

Lab 3

Eventually Consistent Blackboard

Consistency in our Blackboard so far

- Simple board (Lab 1):
 - Inconsistent
 - It could happen that different boards show messages in a different order
- With a central leader (Lab 2):
 - Consistent
 - The leader gathered all the messages and decided the order
 - Not scalable
- These are just two extreme cases in the spectrum of consistency trade offs...

Consistency Trade-offs

- Balance between strictness of consistency and efficiency/scalability
 - How “much” consistency we need, depends on the application

Loose Consistency

Strict Consistency



Efficient and
scalable

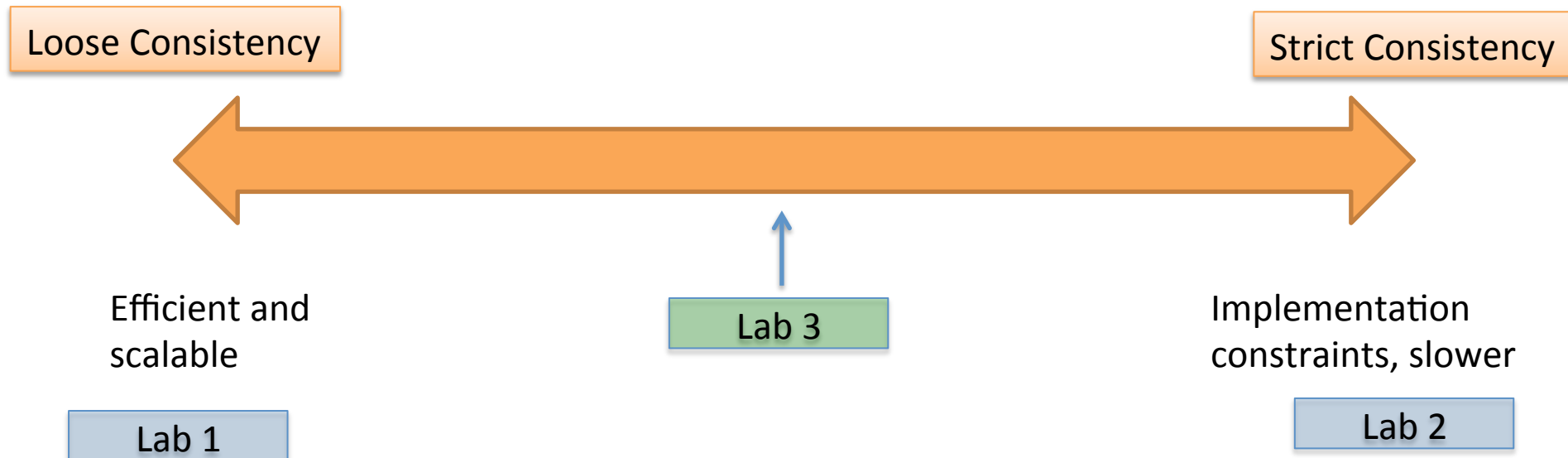
Implementation
constraints, slower

Lab 1

Lab 2

Lab 3: Eventual Consistency

- Lab 3 will be somewhere in the middle.
 - Messages can appear in different order **temporarily**.
 - Eventually, they will converge to the same order.



Eventual Consistency

- All replicas *eventually converge* to the same value
- A protocol for eventual consistency:
 - *Writes* are eventually applied in total order
 - same order on all replicas
 - lead to the same value
 - eventual consistency
 - *Reads* might not see most recent writes in total order
- Used in many applications like Google File System (GFS) and Facebook's Cassandra .
- More information on Friday's lecture.

Design Considerations

- In eventually consistent (large) data-stores:
 - Write-write conflicts are rare (but **not** in our board!!!)
 - 2 processes writing the same data is a rare case
 - It can be handled through simple mutual exclusion
 - Read-write conflicts are more frequent
 - 1 process reading a value while another process is writing a value to the same variable
 - Eventually Consistent designs should efficiently resolve such conflicts

Assumptions/Requirements

- Boards are distributed:
 - No centralized leader, no ring topology
 - You can base on Lab 1's code (if you want)
- Each post is updated to the local board, then propagated to other boards
- All boards are eventually consistent.
- Delete and modify should also be supported

Sending messages: How-to hints!

- Use logical clocks
 - Each post has a sequence number:
 - Sequence number of a new post:
the last sequence number received + 1
 - On a vessel:
 - Posts are ordered by sequence numbers
 - If two posts have the same sequence number, break ties with some rule (e.g. prioritize highest IP address)

Issues with Delete/Modify

- Consider these cases:
 - Two hosts try to delete the same message.
 - A host receives an update to delete/modify a message from the blackboard, but that message has not even arrived yet on that host .
 - It can happen in eventual consistency!
- Your algorithm needs to take care of those cases.
 - keep a history of some old updates and act accordingly when new updates arrive.

Task 1: Implement Eventual Consistency

Video/report in which you:

- Demonstrate that your Blackboard is eventually consistent
 - Inconsistent for a while
 - Then becomes consistent
 - Delete and Modify still work
- Briefly discuss pros + cons of this design
- Discuss the cost of your solution, as in Lab 2

Measurements

- You will measure the times it takes for the blackboard to reach consistency.
- Scenario:
 - N number of vessels.
 - 1000 posts on the blackboard in total (e.g. if $N=4$ send 250 posts on every vessel) at the same time
 - Measure the total time to reach consistency state:
 - the longest time among all vessels.
 - We want to see how that time changes as we add more vessels in the system.

Measurements (cont.)

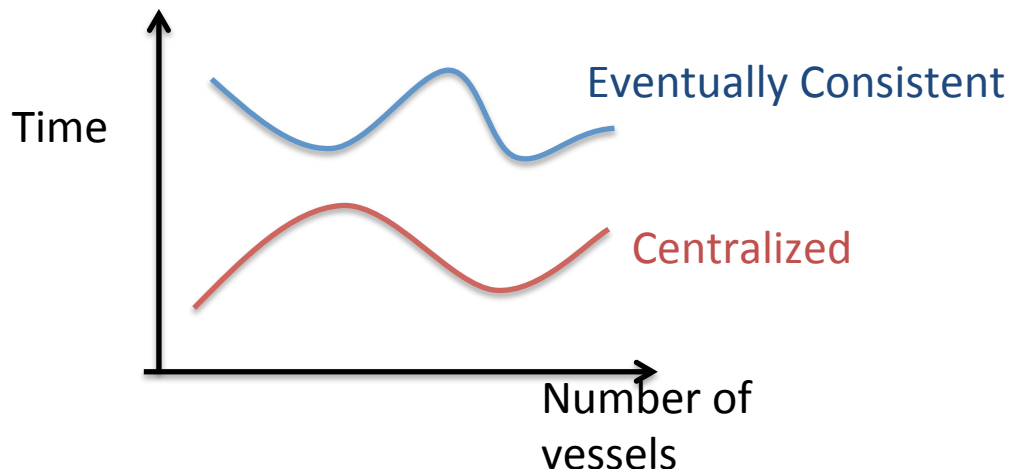
- How to measure?
 - Post messages to vessels at (almost) the same time
 - Write a shell script or a program to automatically post messages
 - Example:

```
for i in `seq 1 250`; do
    curl -d 'entry=t'${i} -X 'POST' 'http://ip:port/entries' &
done
```

This will post entry=t1, entry=t2, ..., entry=t250 to one vessel
 - Record the time for the Blackboard to reach consistency.
 - Time in Remy: using elapsed time (RemyTutorial example 1.4)

What to plot

- Do the same measurements for your eventually consistent blackboard and your centralized one from Lab 2.
- Plot the time it takes for all the vessels to get all the messages, as the number of vessels increases:



Task 2: Measurements

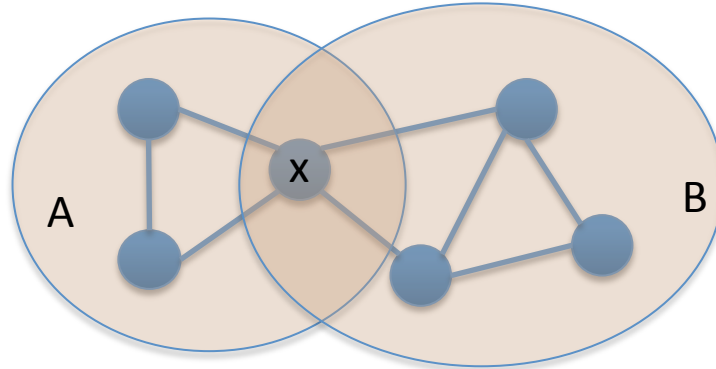
- Run the scenario for Lab2 and Lab3 and take the measurements.
- Use 2,4,6,8 and 10 vessels.
- Plot a graph: time to reach consistency as a function of the number of vessels.
- You can use any tool of choice such as openoffice, excel, google docs, matlab, matplotlib, ...
- Submit the graph, as a pdf file.
- Write a few words explaining the results you see.

Optional: Network Segmentation

- We will try to simulate a scenario where the network is segmented in two parts, then merged again.
- Completely optional, awards up to 5 bonus points.

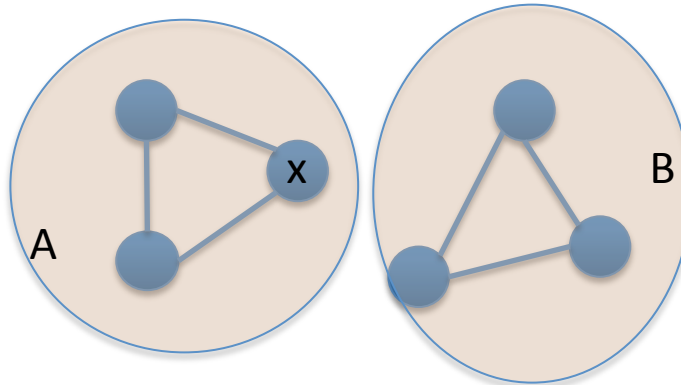
An example

- Phase 1



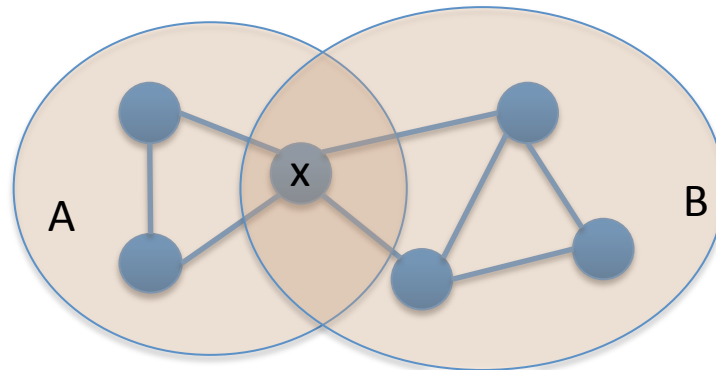
The network is a single segment. Node X communicates with both groups.

- Phase 2



Node X stops communicating with group B. Now there is consistency **only** within each group.

- Phase 3



Node X starts communicating with group B again. Groups A and B will have **inconsistent blackboards** unless they exchange their history.

The challenge

- The blackboards must reach consistency when the two disjointed groups are merged together (4 points).
- Note that the topology is **no longer “All-to-All”**
- You can choose the topology, by having it hardwired in your code, as long as you implement the segmentation in two parts..
- Measure and plot the time it takes to merge, as we add more vessels in the system(1 point).

Optional Task: Network Segmentation

- Video/report: Demonstrate that the network:
 - gets segmented, with every segment having a different view of the blackboard
 - Reaches eventual consistency, when the two segments are merged again.
- Pdf with the graph about the time it takes to merge.
- Briefly discuss how your solution works.

Summary

- Video/Report:
 - The consistent blackboard, show that it works with posts, delete and modify.
 - Optional: Network Segmentation.
- Pdf:
 - The graph comparing consistency time between Lab2 and Lab3.
 - Optional: Graph of the time it takes to merge the segments.