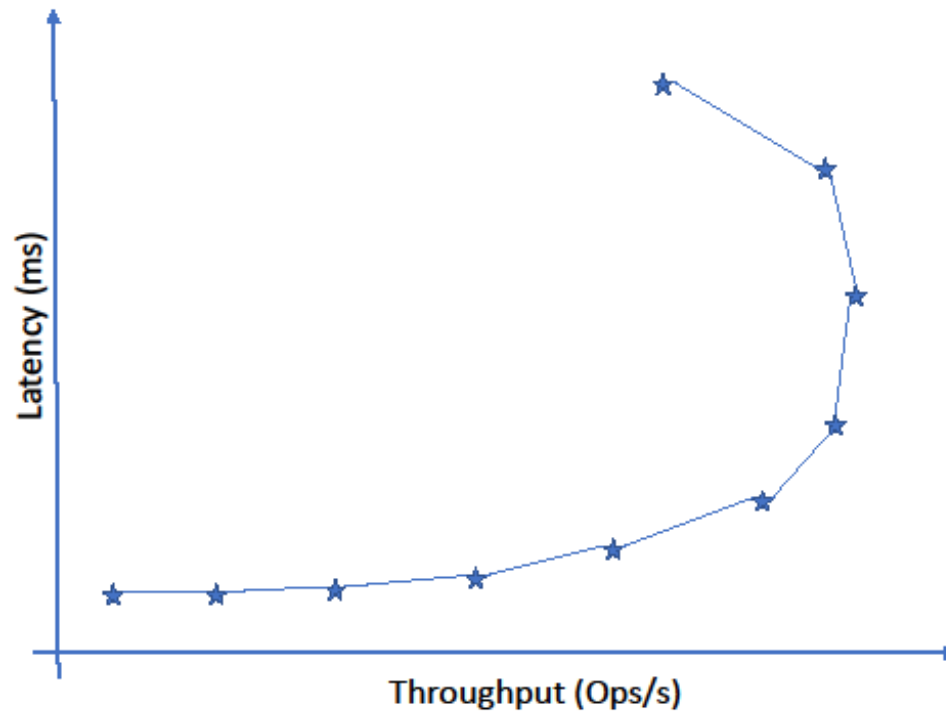
Evaluation

- Evaluate the latency and throughput of your replicated system from the perspective of the clients (emulated by YCSB).
- Evolution of these performance indicators as the number of clients increase (more clients induce more load because clients operate in a closed-loop).
- Each YCSB instance emulated multiple clients by executing multiple threads (each thread represents a client).
- Be careful to avoid saturating clients (there is a limit of threads per YCSB instance that depends on the machine you are using to run it).

Evaluation

- When using multiple YCSB instances simultaneously:
 - Average the latency observed by clients on each instance.
 - Sum the throughput of each YCSB instance.
- Increase the number of clients and plot the observed performance indicators for each system and connect increasing number of clients for each system with a line.

Example of throughput-latency plot



Evaluation

- Evaluation will not cover membership changes or (planned) failures.
- Experiments would ideally consider nodes running on different machines, and clients executing on other machines of the cluster.
 - No Docker script given, due to confusion over scripts last time
- The total number of client threads that you have to use might depend on your implementation (better implementations might need more clients to saturate the system and observe the performance wall).

Appendix: Paxos/Multi-Paxos

- You will need to implement agreement via Paxos/Multi-Paxos
- Here are some extra details that you should consider.

Agreement: Paxos / Multi-Paxos

- In Paxos you still need to have multiple instances, one for each position of the sequence of commands in the state machines.
- Suggestion:
 - Have a map in which the key is the Paxos instance, and whose value is a Java class (implemented by you) with all state variables of Paxos.
 - Paxos messages should always be tagged with the instance identifier.
 - When you process a message you should only modify the state variables of Paxos for that instance.

Agreement: Paxos / Multi-Paxos

- In Paxos there is no state that is shared between different instances.
- You have also to be careful about membership changes, since the size of quorums used in Paxos are a majority that depends on how many replicas you have.
- A replica that is not part of the system should never participate in Paxos instances (which implies that Paxos needs to have a notion of membership).

Agreement: Paxos / Multi-Paxos

- A process accepts another as a leader when it replies to a prepare message with a prepare ok.
- From that point onward, it cannot process accept messages from a different process in any instance after that one.
- A process becomes leader when it collects a majority of prepare ok messages.

Agreement: Paxos / Multi-Paxos

- In Multi-Paxos you have the notion of a leader. The leader has to be communicated to the State Machine Replication protocol, because it modifies its state.
- Notice that a replica becomes leader in a particular instance (this might need to be also communicated to the State Machine Replication Protocol).

Agreement: Paxos / Multi-Paxos

- In Multi-Paxos leader elections reflect on multiple instances, and PrepareOk messages have to carry information for all instances executed after the current one.
- You have also to be careful about membership changes, since the size of quorums used in Multi-Paxos are a majority that depends on how many replicas you have.
- A replica that is not part of the system should never participate in Multi-Paxos instances (which implies that Multi-Paxos also needs to have a notion of membership).

Details

- You might receive messages in Paxos or Multi-Paxos for an instance that you have not yet started.
- Be careful, since if you do not know yet what was decided in the previous round, you might have been excluded or a new replica might have joined the system which affects the size of quorums.
- In Paxos, you can process these messages only after you get an initial proposal from your local state machine protocol (but this might require that you have a timer in the state machine to propose special NO-OP operations that have no effect in the system).
- You can also go ahead and process the message immediately, but again you must be careful about membership changes.