

# EpTO TO DEFINE

X

Master Thesis

Jocelyn Thode

Université de Fribourg

27 December 2016

**Abstract**—EpTO is one of the recently introduced total order algorithms for large-scale distributed systems and claims to provide total order and scalability at the same time. In this paper, we verify this claim by implementing EpTO [1] and evaluate its reliability in real-world conditions by comparing it to a deterministic total order algorithm named JGroups. We first compare them in a perfect environment assessing the scalability in terms of events throughput and peer numbers. We then compare them using a synthetic churn trace and a real one using the same metrics. Our results shows EpTO is performing as expected, although further testing on a bigger cluster is required to evaluate how it performs when JGroups no longer does.

## I. INTRODUCTION

Creating an algorithm providing scalability, integrity and validity, along with a total ordering for the events through all peers in a distributed system has been one of the hot topics in distributed systems research for many years. One of the recently designed algorithms on this topic is EpTO [1]. EpTO is an algorithm that claims to provide integrity, validity and total order in a large-scale distributed system. In addition, EpTO is designed to work without a global clock, removing the need to synchronize clocks precisely on every peer and thus is well-suited for dynamic large-scale distributed systems.

There are many other algorithms for disseminating and ordering events in a distributed system. There are some deterministic algorithms, which guarantee total order, agreement or other strong properties. Unfortunately, these types of algorithms are not scalable enough to be used in a large-scale distributed system [2, 3].

The problem with existing deterministic total ordering protocols is that they need some sort of agreement between all peers in the system. This causes a massive amount of network traffic and overhead on the system. Moreover, an agreement feature for an asynchronous system requires to explicitly maintain a group and have access to a failure detector [4, 5]. Due to faults and churn in large-scale distributed systems, the failure detector turns into a bottleneck for the system and thus limits the algorithm’s scalability.

As an alternative to deterministic algorithms, there are probabilistic algorithms, focusing on scalability and resiliency against failures using a probabilistic dissemination approach [6, 7, 8, 9, 10, 11, 12, 13]. These algorithms guarantee the dissemination of events in the system with high probability. This way, there is no need for failure detectors and redundant traffic, making these algorithms highly scalable. As these algorithms focus on reliability of dissemination, they often have to ignore other properties such as total ordering.

This is where EpTO shines. EpTO, by mixing these two types of algorithms, provides total order along with scalability, validity and integrity. EpTO consists of two distinct parts. The first part is probabilistic dissemination.

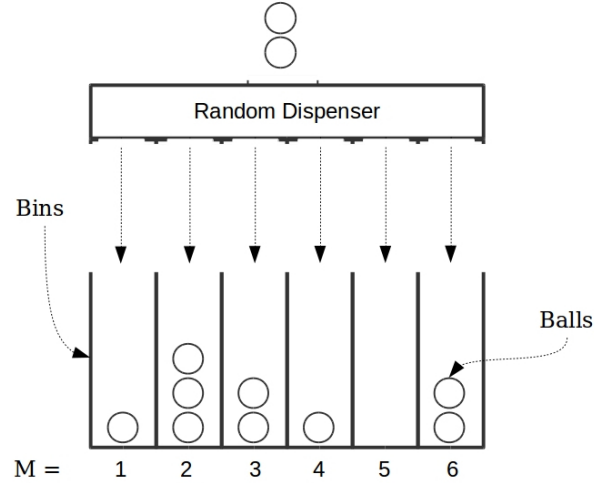


Figure 1. Balls-and-Bins [bnb].

EpTO guarantees that all peers will receive an event with arbitrarily high probability. The second part of EpTO is deterministic ordering. Once peers have received all events, they deterministically order them using the events timestamp, and in case of a tie use the broadcaster ID of the events.

To model the first part, EpTO is using a *balls-and-bins* approach [13]. The *balls-and-bins* problem is a basic probabilistic problem: consider  $n$  balls and  $m$  bins where we consequently throw balls into a bin uniformly at random and independently from other balls. In this scenario, one of the natural questions that comes to mind is: what is the minimum number of balls that should be thrown, so that every bin has at least one ball with high probability?

Using the balls-and-bins approach we model processes as bins and events as balls and calculate how many balls need to be *thrown* such that every bin contains at least one ball with arbitrarily high probability. With this approach the number of messages transmitted per process per round is logarithmic in the number of processes, therefore the number of messages sent in the system is low and uniform over all processes. Thanks to these approaches, EpTO becomes highly scalable and resilient, while still providing total order.

Until now, the creators of EpTO have only tested this algorithm in a simulated environment. In this work, we implement EpTO in pure Kotlin and introduce a benchmarking solution to show that EpTO is suitable for real-world large-scale distributed systems. We compare EpTO to the deterministic total order algorithm provided by JGroups [14] in both stable and unstable environments. EpTO uses a PSS CYCLON [15] to obtain a random view, coupled with a tracker that keeps track of dead and alive peers to obtain a first initial view. These benchmarks can easily be launched and scaled on a cluster using Docker

and Docker Swarm. Furthermore, we can inject synthetic or real churn<sup>1</sup> in the system.

The next sections are dedicated to present our benchmarking solution and to present the results obtained. section II presents the different kinds of ordering. section III defines the terms used in the paper. section IV presents EPTO and the architecture of our implementation. section V shows and explains the results obtained as well as the limitations our project faced. We then present possible future work in section VI. We finally conclude in section VII.

## II. ORDERING ALGORITHMS

**Jocelyn** ▶ *Do we still want this section?* ◀ Distributed systems and centralized systems alike, need to preserve the temporal order of events produced by concurrent processes in the system. When there are separated processes that can only communicate through messages, you cannot easily order these messages. Therefore we need ordering algorithms to overcome this problem.

We have two types of ordering algorithms [3]: the partial order algorithms and the total order algorithms.

### A. Partial Order Algorithms

Assuming  $S$  is partially ordered under  $\leq$ , then the following statements hold for all  $a, b$  and  $c$  in  $S$ :

- Reflexivity:  $a \leq a$  for all  $a \in S$ .
- Antisymmetry:  $a \leq b$  and  $b \leq a$  implies  $a = b$ .
- Transitivity:  $a \leq b$  and  $b \leq c$  implies  $a \leq c$ .

### B. Total Order Algorithms

A totally ordered set of events is a partially ordered set which satisfies one additional property:

- Totality (trichotomy law): For any  $a, b \in S$ , either  $a \leq b$  or  $b \leq a$ .

In other words, total order is an ordering that defines the exact order of every event in the system. On the other hand, partial ordering only defines the order between certain key events that depend on each other. Partial order can be useful since it is less costly to implement. However, in some cases the order of all events is important, for example when we need to know exactly which operations have been invoked in which order. We then have to use total order, otherwise we can end up in an inconsistent state.

## III. DEFINITIONS

A process or peer is defined as an actor in our system running the application that needs total order. Each process will communicate with other processes in the distributed system, exchange events, and order them together.

An event is defined as data sent at a given time by a peer. For example, we could imagine a system where each

process publishes some data to other peers. The moment where we publish data combined with the data is called an event.

A ball is a set of events bundled together and sent as one package. We use balls in EPTO to reduce network traffic.

We define EPTO scaling well as it is defined in [1]: “The number of messages transmitted per process per round is logarithmic in the number of processes, ...”. The number of rounds needed to deliver is logarithmic as well.

Since EPTO uses a probabilistic agreement instead of a deterministic agreement, there is a very low arbitrary probability a peer does not receive an event. In this instance there will be a hole in the sequence of delivered events but even in this case, the order of the delivered events will be protected by EPTO’s deterministic ordering algorithm, thus the total order property is preserved.

We write the cluster parameters as  $(n, e)$ , where  $n$  designates the total number of peers and  $e$  denominates the global event throughput per second. We tested EPTO and JGROUPS with three different cluster parameters:

- (50, 50): 50 peers with a global event throughput of 50 events per second.
- (50, 100): 50 peers with a global event throughput of 100 events per second.
- (100, 50): 100 peers with a global event throughput of 50 events per second. These parameters were also used for all synthetic churns.

## IV. EPTOTESTER

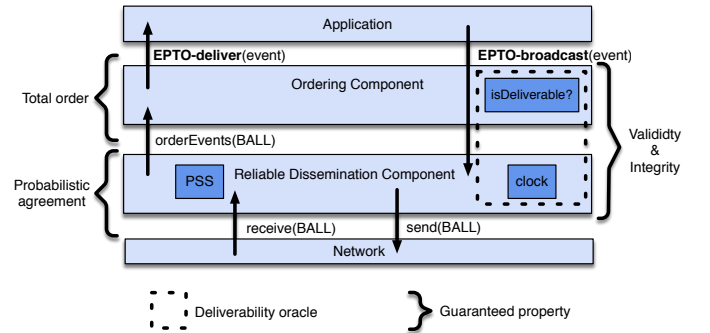


Figure 2. EPTO architecture [1].

EPTOTester is a practical implementation of EPTO designed to assess the claims made in [1]. Although the implementation was written with benchmarking in mind, the code could be adapted to be used in a real application with only minimal changes to the sources.

### A. Architecture

Figure 2 illustrates the architecture of a single replica. An application extending our Application class can EPTO-broadcast and EPTO-deliver events. Events broadcasted to the Dissemination component are sent over the network every  $\delta$  period. Every ball received from the network is unwrapped and its events are analyzed by the

<sup>1</sup>Using the Failure Trace Archive databases

Dissemination component to find out whether they need to be propagated further or not according to their Time To Live (TTL). They are then sent to the ordering component so that EPTO can determine whether to deliver these events or not and in which order. The order is based on the timestamp of the events given by the logical clock and their broadcaster ID in case of a tie. The network layer communication is done using UDP and a PSS CYCLON to obtain a random view of peers. To obtain an initial random view, we contact an independent tracker that keeps track of dead and alive nodes in the system. We want to emphasize that the tracker is not required. In practice it works well, but using a DHT is certainly a possibility.

### B. Implementation

We implement the data sent as randomly generated UUIDs. We implement our own PSS CYCLON operating on its own port. We use JGROUPS 3.6.11. When implementing JGROUPS we use TCPGOSSIP provided by the JGROUPS library instead of the traditional MULTICAST option to coordinate peers. This is not a problem as in a real WAN JGROUPS could not rely on MULTICAST. Finally, EpToTester is open-sourced on Github<sup>2</sup>.

### C. Deployment

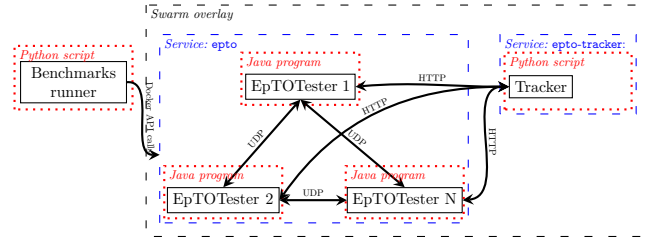


Figure 3. EpToTester architecture

To deploy EPTO effectively we need two different Docker Swarm services. One containing the tracker and one containing all EPTO replicas. Both of these services use the same network overlay to communicate. This is achieved using Docker and especially the new Docker Swarm introduced in Docker 1.12. This lets us have a unified way of deploying our benchmarks locally or remotely on vastly different clusters with minimal modifications. Every benchmark is started through a Python script available on the master mode. This script handles everything from starting benchmarks, scaling the cluster during a churn and collecting results. All EPTO parameters are customizable by using options provided through this script. Gradle is used to automate the image building in the python scripts. Finally, a script is provided to push the images to a repository accessible to the remote cluster and to push the benchmarking script on the master node. JGROUPS deployment is much the same. The entire EpToTester architecture is shown in Figure 3.

Jocelyn

►Figure 3 is inspired from the one in ErasureBench. How should I indicate this?◄

### D. Fault Injection

Our framework supports the ability to inject synthetic and real traces thanks to the work done in [16]. The synthetic churn provided by it was improved to support adding and removing nodes at the same time.

## V. EVALUATION

Jocelyn ►specify somewhere the problems encountered with JGroups◄ In this section we present our results. In the following benchmarks EPTO uses a constant  $c$  to set the probability of a hole appearing. We set  $c = 4$  in all our benchmarks so as to not experience any holes as JGROUPS does not produce holes under normal conditions. We use a  $\delta$  period of 100 ms for tests with no churn or synthetic churn and 250 ms for tests following a real trace.

We used specific settings recommended by JGROUPS for a big cluster.

Every JGROUPS test run with churn is run once killing the coordinator and once not killing it. We use the following nomenclature to differentiate both tests:

**JGroups-coord:** The coordinator is killed once.

**JGroups-nocoord:** The coordinator is purposely kept alive throughout the experiment.

**Testbed.** Jocelyn ►Most of this paragraph is copied from ErasureBench as we use the same cluster and Docker as well. Is this fine as is?◄ Our cluster consists of a cluster of 20 machines interconnected by a 1 Gb/s switched network. Each machine has an 8-core Intel Xeon CPU and 8 GB of RAM. We deploy 12 virtual machines on these hosts. Each virtual machine has access to 4 CPUs and 4 VCPUs as well as 7 GB of RAM. We use KVM as our hypervisor. Each VM uses Debian as its O/S. EpToTester is packaged as docker images. We use Docker 1.12 and its new functionality Docker Swarm to orchestrate our services. Each docker container has a memory restriction of 300 MB of RAM to be able to run them all together.

**Experiment parameters.** Every benchmark except the one following a real trace are run 10 times, each during 20 min. When there is synthetic churn, the churn starts 30 s after the benchmark and run for 17 min. The benchmarks following a real trace are run 5 times. The trace is speed-up 2x, which means we follow 2 h worth of trace in 1 h. Every benchmark run with churn is run with the (100, 50) parameters. When we have a global even throughput lower than the number of actual peers, we use a uniform random number generator to artificially reduce the global event throughput. For example, using (100, 50) parameters, each peer has a 50% probability of sending an event at each second.

We compared EPTO and JGROUPS based on their bandwidth, the time to deliver every expected events, their local dissemination stretch and finally the number of events sent.

<sup>2</sup><https://github.com/jocelynhode/EpToTester>

### A. Bandwidth

The initial bandwidth peaks observed for EPTO are due to the PSS initialization. When a node starts, it received an initial view and then quickly run the PSS four times to obtain a random enough view. In Figure 4 we observe

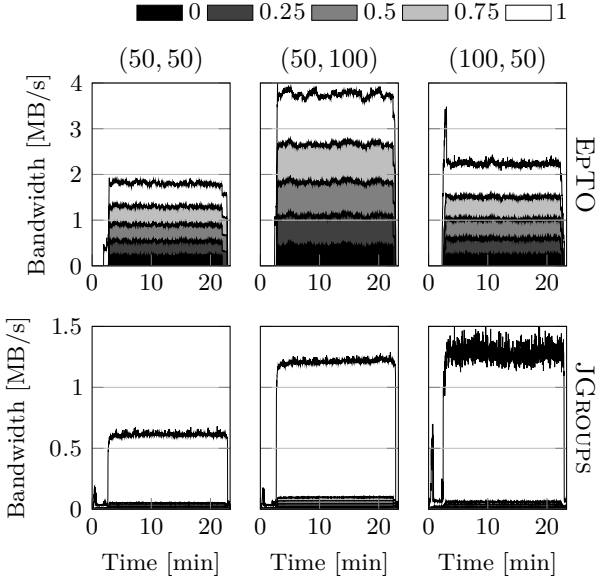


Figure 4. Bandwidth percentiles of a node during an experiment that EPTO has a worse baseline compared to JGROUPS. It uses a median bandwidth of approximately 1 MB/s for (50, 50) whereas JGROUPS uses a median bandwidth of less than 0.20 MB/s. However, in JGROUPS most of the work is done solely by the coordinator. We can clearly see this as the 100th percentile is much higher than the rest and uses approximately 0.60 MB/s.

Comparing EPTO and JGROUPS in terms of bandwidth when we increase the number of events sent per second, we can see the bandwidth doubling in both cases. In lower peers scenario such as the ones presented in Figure 4 JGROUPS is clearly at an advantage. Since EPTO has a worse baseline we will reach the maximum bandwidth possible much quicker when increasing the event throughput.

Comparing EPTO and JGROUPS in terms of bandwidth when we increase the number of peers, EPTO scales better than JGROUPS. Where JGROUPS basically has to double the bandwidth usage of the coordinator, EPTO only increases it marginally Jocelyn ►logarithmic I think as was expected◀. Thus in a scenario where we have many peers EPTO will be more efficient than JGROUPS at not reaching the maximum bandwidth.

In Figure 5 We analyze two different synthetic churns. In the first one we kill one node per minute. In the second one, we still kill one node per minute, but we immediately create a new one. For JGROUPS we ran the benchmarks once without killing the coordinator and once killing it.

We can see that the churn doesn't affect EPTO at all when there are only nodes leaving. We have small peaks

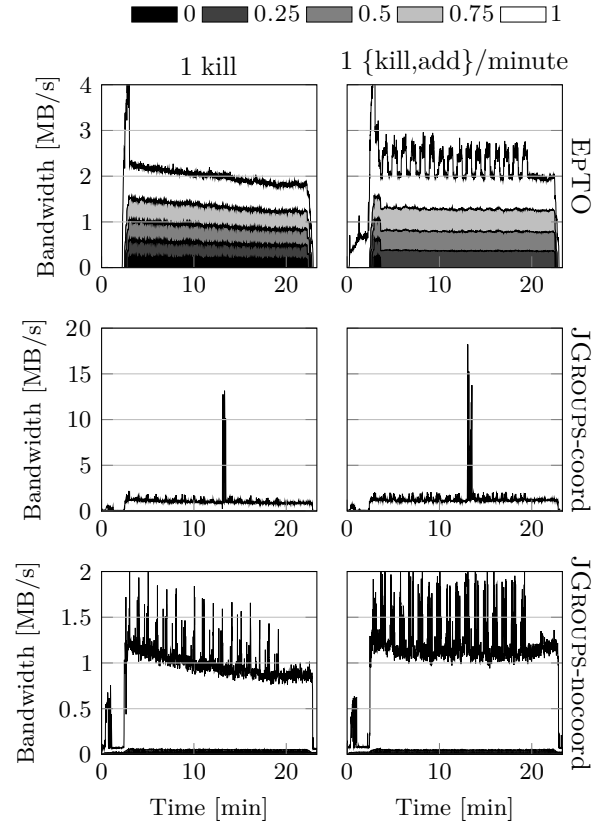


Figure 5. Bandwidth percentiles of a node during an experiment with churn

when adding a node to 3 MB/s or less. Probably due to running the PSS initialization method on top of having one more node spreading rumors in the system. This is confirmed at the end of the plot where EPTO goes back to a normal Bandwidth after stabilization.

On the other hand, when killing the coordinator in JGROUPS we can see a huge spike in bandwidth, going from 1.20 MB/s to more than 15 MB/s. This is due to how JGROUPS operates when selecting a new coordinator.

Even when not killing the coordinator, JGROUPS suffers from the churn. We can see that each time the view changes, it generates an almost 100% increase in bandwidth usage. This is due to JGROUPS having to update the view and propagate it to every peer.

### B. Total GigaBytes sent/received

In all cases, JGROUPS receives more bytes than it sends. This is due to the fact that we do not measure the bandwidth on the TCPGOSSIP, which is in charge of view changes.

In Table I, EPTO has a worse baseline than JGROUPS. This is expected as EPTO sends  $c \cdot n \cdot \log_2 n$  messages per events and JGROUPS sends at least  $n$  messages per event so we should have at least  $c \cdot \log_2 n$  more messages sent in EPTO if JGROUPS is perfect. Here we are well within this ratio.

Table I  
TOTAL GB SENT/RECEIVED IN A STABLE SYSTEM

Protocol		Cluster parameters		
		(50, 50)	(50, 100)	(100, 50)
EPTO	Receive	$10.84 \pm 0.16$	$22.31 \pm 0.39$	$26.01 \pm 0.27$
	Sending	$10.84 \pm 0.16$	$22.31 \pm 0.39$	$26.01 \pm 0.27$
JGROUPS	Receive	$0.78 \pm 0.03$	$1.45 \pm 0.01$	$1.88 \pm 0.01$
	Sending	$0.77 \pm 0.03$	$1.44 \pm 0.01$	$1.84 \pm 0.01$

Table II  
TOTAL GB SENT/RECEIVED WITH A SYNTHETIC CHURN

Protocol		Churn parameters	
		1 kill/minute	1{kill,add}/minute
EPTO	Receive	$21.00 \pm 0.24$	$26.32 \pm 0.32$
	Sending	$21.21 \pm 0.25$	$26.57 \pm 0.32$
JGROUPS-coord	Receive	$1.47 \pm 0.02$	$1.75 \pm 0.02$
	Sending	$1.43 \pm 0.02$	$1.70 \pm 0.02$
JGROUPS-nocoord	Receive	$1.45 \pm 0.01$	$1.73 \pm 0.02$
	Sending	$1.41 \pm 0.01$	$1.68 \pm 0.02$

In Table II we see that JGROUPS total bandwidth usage is smaller when there is churn. One hypothesis for this is that a JGroups replica takes a longer time to start up compared to stopping a replica. Therefore the overall benchmark has a longer time with less than 100 replicas. We also do not see a difference whether we kill the coordinator or not. This can be explained by the fact that before the detection of the faulty coordinator JGROUPS is forced to a stop for time up to 20s. The big spike afterwards compensates for this hole.

### C. Local Times

In Figure 6, JGROUPS delivers all events quicker than EPTO in all scenarios, even when churn is involved as is shown in Figure 7. However, EPTO is not too far behind. The difference between EPTO and JGROUPS is likely to be even smaller when running them in a real WAN network due to the latency. EPTO in our configuration has a  $\delta$  period of 100 ms and is thus handicapped against JGROUPS in a LAN environment, because it only increments the TTL of an event every 100 ms. We expect EPTO to outperform JGROUPS in this situation when we have a high number of peers in the system.

### D. Global Times

We computed global times as well. They are represented in Figure 8 and Figure 9. These global times are of less interest than their local counterpart as the differences in clocks between hosts can skew this measurement.

Nonetheless, here too we can see that EPTO is consistently slower than JGROUPS for the same reason as stated in subsection V-C.

### E. Local Dissemination stretch

In Figure 10, We see the percentiles of the local dissemination stretch. The local dissemination stretch is the time measurement between the sending of an event by a peer and the delivery of this event locally.

JGROUPS is usually much faster than EPTO in a perfect environment. This is expected as the benchmarks involve a small number of nodes and are performed in a LAN environment with minimal latency. The median dissemination stretch of JGROUPS is around 7 ms where as the median dissemination stretch of EPTO is around 630 ms for (100, 50). When increasing the number of peers, JGROUPS starts to have longer delivery times for some outliers. A small portion of the local dissemination stretches are really fast (1 ms or lower). This happens when the coordinator itself sends an event. Since it has to deliver it to himself the local dissemination stretch is intra-process and thus extremely fast.

In Figure 11 We see a completely different picture. When under churn, the 95th percentile of JGROUPS is at 31 ms compared to 14 ms when there is no churn. The highest percentiles are at more than 10 s. This effect is due to the coordinator dying as we clearly see that it does not happen when we do not kill it.

The median is bigger at around 9 ms, whether we kill the coordinator or not. This shows that there are some degradation in JGROUPS local dissemination stretch when under churn.

On the contrary, EPTO performs very well under churn. The median degrade a bit at 650 ms with the 99th percentile being at 1030 ms compared to 982 ms when no churn is happening.



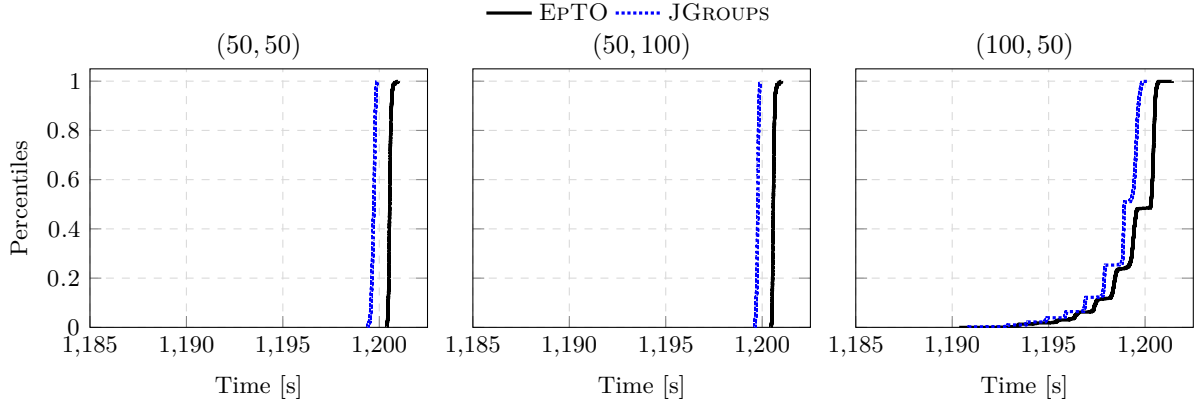


Figure 6. Local dissemination times

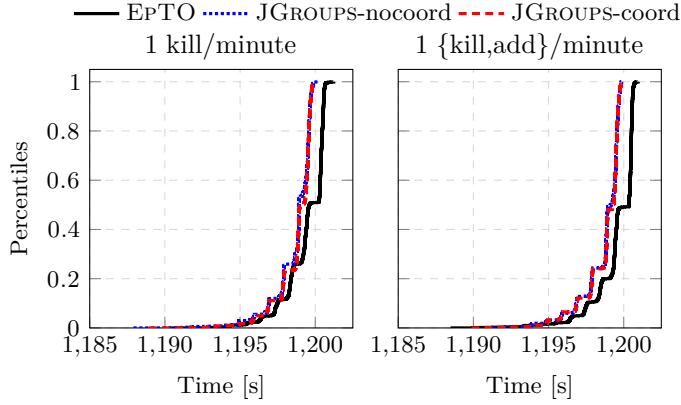


Figure 7. Local dissemination times with churn

#### F. Events sent

The variation observed in the different tests are due to the fact that we run an experiment for a period of time and that in some configurations a randomness decides whether an event is sent or not.

Table III  
TOTAL EVENTS SENT IN A STABLE ENVIRONMENT

Protocol	Cluster parameters		
	(50, 50)	(50, 100)	(100, 50)
EPTO	59993.8 $\pm$ 3.3	119898.2 $\pm$ 9.7	59913.0 $\pm$ 164.3
JGROUPS	59961.9 $\pm$ 10.9	119885.7 $\pm$ 5.0	60023.1 $\pm$ 287.1

In Table III we see that both EPTO and JGROUPS deliver the same amount of events. This is expected in a perfect environment.

In Table IV When only killing nodes, EPTO and JGROUPS again deliver the same amount of events. When killing and adding nodes, JGroups seems to deliver a smaller amount of nodes, however it does not look significant enough to draw any conclusion. Jocelyn ► *I didn't run any statistical analysis* ◀

Table IV  
TOTAL EVENTS SENT WITH A SYNTHETIC CHURN

Protocol	Cluster parameters	
	1 kill/minute	1 {kill,add}/minute
EPTO	53898.5 $\pm$ 133.9	59798.6 $\pm$ 140.1
JGROUPS-coord	53834.7 $\pm$ 175.5	59507.9 $\pm$ 240.9
JGROUPS-nocoord	53830.5 $\pm$ 200.3	59450.5 $\pm$ 175.1

#### G. Real Trace

Figure 12 confirms the results from the synthetic churn. EPTO is not affected by the churn. As before, JGROUPS is affected in the same situations as when following a synthetic churn. The spikes visible on JGROUPS-coord and JGROUPS-nocoord show the different view changes occurring.

Table V  
TOTAL GB SENT/RECEIVED

Protocol	Churn parameters	
	Real Trace	
EPTO	Receive	81.41 $\pm$ 1.08
	Sending	82.67 $\pm$ 1.08
JGROUPS-coord	Receive	5.56 $\pm$ 0.08
	Sending	5.40 $\pm$ 0.08
JGROUPS-nocoord	Receive	5.58 $\pm$ 0.05
	Sending	5.43 $\pm$ 0.05

Table V also confirms our earlier results. We run the experiment for approximately 3 times the synthetic period and we get approximately 3 times more GB used.

The local and global dissemination times shown in respectively Figure 13 and Figure 14 still show JGROUPS to outperform EPTO in this scenario. Again we want to emphasize that EPTO would more than probably outperform JGROUPS in a bigger size cluster.

In Figure 15, EPTO as worse outliers than before and so does JGROUPS-coord. Table V is interesting. We see that on average JGROUPS seems to deliver more events

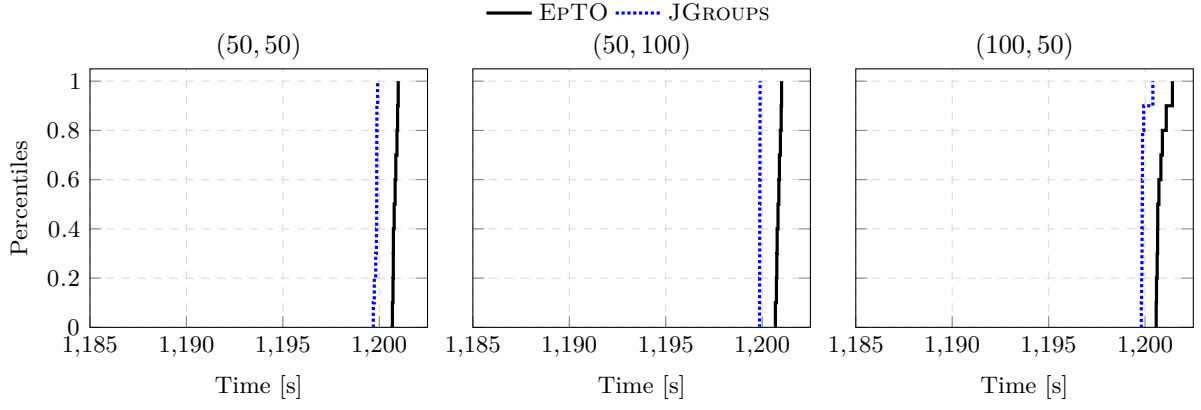


Figure 8. Global dissemination times

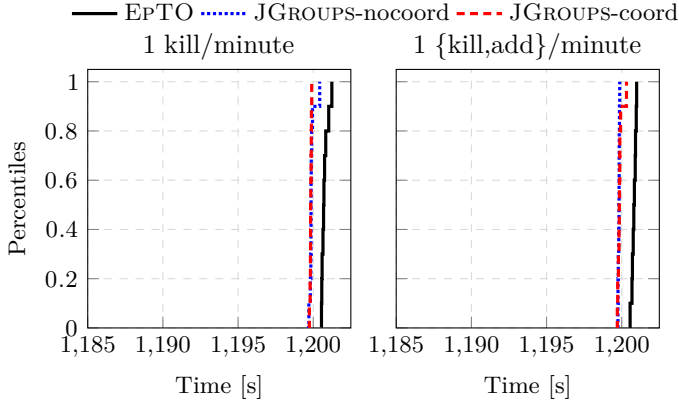


Figure 9. Global dissemination times with churn

Table VI  
TOTAL EVENTS SENT DURING A REAL TRACE

Protocol	
EPTo	165 844.2 $\pm$ 210.2
JGROUPS-coord	166 183.0 $\pm$ 1368.1
JGROUPS-nocoord	166 585.8 $\pm$ 824.9

than EPTo, but its standard deviation is way higher than that of EPTo, suggesting EPTo to be more stable and consistent.

#### H. Problems encountered

The EPTo simulation and theory is assuming we can have balls of infinite size. In practice, packet loss happens. We must therefore limit the ball size to a practical number. We use balls with a maximum size of 1432 B. We choose this number as the default MTU on Ethernet is 1500 B. We then subtract the bytes required for the IPv4 and UDP headers.

When running JGROUPS we sometimes get failed runs, where either some peers will experience holes or peer will deliver out of order events. This only happens when there is churn happening and when the coordinator dies. We

have not investigated these problems further. This points to either JGROUPS failing when there is too much churn or a bug present in the SEQUENCER implementation.

#### VI. FUTURE WORK

**Jocelyn** ► *Speak about running test on a bigger infrastructure, because the small amount of nodes does not let us see the full potential of EPTo. Then speak about the ongoing work of Push-Pull* ◀ Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

#### VII. CONCLUSION

**Jocelyn** ► *Do later* ◀ Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque.



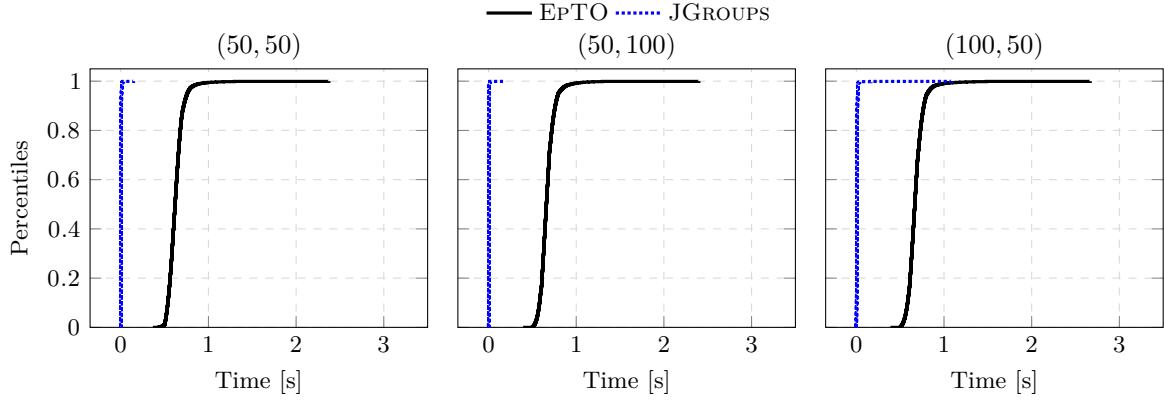


Figure 10. Local dissemination stretch

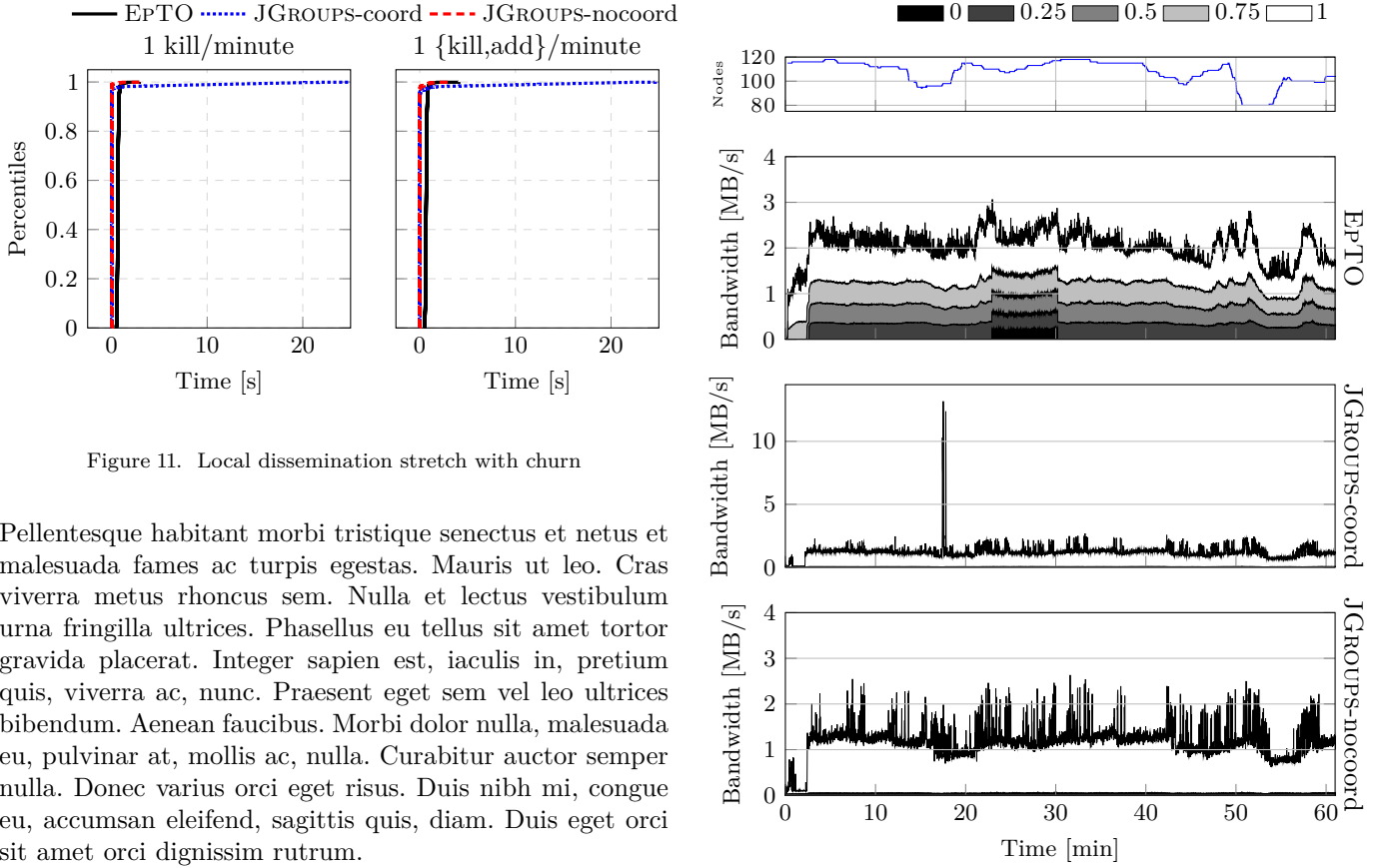


Figure 11. Local dissemination stretch with churn

Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque

Figure 12. Throughput percentiles of a node during an experiment with churn

ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus

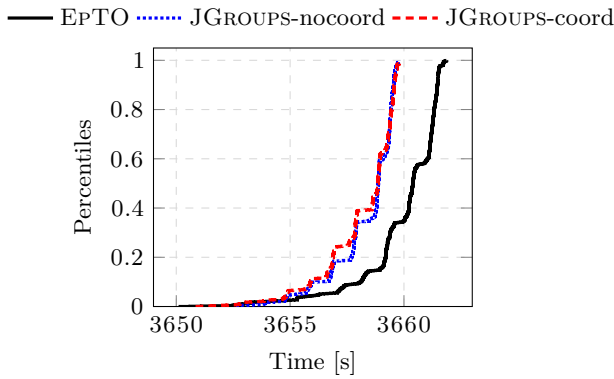


Figure 13. Local dissemination times

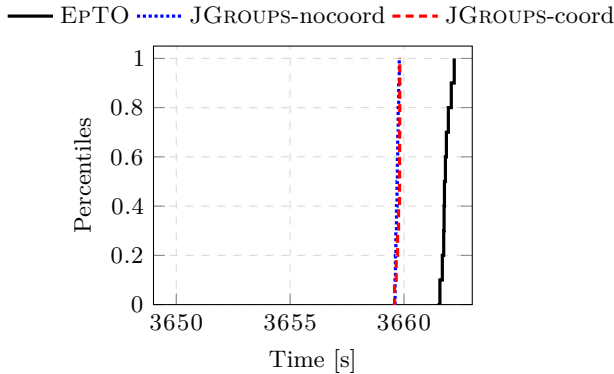


Figure 14. Global dissemination times with churn

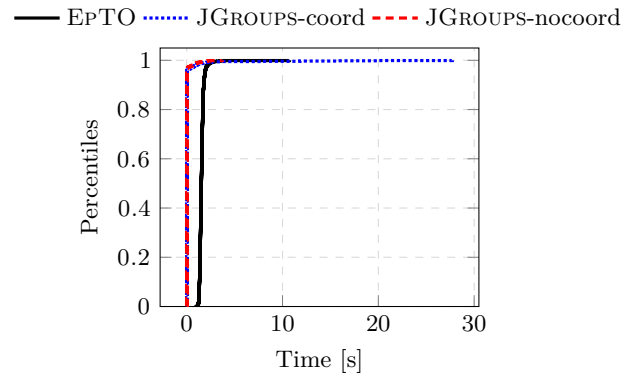


Figure 15. Local dissemination stretch with churn

pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

#### REFERENCES

- [1] M. Matos, H. Mercier, P. Felber, R. Oliveira, and J. Pereira, "EpTO: An epidemic total order algorithm for large-scale distributed systems," in *Proceedings of the 16th Annual Middleware Conference*, ACM, 2015, pp. 100–111.
- [2] X. Défago, A. Schiper, and P. Urbán, "Total order broadcast and multicast algorithms: Taxonomy and survey," *ACM Computing Surveys (CSUR)*, vol. 36, no. 4, pp. 372–421, 2004.
- [3] L. Lamport, "Time, clocks, and the ordering of events in a distributed system," *Communications of the ACM*, vol. 21, no. 7, pp. 558–565, 1978.
- [4] T. D. Chandra, V. Hadzilacos, and S. Toueg, "The weakest failure detector for solving consensus," *Journal of the ACM (JACM)*, vol. 43, no. 4, pp. 685–722, 1996.
- [5] T. D. Chandra and S. Toueg, "Unreliable failure detectors for reliable distributed systems," *Journal of the ACM (JACM)*, vol. 43, no. 2, pp. 225–267, 1996.
- [6] K. P. Birman, M. Hayden, O. Ozkasap, Z. Xiao, M. Budiu, and Y. Minsky, "Bimodal multicast," *ACM Transactions on Computer Systems (TOCS)*, vol. 17, no. 2, pp. 41–88, 1999.
- [7] N. Carvalho, J. Pereira, R. Oliveira, and L. Rodrigues, "Emergent structure in unstructured epidemic multicast," in *Dependable Systems and Networks, 2007. DSN'07. 37th Annual IEEE/IFIP International Conference on*, IEEE, 2007, pp. 481–490.
- [8] A. Demers, D. Greene, C. Hauser, W. Irish, J. Larson, S. Shenker, H. Sturgis, D. Swinehart, and D. Terry, "Epidemic algorithms for replicated database maintenance," in *Proceedings of the sixth annual ACM Symposium on Principles of distributed computing*, ACM, 1987, pp. 1–12.
- [9] P. T. Eugster, R. Guerraoui, S. B. Handurukande, P. Kouznetsov, and A.-M. Kermarrec, "Lightweight probabilistic broadcast," *ACM Transactions on Computer Systems (TOCS)*, vol. 21, no. 4, pp. 341–374, 2003.
- [10] P. Felber and F. Pedone, "Probabilistic atomic broadcast," in *Reliable Distributed Systems, 2002. Proceedings. 21st IEEE Symposium on*, IEEE, 2002, pp. 170–179.
- [11] M. Hayden and K. Birman, "Probabilistic broadcast," Cornell University, Tech. Rep., 1996.
- [12] C. Kim, J. Ahn, and C. Hwang, "Gossip based causal order broadcast algorithm," in *Computational Science and Its Applications-ICCSA 2004*, Springer, 2004, pp. 233–242.
- [13] B. Koldehofe, "Simple gossiping with balls and bins," in *Studia Informatica Universalis*, 2002, pp. 109–118.
- [14] (2016). JGroups - a toolkit for reliable messaging, [Online]. Available: <http://jgroups.org/> (visited on 12/23/2016).
- [15] S. Voulgaris, D. Gavidia, and M. van Steen, "Cyclon: Inexpensive membership management for unstructured p2p overlays," *Journal of Network and Systems Management*, vol. 13, no. 2, pp. 197–217, 2005.
- [16] S. Vaucher, H. Mercier, and V. Schiavoni, "Have a seat on the erasurebench: Easy evaluation of erasure coding libraries for distributed storage systems," in *2016 IEEE 35th Symposium on Reliable Distributed Systems Workshops (SRDSW)*, 2016, pp. 55–60.