

# From Clarity to Efficiency for Distributed Algorithms

Y. Annie Liu

Computer Science Department  
State University of New York at Stony Brook

joint work with  
Scott Stoller, Bo Lin, and Michael Gorbovitski

# Age of distributed programming

---

search engines

social networks

cloud computing

mobile computing



# Programming algorithms

---

significant advances in programming languages:

... ALGOL ... C++ ... Java ... Python ... Prolog ...

- statements: assignments, conditionals, loops
- expressions: arithmetic, Boolean, other data (sets)
- subroutines: functions, procedures (recursion)
- logic rules: predicates, deduction, though less used
- objects: keep data and do operations (organization)

that's mostly sequential and centralized.

# Concurrent programming

---

threads: multiple threads accessing shared data

threads as concurrent objects

- is the concurrent programming model in Java
- is adopted by other languages, such as Python and C#

Java made concurrent programming easier.

# Distributed programming

---

low-level or complex libraries, or restricted programming models

- sockets: C, Java, ... most widely used languages
- MPI: Fortran, C++, ... for high performance computing
- RPC: C, ... just about any language, Java RMI
- processes: Erlang, and more theoretically studied languages
- ...

study of distributed algorithms, not for building real applications

- pseudocode, English: high-level but imprecise, not executable
- formal specification languages: precise but lower-level

much less progress

# Our work: DistAlgo

---

a simple and powerful new language: very high-level, executable

- distributed processes as objects, sending messages
- yield points for control flow, handling messages
- **await** and **synchronization conditions as queries of msg history**
- high-level constructs for system configuration

compilation, optimization to generate efficient implementations:

transform expensive synchronization conditions

into efficient handlers as messages are sent and received,

by **incrementalizing queries** — **logic quantifications**

via incremental aggregate ops on sophisticated data structures

experiments with well-known distributed algorithms

including **Paxos** and **multi-Paxos** for distributed consensus

# Example: distributed mutual exclusion

---

Lamport's algorithm: developed to show logical timestamps

$n$  processes access a shared resource, need mutex, go in CS

a process that wants to enter critical section (CS)

- send requests to all
- wait for replies from all
- enter CS
- send releases to all

each process maintains a queue of requests

- order by logical timestamps
- enter CS only if it is the first on the queue
- when receiving a request, enqueue
- when receiving a release, dequeue

safety, liveness, fairness, efficiency

# How to express it

---

two extremes, and many in between

1. English: clear high-level flow; imprecise, informal
2. state machine based specs: precise; low-level control flow  
Nancy Lynch's I/O automata: 1 1/5 pages, most two-column

in between:

- Michel Raynal's pseudocode: still informal and imprecise
- Leslie Lamport's PlusCal: 90 lines (excluding comments and empty lines, by Merz)
- Robbert van Renesse's pseudocode: precise, almost high-level

lack concepts and ways for building real distributed applications  
most of these are not executable at all.

# Original description in English

---

The algorithm is then defined by the following five rules. For convenience, the actions defined by each rule are assumed to form a single event.

1. To request the resource, process  $P_i$  sends the message  $T_m : P_i \text{ requests resource}$  to every other process, and puts that message on its request queue, where  $T_m$  is the timestamp of the message.

2. When process  $P_j$  receives the message  $T_m : P_i \text{ requests resource}$ , it places it on its request queue and sends a (timestamped) acknowledgment message to  $P_i$ .

3. To release the resource, process  $P_i$  removes any  $T_m : P_i \text{ requests resource}$  message from its request queue and sends a (timestamped)  $P_i \text{ releases resource}$  message to every other process.

4. When process  $P_j$  receives a  $P_i \text{ releases resource}$  message, it removes any  $T_m : P_i \text{ requests resource}$  message from its request queue.

5. Process  $P_i$  is granted the resource when the following two conditions are satisfied: (i) There is a  $T_m : P_i \text{ requests resource}$  message in its request queue which is ordered before any other request in its queue by the relation  $<$ . (To define the relation  $<$  for messages, we identify a message with the event of sending it.) (ii)  $P_i$  has received an acknowledgment message from every other process timestamped later than  $T_m$ .

Note that conditions (i) and (ii) of rule 5 are tested locally by  $P_i$ .

# Challenges

---

each process must

- act as both  $P_i$  and  $P_j$  in interactions with all other processes
- have an order of handling all events by the 5 rules, trying to enter and exit cs while also responding to msgs from others
- keep testing the complex condition in rule 5 as events happen

actual implementations need many more details:

- create processes, let them establish channels w each other
- incorporate appropriate clocks (e.g., Lamport, vector) if needed
- guarantee the specified channel properties (e.g., reliable, FIFO)
- integrate the algorithm with the overall application

how to do all of these in an easy and modular fashion?

# Original algorithm in DistAlgo

---

```
1  def setup(s):
2      self.s = s                # set of all other processes
3      self.q = {}              # set of pending requests with logical clock

4  def cs(task):                 # for calling task() in critical section
5      -- request
6      self.c = Lamport_clock()  # rule 1
7      send ('request', c, self) to s  #
8      q.add(('request', c, self))    #
9      await each ('request',c2,p2) in q | (c2,p2) != (c,self) implies (c,self) < (c2,p2)
10     and each p2 in s | some received('ack',c2,p2) | c2 > c # rule 5
11     task()                      # critical section
12     -- release
13     q.del(('request', c, self))    # rule 3
14     send ('release', Lamport_clock(), self) to s  #

15 receive ('request', c2, p2):     # rule 2
16     q.add(('request', c2, p2))    #
17     send ('ack', Lamport_clock(), self) to p2  #

18 receive ('release', c2, p2):    # rule 4
19     q.del(('request', _, p2))    #
```

# Complete program in DistAlgo

---

```
0 class P extends Process:
... # content of the previous slide

20 def run():
    ...
21 def task(): ...
22 cs(task)
    ...

23 def main():
    ...
24 use reliable_channel
25 use fifo_channel
26 use Lamport_clock
27 ps = newprocesses(50,P)
28 for p in ps: p.setup(ps-{p})
29 for p in ps: p.start()
    ...
```

# Optimized program after incrementalization

---

```
0 class P extends Process:
1   def setup(s):
2     self.s = s                                # self.q was removed
3     self.total = size(s)                     # total number of other processes
4     self.ds = new DS()                       # aux DS for maint min of requests by other processes

6   def cs(task):
7     -- request
8     self.c = Lamport_clock()
9     self.responded = {}                     # set of responded processes
10    self.count = 0                           # count of responded processes
11    send ('request', c, self) to s           # q.add(...) was removed
12    await (ds.is_empty() or (c,self) < ds.min()) and count == total # use maintained
13    task()
14    -- release
15    send ('release', Lamport_clock(), self) to s # q.del(...) was removed

16  receive ('request', c2, p2):
17    ds.add((c2,p2))                          # add to the auxiliary data structure
18    send ('ack', Lamport_clock(), self) to p2 # q.add(...) was removed

19  receive ('ack', c2, p2):                   # new message handler
20    if c2 > c:                                # test comparison in condition 2
21      if p2 in s:                             # test membership in condition 2
22        if p2 not in responded:               # test whether responded already
23          responded.add(p2)                   # add to responded
24          count += 1                           # increment count

25  receive ('release', c2, p2):               # q.del(...) was removed
26    ds.del((c2,p2))                           # remove from the auxiliary data structure
```

# Simplified program by un-incrementalization

---

```
0 class P extends Process:
1   def setup(s):
2     self.s = s

3   def cs(task):
4     -- request
5     self.c = Lamport_clock()
6     send ('request', c, self) to s
7     await each received('request',c2,p2)
8         | not received('release',c2,p2) implies (c,self) < (c2,p2)
9         and each p2 in s | some received('ack',c2,p2) | (c,self) < (c2,p2)
9     task()
10    -- release
11    send ('release', Lamport_clock(), self) to s

12 receive ('request', c2, p2):
13   send ('ack', Lamport_clock(), self) to p2
```

# Optimized w/o queue after incrementalization

---

```
0 class P extends Process:
1   def setup(s):
2     self.s = s
3     self.q = {}           # self.q is kept as a set, need no aux DS ds
4     self.total = size(s)  # total number of other processes
5
6   def cs(task):
7     -- request
8     self.c = Lamport_clock()
9     self.responded = {}   # set of responded processes
10    self.count = 0        # count of responded processes
11    self.count2 = size(q)  # count of pending earlier requests
12    send ('request', c, self) to s          # q.add(...) was removed
13    await count2 == 0 and count == total    # use maintained results
14    task()
15    -- release
16    send ('release', Lamport_clock(), self) to s # q.del(...) was removed
17
18  receive ('request', c2, p2):
19    if (c2,p2) not in q:           # test membership in q
20      if (c2,p2) < (c,self):      # test comparison in condition 1
21        count2 +=1                # increment count2
22        q.add((c2,p2))            # q.add is kept, need no aux ds.add
23    send ('ack', Lamport_clock(), self) to p2
24
25  receive ('ack', c2, p2):        # new message handler
26    if c2 > c:                    # test comparison in condition 2
27      if p2 in s:                 # test membership in condition 2
28        if p2 not in responded:   # test whether responded already
29          responded.add(p2)        # add to responded
30          count += 1              # increment count
31
32  receive ('release', c2, p2):
33    if (c2,p2) in q:              # test membership in q
34      if (c2,p2) < (c,self):      # test comparison in condition 1
35        count2 -=1                # decrement count2
36    q.del((c2,p2))                # q.del is kept, need no aux ds.del
```

# Implementation of Lamport's algorithm

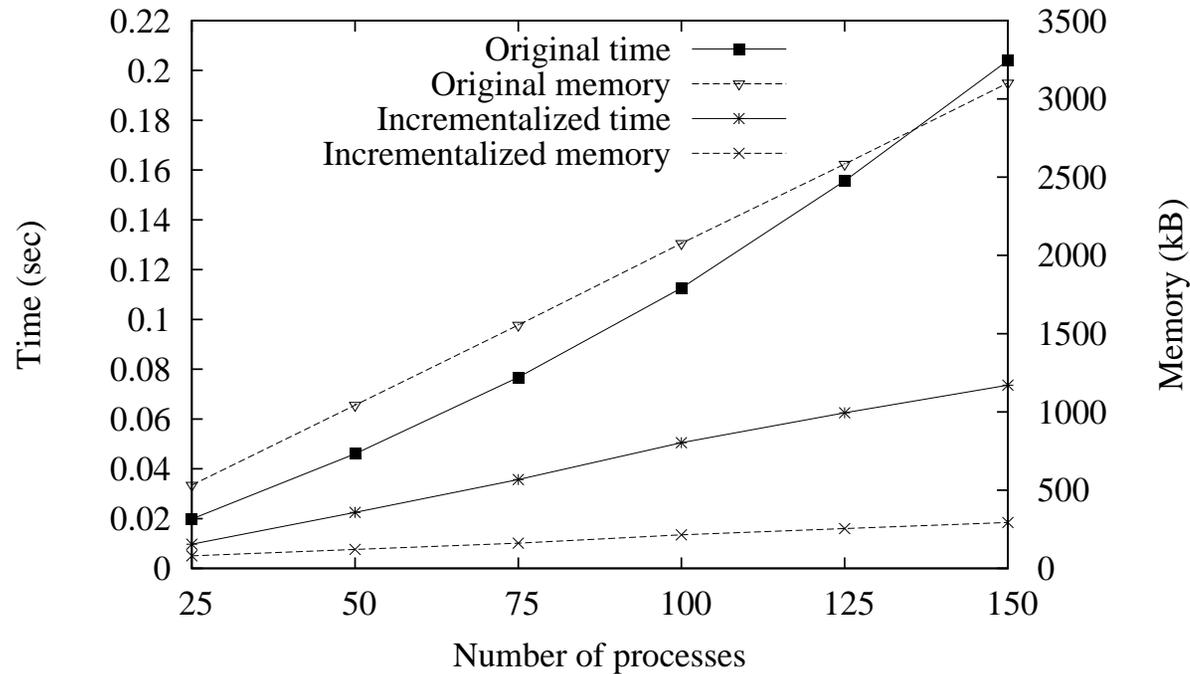
---

---

Language	Dist. programming features used	Total	Clean
C	TCP socket library	358	272
Java	TCP socket library	281	216
Python	multiprocessing package	165	122
Erlang	built-in message passing	177	99
PlusCal	single process simulation w array	134	90
DistAlgo	built-in high-level synchronization	41	32

program size in total number of lines of code,  
and number of lines excluding comments and empty lines

# Performance of generated implementation



running time and memory usage for Lamport's algorithm:

CPU time for each proc to complete a call to `cs(task)`, including time spent handling messages from other processes, averaged over processes and over runs of 30 calls each;

raw size of all data structures created, measured using `pympler`

# Program size for well-known algorithms

---

---

Algorithm	DistAlgo	PlusCal	IOA	Overlog	Bloom
La mutex	31	90	64		
La mutex2	32				
RA mutex	32				
RA token	36				
SK token	39				
CR leader	25		41		
HS leader	41				
2P commit	32	68			85
DS crash	24				
La Paxos	44	83	145	230	157
CL Paxos	72	166			
vR Paxos	132				

number of lines excluding comments and empty lines,  
compared with specifications written by others in other languages

# Compilation and optimization

Algorithm	Compilation time (ms)	DistAlgo size	Generated Python size
La mutex	4.45	31	951
La mutex2	4.71	32	955
RA mutex	3.92	32	949
RA token	4.26	36	952
SK token	4.72	39	954
CR leader	3.15	25	939
HS leader	5.87	41	957
2P commit	5.91	33	978
DS crash	3.42	24	940
La Paxos	9.12	44	1003
CL Paxos	13.06	72	1044
vR Paxos	21.60	132	1099
Incrementalized	Compilation time (ms)	DistAlgo size	Generated Python size
La mutex	4.99	43	960
2P commit	6.82	55	1001
La Paxos	7.61	59	999
CL Paxos	12.35	81	1024

# Grad and undergrad projects in DistAlgo

---

---

Project	Description	Notes
Leader	ring, randomized; arbitrary net	3 algorithms
Narada	overlay multicast system	
Chord	distributed hash table (DHT)	
Kademlia	DHT	
Pastry	DHT	
Tapestry	DHT	
HDFS	Hadoop distributed file system	part
UpRight	cluster services	part
AODV	wireless mesh network routing	python
OLSR	optimized link state routing	python

part: omitted replication, but done in our impl. of vR Paxos

python: in Python, but knew it would be easier in DistAlgo

each is about 300-600 lines of code, took about half a semester.

# Summary and conclusion

---

programming distributed algorithms: **much to do**

need both **clarity** (both high-level and precise) and **efficiency**

**DistAlgo**: a new language, simple, powerful

distributed processes and sending messages

yield points and handling messages

**await** and **synchronization conditions** as queries of msg history

high-level constructs for system configuration

**new optimization**: powerful

transform expensive synchronization conditions

into efficient handlers as messages are sent and received,

**by incrementalizing queries** — **logic quantifications**

via incremental aggregate ops on sophisticated data structures

experiments with well-known distributed algorithms

including **Paxos** and **multi-Paxos** for distributed consensus

Thanks!

# Example: two-phase commit

---

a coordinator and a set of cohorts try to commit a transaction

phase 1:

- coordinator sends a prepare to all cohorts.
- each cohort replies with a ready vote if it is prepared to commit, or else replies with an abort vote and aborts.

phase 2:

- if coordinator receives a ready vote from all cohorts, it sends a commit to all cohorts; each cohort commits and sends a done to coordinator; coordinator completes when receives a done from all cohorts.
- if coordinator receives an abort vote from any cohort, it sends an abort to all cohorts who sent a ready vote; each cohort who sent a ready vote aborts.

agreement, validity, weak termination,  $4n-4$  msgs

# How to express it

---

two extremes, and many in between

1. English: clear high-level flow; imprecise, informal
2. state machine based specs: precise; low-level control flow  
Nancy Lynch's I/O automata: (book p183-184, but  $2n-2$  msgs)

in between:

- Michel Raynal's pseudocode: still informal and imprecise
- Leslie Lamport's PlusCal: still complex  
(P2TwoPhase, 68 lines excluding comments and empty lines)
- Robbert van Renesse's pseudocode: precise, almost high-level

lack concepts and ways for building real distributed applications  
most of these are not executable at all.

# Original description in English

---

## Phase 1:

Summary of the protocol [KBL06 DB and TP]

1. The coordinator sends a *prepare message* to all cohorts.
2. Each cohort waits until it receives a *prepare message* from the coordinator. If it is prepared to commit, it forces a prepared record to its log, enters a state in which it cannot be aborted by its local control, and sends “ready” in the *vote message* to the coordinator.

If it cannot commit, it appends an abort record to its log. Or it might already have aborted. In either case, it sends “aborting” in the *vote message* to the coordinator, rolls back any changes the subtransaction has made to the database, release the subtransaction’s locks, and terminates its participation in the protocol.

## Phase 2:

1. The coordinator waits until it receives votes from all cohorts. If it receives at least one “aborting” vote, it decides to abort, sends an *abort message* to all cohorts that voted “ready”, deallocates the transaction record in volatile memory, and terminates its participation in the protocol.

If all votes are “ready”, the coordinator decides to commit (and stores that fact in the transaction record), forces a commit record (which includes a copy of the transaction record) to its log, and sends a *commit message* to each cohort.

2. Each cohort that voted “ready” waits to receive a message from the coordinator. If a cohort receives an *abort message*, it rolls back any changes the subtransaction has made to the database, appends an abort record to its log, releases the subtransaction’s locks, and terminates its participation in the protocol.

If the cohort received a *commit message*, it forces a commit record to its log, releases all locks, sends a *done message* to the coordinator, and terminates its participation in the protocol.

3. If the coordinator committed the transaction, it waits until it receives *done message* from all cohorts. Then it appends a completion record to its log, deletes the transaction record from volatile memory, and terminates its participation in the protocol.

# Original algorithm in DistAlgo

---

```
1 class Coordinator extends Process:
2   def setup(tid, cohorts): pass # transaction id and cohorts
3   def run():
4     send ('prepare',tid) to cohorts
5     await each c in cohorts | received('vote',_,tid) from c
6     if each c in cohorts | received('vote','ready',tid) from c:
7       send ('commit',tid) to cohorts
8       await each c in cohorts | received('done',tid) from c
9       print('complete'+tid)
10    else:
11      s = {c in cohorts | received('vote','ready',tid) from c}
12      send ('abort',tid) to s
13      print('terminate'+tid)

14 class Cohort extends Process:
15   def setup(f): pass # failure rate
16   def run():
17     await(False)
18   receive ('prepare',tid) from c:
19     if prepared(tid):
20       send ('vote','ready,tid) to c      # await commit or abort here?
21     else:
22       send ('vote','abort',tid) to c
23       abort(tid)
24   receive ('commit',tid) from c:
25     commit(tid)
26     send ('done',tid) to c
27   receive ('abort',tid):
28     abort(tid)
29   def prepared(tid): return randint(0,100) > f
30   def abort(tid): print('abort'+tid)
31   def commit(tid): print('commit'+tid)
```

# Complete program in DistAlgo

---

```
0  from random import randint

... # content of the previous slide

32 def main():
33     cs = createprocs(Cohort,25,(10))      # create 25 cohorts
34     c = createprocs(Coordinator,1,(0,cs)) # create 1 coordinator
35     startprocs(cs)                       # start cohorts
36     startprocs(c)                        # start coordinator
```

# Optimized after incrementalization (part 1)

---

```
1 class Coordinator extends Process:
2   def setup(tid, cohorts):
3     ncohorts = size(cohorts) # number of cohorts
4     svoted = {}             # set of voted cohorts
5     nvoted = 0              # number of voted cohorts
6     sready = {}            # set of ready cohorts
7     nready = 0              # number of ready cohorts
8     sdone = {}             # set of done cohorts
9     ndone = 0              # number of done cohorts
10
11  def run():
12    send ('prepare',tid) to cohorts
13    await nvoted == ncohorts # replaced universal quantification
14    if nready == ncohorts:   # replaced universal quantification
15      send ('commit',tid) to cohorts
16      await ndone == ncohorts # replaced universal quantification
17      print('complete'+tid)
18    else:
19      s = sready             # replaced set query
20      send ('abort',tid) to s
21      print('terminate'+tid)
```

# Optimized after incrementalization (part 2)

---

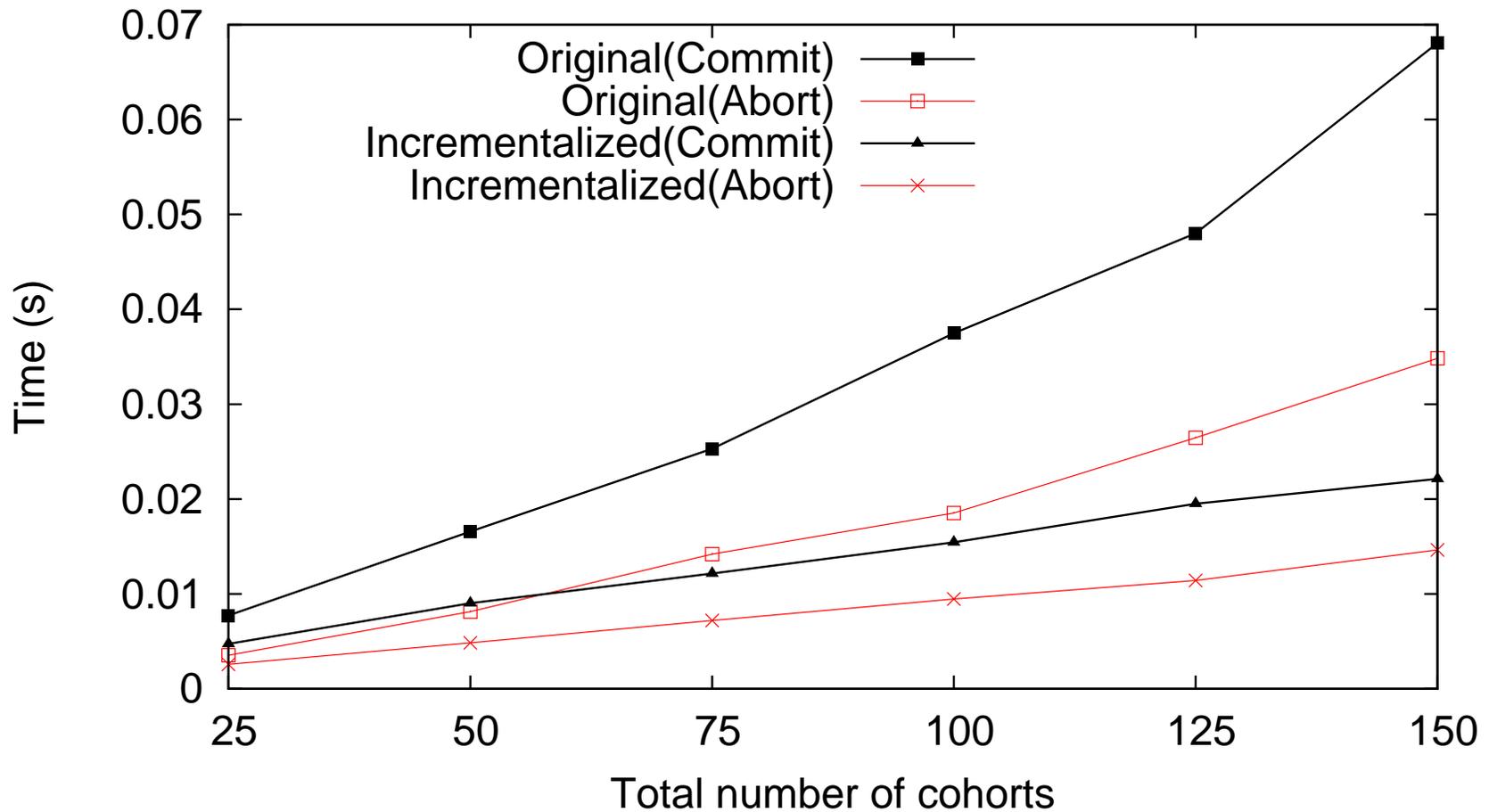
```
    # new message handler
21  receive ('vote',v,tid) from c:
22      if c in cohorts:
23          if c not in svoted:
24              svoted.add(c)
25              nvoted += 1
26          if v == 'ready':
27              if c not in sready:
28                  sready.add(c)
29                  nready += 1

    # new message handler
30  receive ('done',tid) from c:
31      if c in cohorts:
32          if c not in sdone:
33              sdone.add(c)
34              ndone += 1

35  class Cohort extends Process:
52      ... # no change

53  def main():
57      ... # no change
```

# Performance of generated implementation



for two-phase commit, for failure rates of 0 (Commit) and 100 (Abort), averaged over 50 rounds and 15 independent runs.

# DistAlgo: distributed procs and sending msgs

## process definition

```
class P extends Process: class_body
```

defines class P of process objects, with private fields

## process creation

```
new P(...,s)                                newprocesses(n,P)
```

creates a new proc of class P on site s, returns the proc

## sending messages

```
send m to p                                  send ms to ps
```

sends message m to process p

tuples or objects for messages;

first component or class indicates the kind of the message

# DistAlgo: control flows and receiving msgs

---

## label for statement

-- l

defines program point l where the control flow can yield to handling of certain messages and resume afterwards

## handling messages received

receive m from p at l: stmt                      receive ms at ls

allows handling of ms at ls; default is at all labels

## synchronization

await bexp timeout time

awaits value of bexp to be true, or time seconds have passed

can query sequences of messages received and sent

# Optimization: incrementalization

---

1. identify all expensive queries
2. determine all updates to the parameters of these queries
3. transform queries into efficient incremental comp. at updates

new: systematic handling of

1. quantifications for synchronization as expensive queries
2. updates caused by sending, receiving, and handling of msgs in the same way as other updates in the program

transform expensive synchronization conditions into efficient tests and incremental updates as msgs are sent and received

sequences received and sent will be removed

only values needed for incremental computation of synchronization conditions will be stored and incrementally updated

# Incrementalization: example

---

expensive computation of synchronization condition:

each ('request',c2,p2) in q | (c2,p2) != (c,self) implies (c,self) < (c2,p2)  
and each p2 in s | some received('ack',c2,p2) | (c,self) < (c2,p2)

all updates to variable used by expensive computation:

```
4     self.q = {}
9     q.add(('request', c, self))
14    q.del(('request', c, self))
17    q.add(('request', c2, p2))
20    q.del(('request', _, p2))

7     self.c = Lamport_clock()

3     self.s = s

new   received.add(('ack', c2, p2))
```

incrementalize: how?

# Incrementalization: at a high-level

---

aggregates:

```
( {(c2,p2) : ('request',c2,p2) in q | (c2,p2) != (c,self)} == {} or  
  (c,self) < min({(c2,p2) :  
    ('request',c2,p2) in q | (c2,p2) != (c,self)}) )
```

and

```
size({p2 in s : ('ack',c2,p2) in received | c2 > c}) == size(s)
```

- introduce variables to store values of 4 aggregates
- transform the aggregates to use introduced variables
- incrementally maintain stored values at each update

synchronization condition will become:

```
(ds.is_empty() or (c,self) < ds.min()) and count == total
```

# Implementation

---

**syntax and parsing:** build on python parser.

provide two options for yield points and message handlers:

- modify python ASDL to take succinct syntax
- use unmodified python exploiting existing python syntax

**compilation:**

- generate python code, using multi-processing and socket
- could generate C, Java, and Erlang too
- plan to generate PlusCal and others for verification

**optimization:**

- transform quantifications to aggregates & comprehensions
- build an interface to incrementalizer
- use InvTS for incrementalization