

CPSC 416 - Distributed program analysis and Invariant Inference

March 28, 2018

Lecturer Stewart Grant

The problem: distributed systems are complex!

Difficulty	Factors
Understandability	<ul style="list-style-type: none">● Concurrency● Decentralized State● No centralized clock
Debug	<ul style="list-style-type: none">● Nondeterminism● Changing environment (network)
Test	<ul style="list-style-type: none">● State space is massive● Exhaustive testing is impractical● Configuration space is even larger● Cost of large scale deployments

How do you know that a distributed system works?

- Logging
 - Open log in emacs/vi, brew coffee, get comfortable!
 - Maybe use ShiViz on the logs if you are debugging protocol issues
- Test as much as you can (Unit/Integration/**Stress**)
- Mathematically prove correct?
 - (No one does that really)

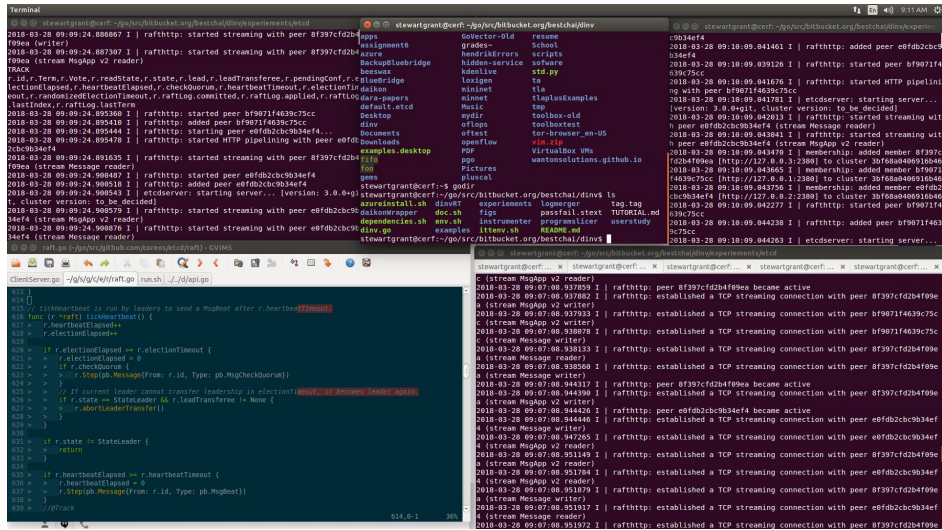


Figure (1) A typical distributed systems developers desktop [my desktop]

What other techniques are available?

Static analysis:

Analyze a programs source code, without running the program.

- Type Checker
- Linter
- Symbolic Execution

Complete but **over-approximate** and **expensive**.

Dynamic analysis:

Analyze a programs behavior as it runs, usually by logging.

- Testing
- Profiling
- Deadlock detection
- Memory profiling (valgrind)

Incomplete but **scalable**

Today's lecture

- Program analysis background
 - Static analysis
 - Dynamic analysis
- Dinv's tool and analyses
 - Data invariants
 - Static: program slicing
 - Dynamic: distributed lattice construction
- Answer any Dinv questions you might have

Program Properties: Data invariants

- An invariant is a property that holds on data at all times
- A data invariant can hold between 1 or more variables
- Data invariants are type dependent

Knowledge of a programs invariants is important for understanding if it is correct or faulty.

Example Program:

```
var sum = 0
for i:=0;i<TOTAL;i++){
    sum += i
}
```

Example Invariants:

```
i < TOTAL // loop invariant
i >= 0
sum >= i
```

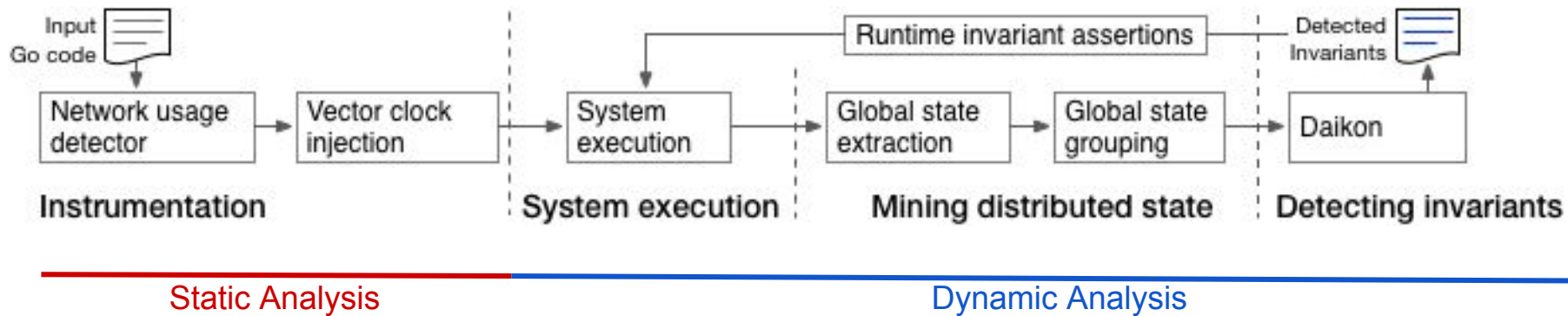
What is a Distributed Data Invariant?

- Distributed data invariants hold across 1 or more nodes in a distributed system
- Some hold globally at all times
- Some are protocol specific

Ex) Distributed Key Value Store Invariant. No two nodes serve the same keys.

$$\forall \text{Nodes } i, j, \text{Keys}_i \neq \text{Keys}_j$$

Dinv Overview

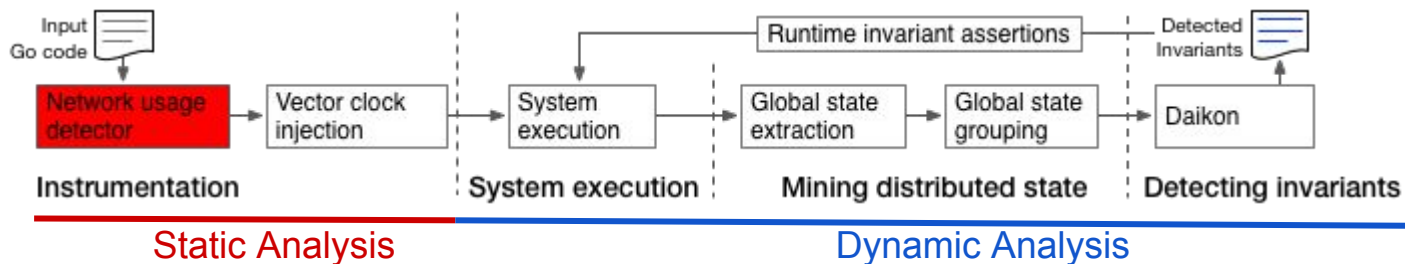


1. Distributed Invariant Inference Challenges

- a. **What state should be logged and when?**
- b. How to infer distributed invariants from logged state?
- c. How to enforce distributed invariants?

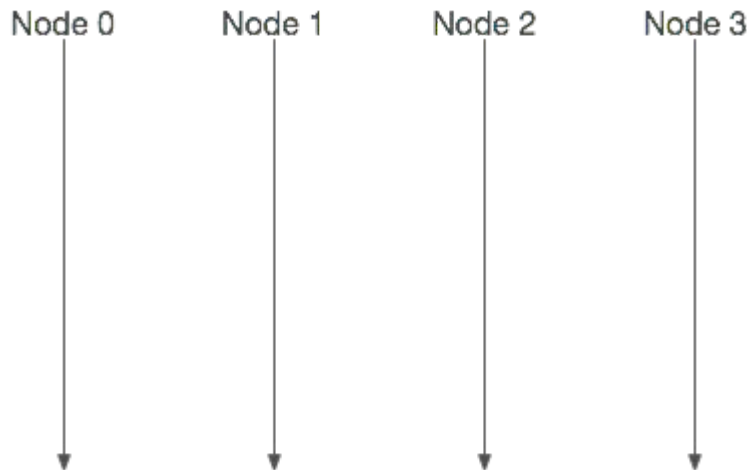
What variables should be logged?

- **Massive** variable state space
- Exponentially larger invariant state space
- Arbitrary distributed invariants be minimized



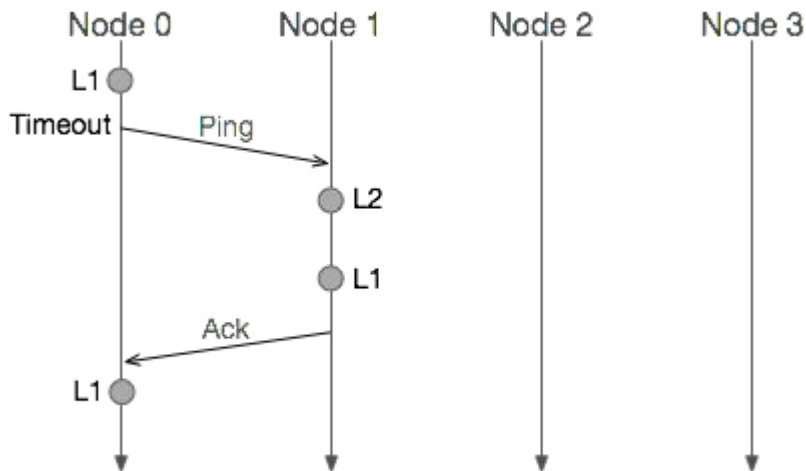


Example Code: Serf



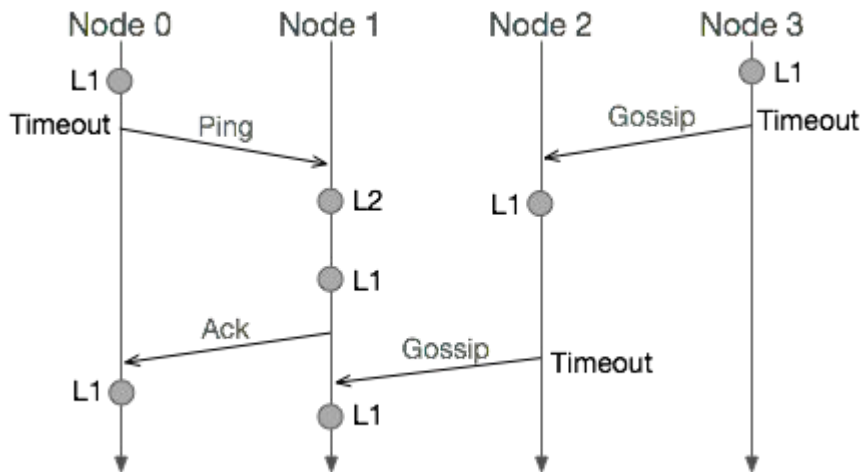
```
1 func (s serfNode) serf(conn UDPConnection) {  
2   for true {  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19   }}}}
```

Example Code: Serf



```
1 func (s serfNode) serf(conn UDPConnection) {
2   for true {
3     msg := conn.Read()
4     switch msg.Type {
5     case PING:
6       conn.WriteToUDP("ACK", msg.Sender)
7       break
8
9
10    }
11    timeout := s.CheckForTimeouts()
12    switch timeout.Type {
13    case PING:
14      conn.WriteToUDP("PING", timeout.Node)
15      break
16
17
18
19  }}}
```

Example Code: Serf



```

1 func (s serfNode) serf(conn UDPConnection) {
2   for true {
3     msg := conn.Read()
4     switch msg.Type {
5     case PING:
6       conn.WriteToUDP("ACK", msg.Sender)
7       break
8     case GOSSIP:
9       s.Events = append(s.Events, msg.Event)
10    }
11    timeout := s.CheckForTimeouts()
12    switch timeout.Type {
13    case PING:
14      conn.WriteToUDP("PING", timeout.Node)
15      break
16    case GOSSIP:
17      gossip(s.Events)
18      break
19    }
20  }
21 }

```

What state should be logged and when?

Insight: *Important distributed state must have dataflow to and from the network.*

Technique: Program slicing [Ottenstein 84]

```
1 func (s serfNode) serf(conn UDPConnection) {
2   for true {
3     msg := conn.Read()
4     switch msg.Type {
5     case PING:
6       //@dump L2
7       conn.WriteToUDP("ACK", msg.Sender)
8       break
9     case GOSSIP:
10      s.Events = append(s.Events, msg.Event)
11    }
12    dinv.Dump("L1",msg.Type,msg.Sender,msg.Event,s.Events) L1
13    timeout := s.CheckForTimeouts()
14    switch timeout.Type {
15    case PING:
16      conn.WriteToUDP("PING",timeout.Node)
17      break
18    case GOSSIP:
19      gossip(s.Events)
20      break
21    }}}
```

What state should be logged and when?

Insight: *Important distributed state must have dataflow to and from the network.*

Technique: Program slicing [Ottenstein 84]

- Transitively track assignments to variables
- A slice is the complete set of statements over which marked data flows

```
1 func (s serfNode) serf(conn UDPConnection) {
2   for true {
3     msg := conn.Read()
4     switch msg.Type {
5     case PING:
6       //@dump L2
7       conn.WriteToUDP("ACK", msg.Sender)
8       break
9     case GOSSIP:
10      s.Events = append(s.Events, msg.Event)
11    }
12    dinv.Dump("L1",msg.Type,msg.Sender,msg.Event,s.Events) L1
13    timeout := s.CheckForTimeouts()
14    switch timeout.Type {
15    case PING:
16      conn.WriteToUDP("PING",timeout.Node)
17      break
18    case GOSSIP:
19      gossip(s.Events)
20      break
21  }}}}
```

What state should be logged and when?

Insight: *Important distributed state must have dataflow to and from the network.*

Technique: Program slicing [Ottensstein 84]

- Transitively track assignments to variables
- A slice is the complete set of statements over which marked data flows

```
1 func (s serfNode) serf(conn UDPConnection) {
2   for true {
3     msg := conn.Read()
4     switch msg.Type {
5     case PING:
6       //@dump L2
7       conn.WriteToUDP("ACK", msg.Sender)
8       break
9     case GOSSIP:
10      s.Events = append(s.Events, msg.Event)
11    }
12    dinv.Dump("L1",msg.Type,msg.Sender,msg.Event,s.Events) L1
13    timeout := s.CheckForTimeouts()
14    switch timeout.Type {
15    case PING:
16      conn.WriteToUDP("PING",timeout.Node)
17      break
18    case GOSSIP:
19      gossip(s.Events)
20      break
21  }}}}
```

Q: What do you think the challenges of dataflow analysis are?

Some Answers:

- Aliasing (when one bit of data can be confused with many)
- Pointer analysis
- Interprocedural flow
- Thread interleaving
- Distributed dataflow

Where should state be logged?

Location and frequency of logging correspond to invariant accuracy

Instrumentation Strategy	Location choice	Variable Choice
Function entrances/exits	Auto	Auto
Network calls	Auto	Auto
User-defined annotations	Manual	Manual or auto

How to Instrument with Dinv

Source: `$REPOLOCATION/dinv/examples/helloDinv/ClientServer.go`

Pre Instrumentation:

Two annotations:

@Track & @Dump

Track Recommended
(Reduces Output Size)

```
50 func client(listen, send string) {
51 > // sending UDP packet to specified address and port
52 > conn := setupConnection(SERVERPORT, CLIENTPORT)
53 > for i := 0; i < MESSAGES; i++ {
54 > > //@track
55 > > outgoingMessage := i
56 > > outBuf := dinvRT.Pack(outgoingMessage)
57 > > _, errWrite := conn.Write(outBuf)
58 > > printErr(errWrite)
59 > > var inBuf [512]byte
60 > > var incommingMessage int
61 > > n, errRead := conn.Read(inBuf[0:])
62 > > printErr(errRead)
63 > > dinvRT.Unpack(inBuf[0:n], &incommingMessage)
64 > > incommingMessage = n - n + incommingMessage
65 > > //fmt.Printf("GOT BACK : %d\n", incommingMessage)
66 > > time.Sleep(1)
67 > }
68 > done <- 1
69 }
```

How to Instrument with Dinv

Source: `$REPOLOCATION/dinv/examples/helloDinv/ClientServer.go`

Instrumentation Command:

```
dinv -i  
-file=ClientServer.go
```

The resulting source code is
Instrumented to log variables.

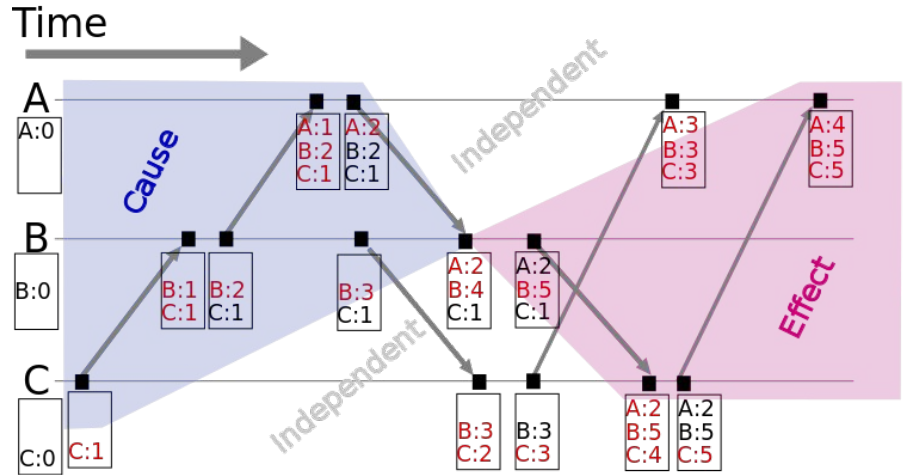
Revert Instrumentation

```
dinv -i -c  
-file=ClientServer.go
```

```
func client(listen, send string) {  
> // sending UDP packet to specified address and port  
> conn := setupConnection(SERVERPORT, CLIENTPORT)  
> for i := 0; i < MESSAGES; i++ {  
>   >   dinvRT.Track("main ClientServer 54 ", "main ClientServer 54 done,main ClientServer 54 isServer  
>   ,main ClientServer 54 isClient,main ClientServer 54 SERVERPORT,main ClientServer 54 MESSAGES,main Cli  
>   entServer 54 CLIENTPORT,main ClientServer 54 cpuprofile,main ClientServer 54 listen,main ClientServer  
>   54 send,main ClientServer 54 conn", done, isServer, isClient, SERVERPORT, MESSAGES, CLIENTPORT, cpuprofile,  
>   listen, send, conn)  
>   >   outgoingMessage := i  
>   >   outBuf := dinvRT.Pack(outgoingMessage)  
>   >   _, errWrite := conn.Write(outBuf)  
>   >   printErr(errWrite)  
>   >   var inBuf [512]byte  
>   >   var incommingMessage int  
>   >   n, errRead := conn.Read(inBuf[0:])  
>   >   printErr(errRead)  
>   >   dinvRT.Unpack(inBuf[0:n], &incommingMessage)  
>   >   incommingMessage = n - n + incommingMessage  
>   >   //fmt.Printf("GOT BACK : %d\n", incommingMessage)  
>   >   time.Sleep(1)  
> }  
> done <- 1  
> }
```

Vector clock refresher

- Distributed systems lack a centralized clock
- Ordering events is therefore hard
- Fundamentally the best that can be done is a *partial order* with happens before
- If **A** receives a message from **B**, the sending event on **B** *happened before* the receive event on **A**

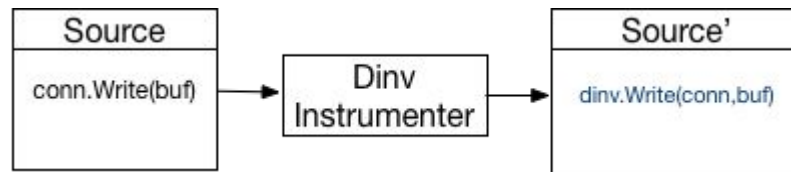


Algorithm:

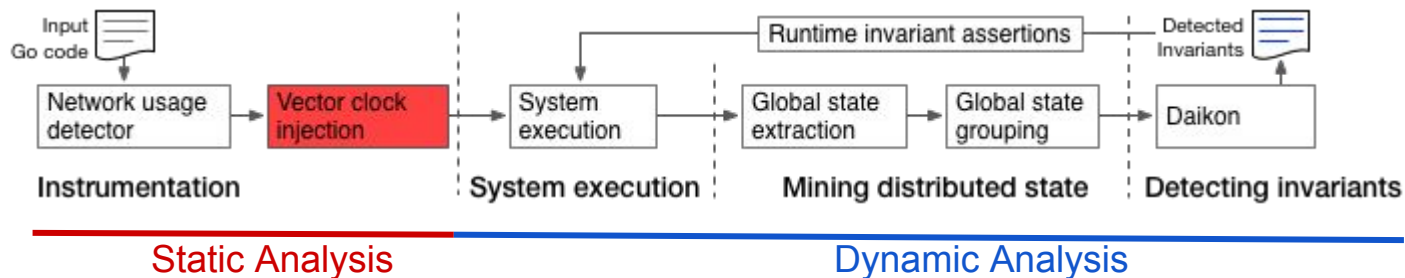
- 1) Increment own index on send & receive
- 2) Take max of all indexes on receive

Tracking time: Vector Clock Instrumentation

- Establish partial event ordering
- Manual and automatic options
- Covers Go standard *net* library



Repository: <https://github.com/DistributedClocks>



Example Vector Clock: Pack/Unpack

Pre-manual Instrumentation:

- Pack line 57

Pack(payload interface{}) []byte

- Unpack line 64

Unpack(buf []byte, toFill interface{})

```
50 func client(listen, send string) {
51 > // sending UDP packet to specified address and port
52 > conn := setupConnection(SERVERPORT, CLIENTPORT)
53 > for i := 0; i < MESSAGES; i++ {
54 > > //@track
55
56 > > outgoingMessage := i
57 > > outBuf := dinvRT.Pack(outgoingMessage)
58 > > _, errWrite := conn.Write(outBuf)
59 > > printErr(errWrite)
60 > > var inBuf [512]byte
61 > > var incommingMessage int
62 > > n, errRead := conn.Read(inBuf[0:])
63 > > printErr(errRead)
64 > > dinvRT.Unpack(inBuf[0:n], &incommingMessage)
65 > > incommingMessage = n - n + incommingMessage
66 > > //fmt.Printf("GOT BACK : %d\n", incommingMessage)
67 > > time.Sleep(1)
68 > }
69 > done <- 1
70
```

VC Instrumentation Options

- Dinv Pack/Unpack take care of marshalling structs!
 - Allows for custom messages to be logged along side vector clocks
- Govector automatically instruments if marshalling is already done
 - Automatic!

GoVector Repository

<https://github.com/DistributedClocks/GoVector>

Command:

```
GoVector -file=filename
```

Method of Injection: AST Rotation

Before:

```
Err = conn.Write(buf)
```

After

```
GoVector.Write(conn.Write,buf)
```


Example Output

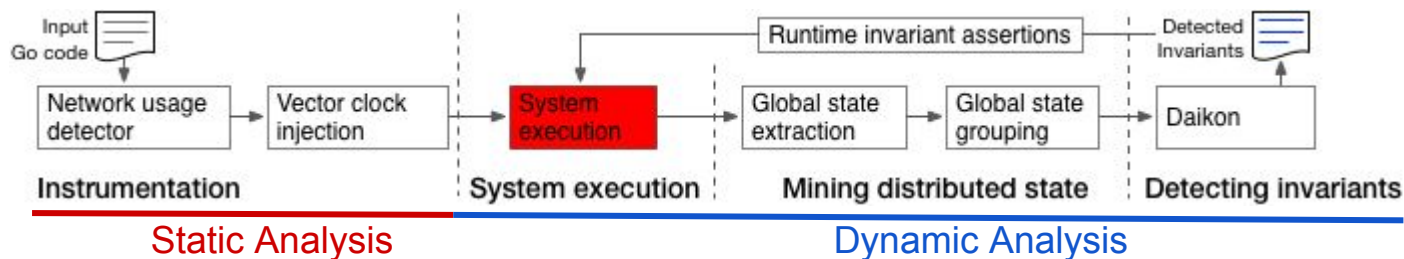
GoVector filename format <nodename>.log-Log.txt

Example Govector output:

```
821589986 {"821589986":1}
Initialization Complete
821589986 {"821589986":2}
Sending from 821589986_main.client+0xaa_/home/stewartgrant/go/src/bitbucket.org/bestchai/dinv/examples/helloDinv/ClientServer.go:57 821589986
821589986 {"821589986":3, "822468001":3}
Received on 821589986_main.client+0x254_/home/stewartgrant/go/src/bitbucket.org/bestchai/dinv/examples/helloDinv/ClientServer.go:64 821589986
821589986 {"822468001":3, "821589986":4}
Sending from 821589986_main.client+0xaa_/home/stewartgrant/go/src/bitbucket.org/bestchai/dinv/examples/helloDinv/ClientServer.go:57 821589986
821589986 {"822468001":5, "821589986":5}
Received on 821589986_main.client+0x254_/home/stewartgrant/go/src/bitbucket.org/bestchai/dinv/examples/helloDinv/ClientServer.go:64 821589986
```


Log Collection (begin dynamic analysis)

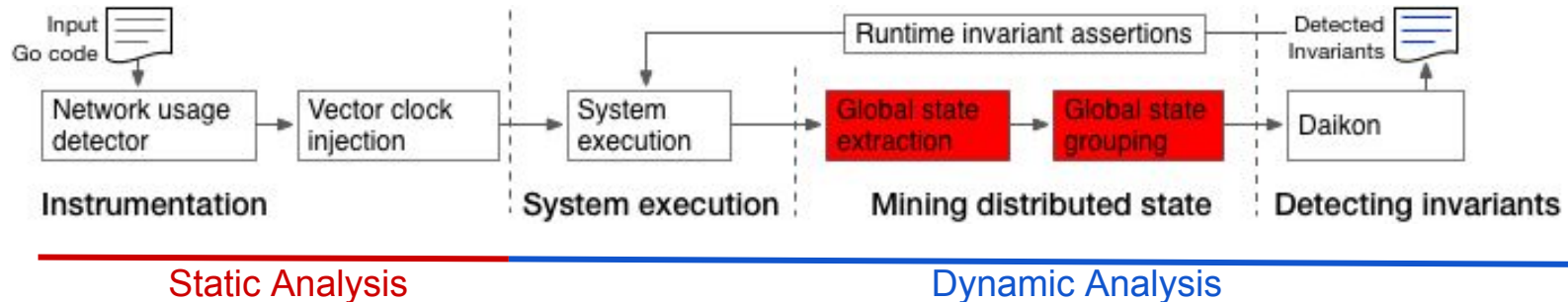
- Analysis - performed on logs collected system execution
- Collection - execute a test suite on an instrumented system
- Quality - of Dinv's invariants improve relative to test exhaustiveness



Dinv Overview

1. Distributed Invariant Inference Challenges

- a. What state should be logged and when?
- b. How to infer distributed invariants from logged state?**
- c. How to enforce distributed invariants?

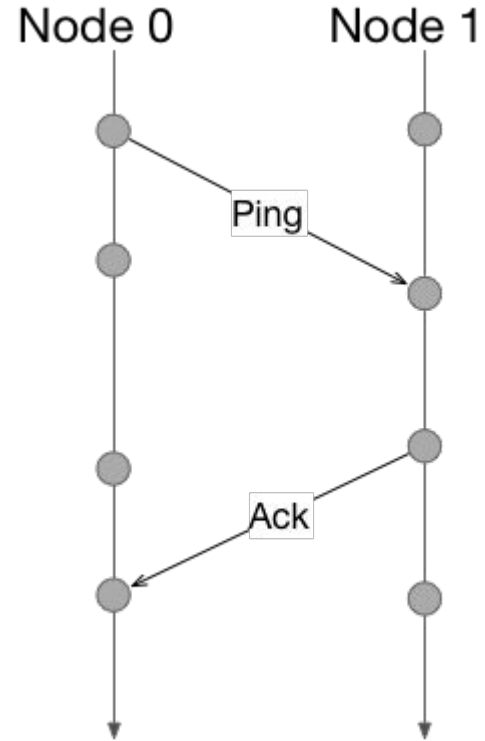


Consistent Cuts

Consistent cut: A partition of an execution, such that causality is preserved.

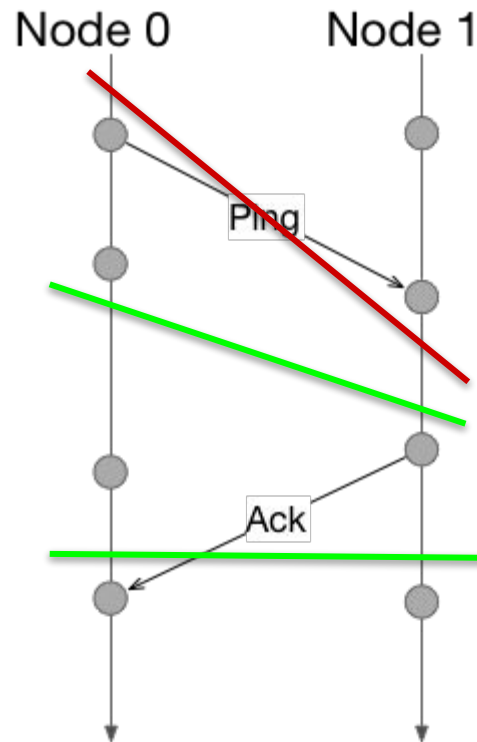
A consistent cut is a global observation of a distributed systems state

Example: Ping and Ack from Serf



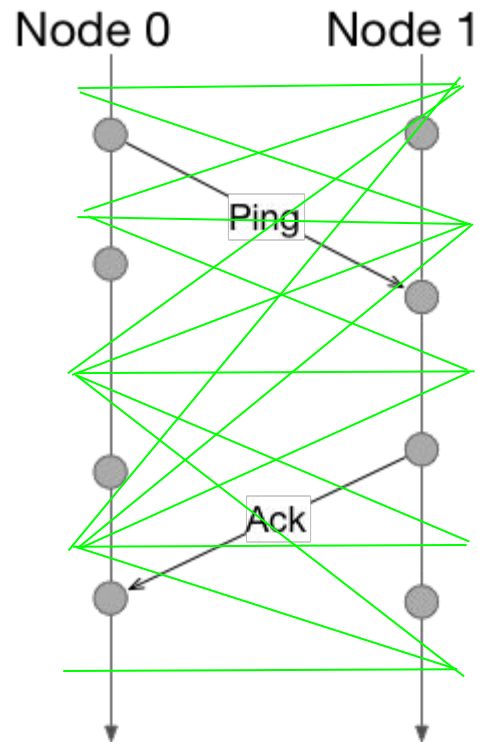
Consistent Cuts

- **Green lines** mark consistent cuts
 - No messages are in flight
 - Message is in flight
- The **red line** is not a consistent cut
 - The ping sent by Node 0 happened before the pings receipt on node 1.



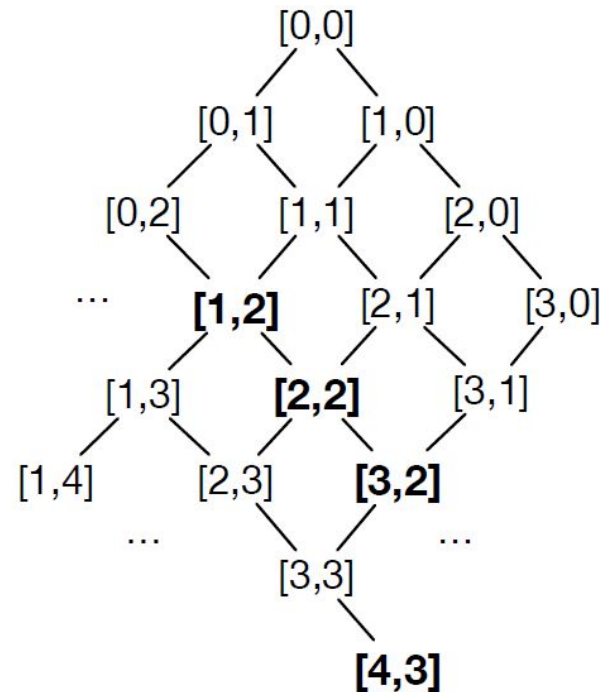
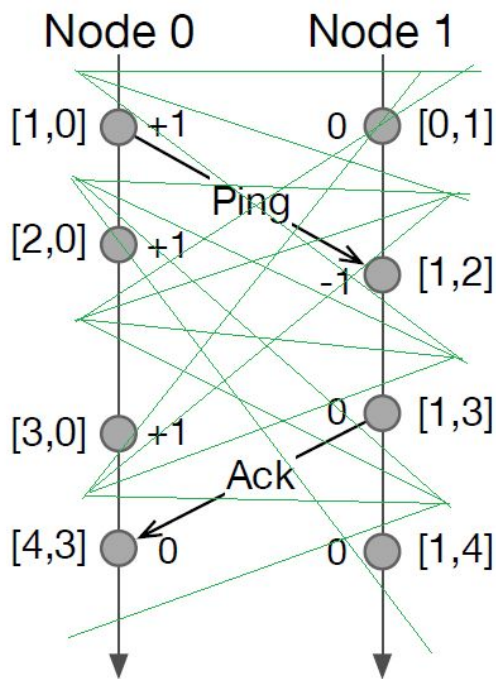
Consistent Cuts

- Executions have an **exponential number** of consistent cuts
- Set of all consistent cuts compose every observable execution path.



Lattice Representation

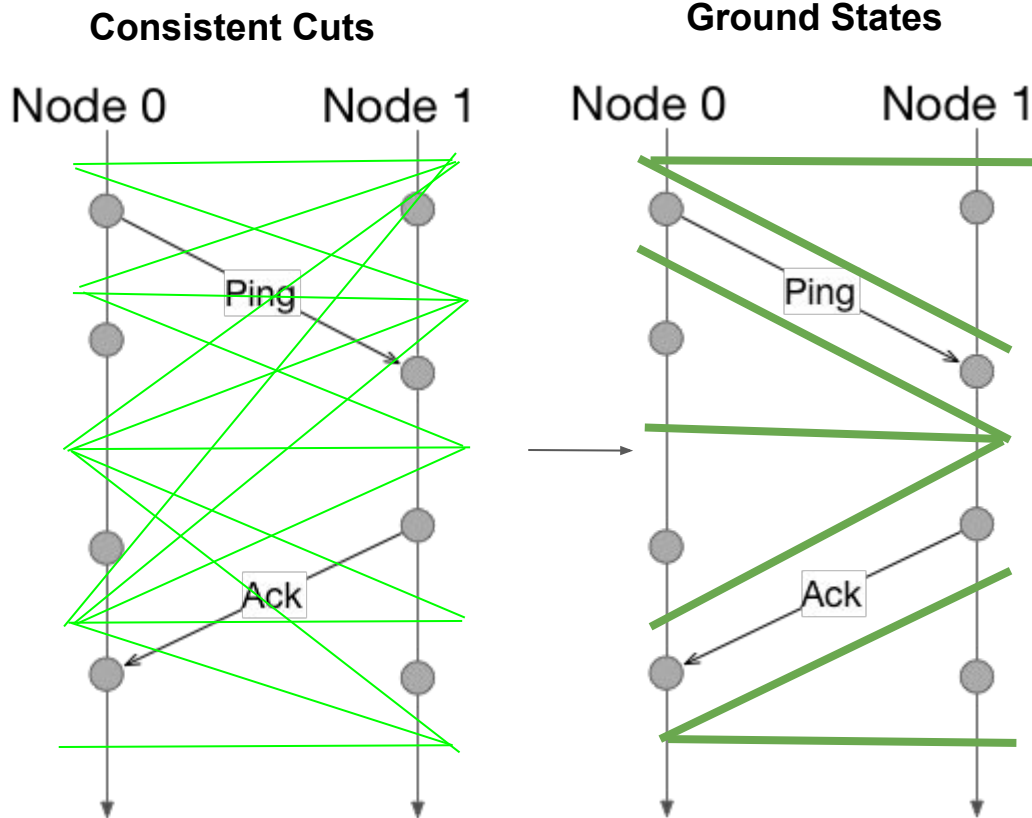
- Cuts are naturally represented as a lattice
- Any path (downward), is a potential execution
- Trillions of points
- Exponentially more paths



Corresponding lattice
(**bold**: no msgs in network)

Ground States

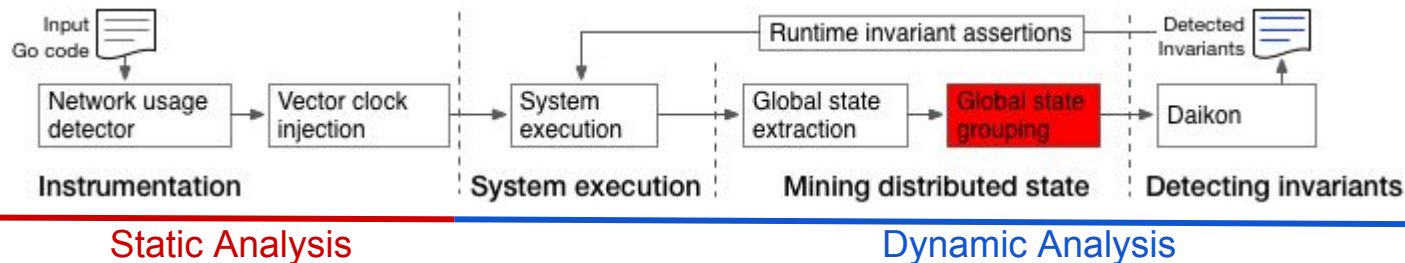
- Consistent cut is massive
- Require sampling heuristic
- **Ground States**: A consistent cut with no in flight messages
- Dramatically collapses search space



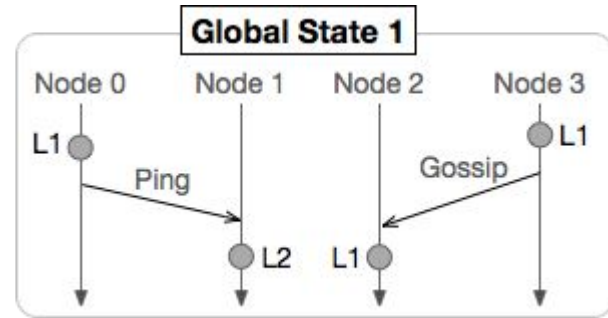
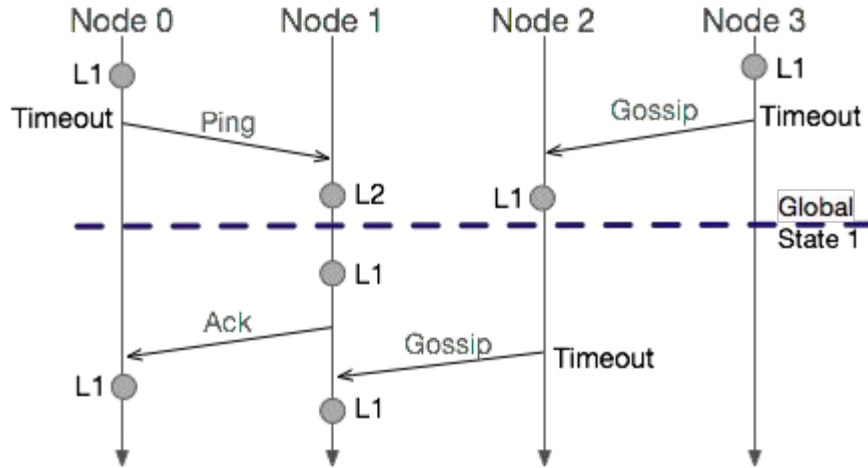
How to infer distributed invariants from global states

- Individual global states are a single instance in time
- Some invariants hold globally, others are protocol specific
- What state should be tested for invariants?

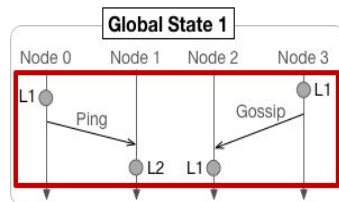
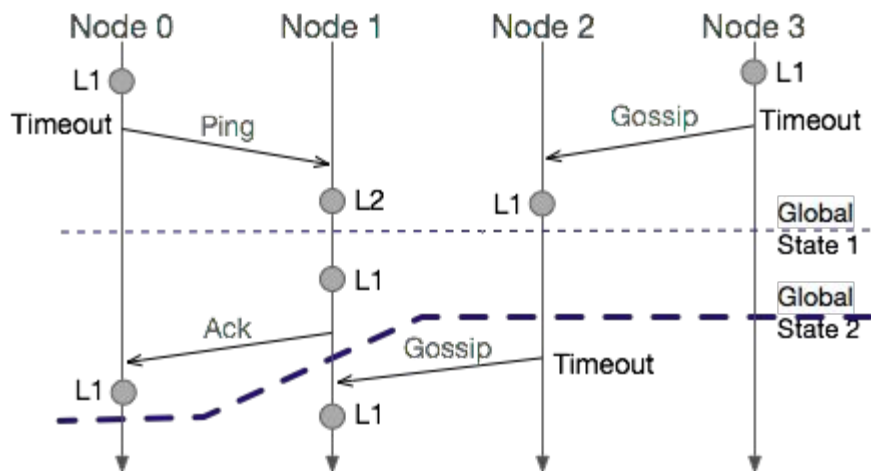
- 1) All - States global merge
- 2) Send-Receive - communication merge
- 3) Total order - transitive merge



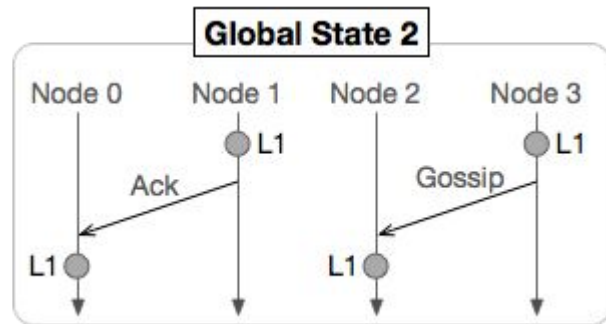
All-State Merging

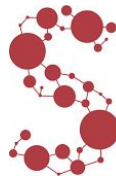


All-State Merging



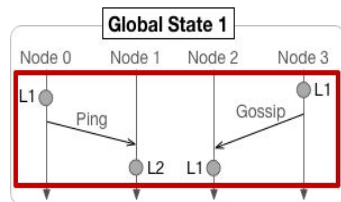
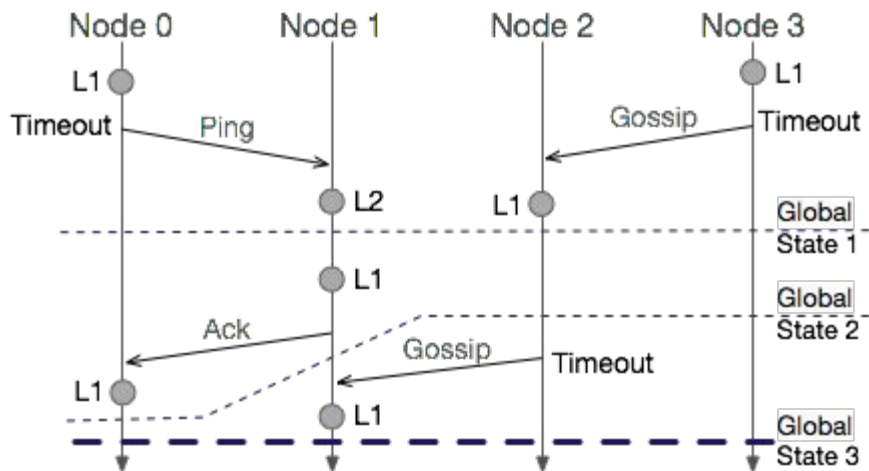
AS: [N0.L1, N1.L2, N2.L1, N3.L1]



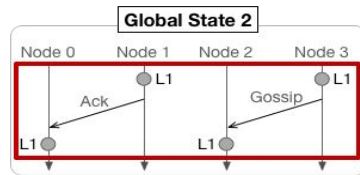


Serf

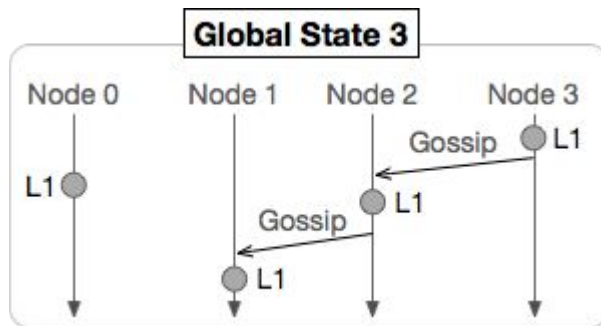
All-State Merging



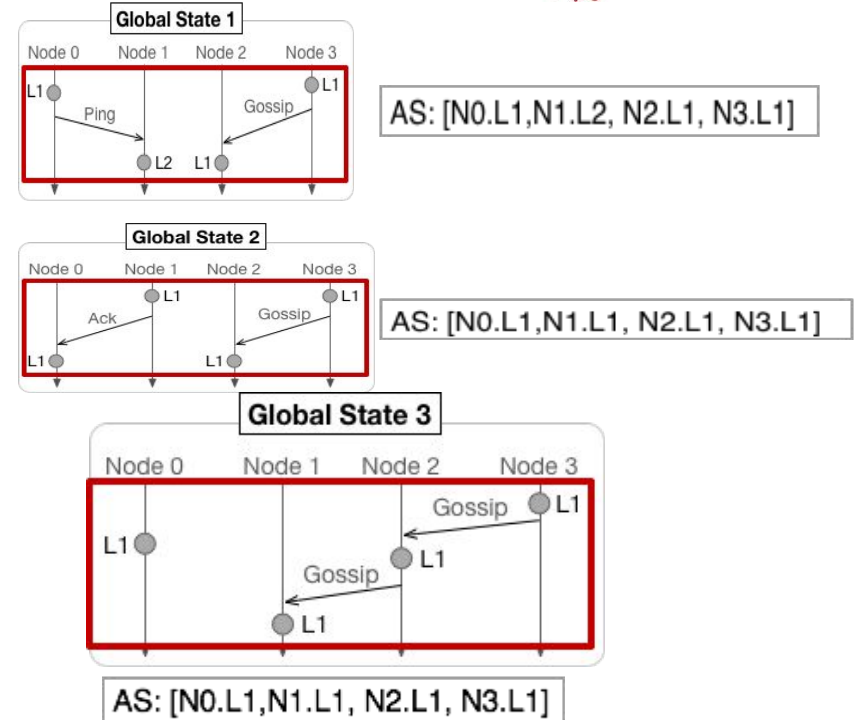
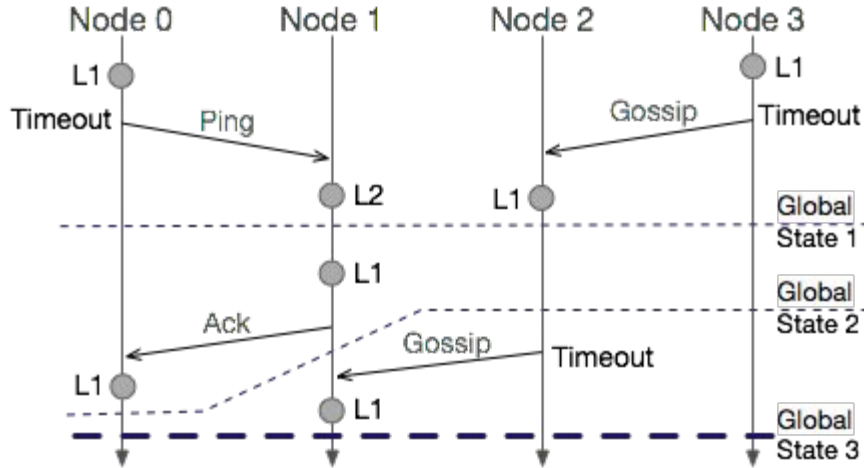
AS: [N0.L1, N1.L2, N2.L1, N3.L1]



AS: [N0.L1, N1.L1, N2.L1, N3.L1]



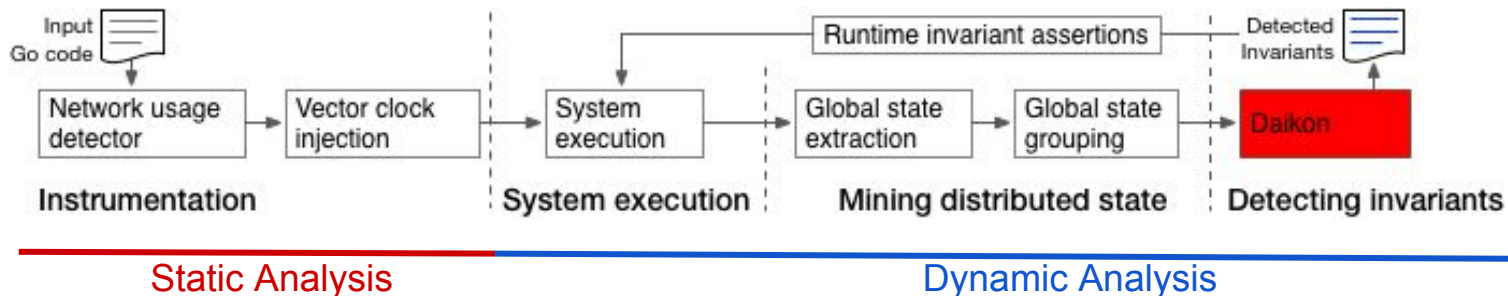
All-State Merging



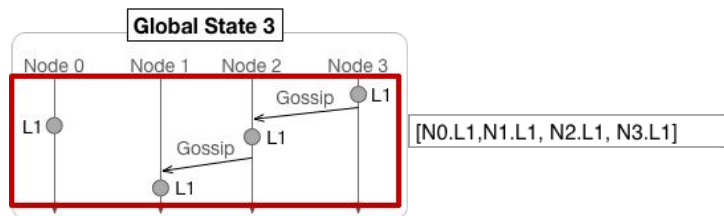


Inferring Invariants

- Daikon tool infers data invariants on data traces
- Insufficient for distributed systems (no partial order)
- Merged states are grouped by IDs and form serialized traces
- Extension for n-ary invariants

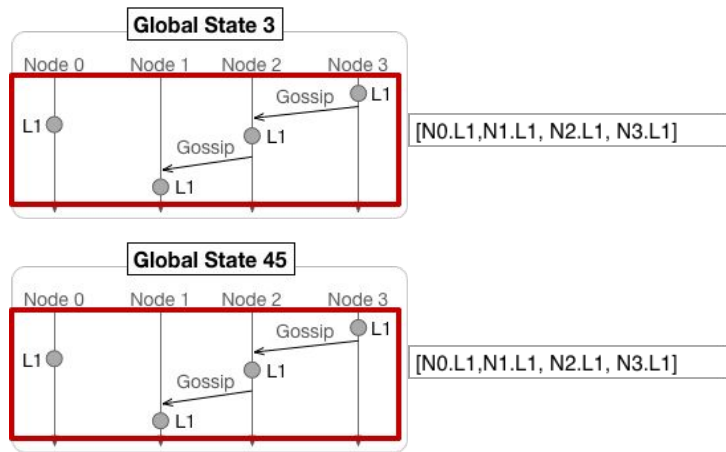


Daikon Bucketing



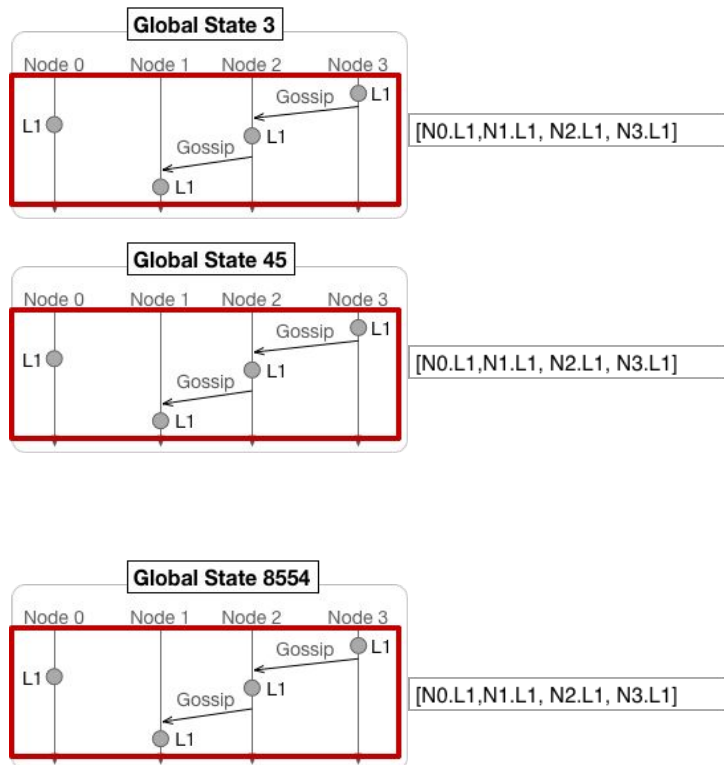
- Distributed States of the same signature are bucketed together

Daikon Bucketing



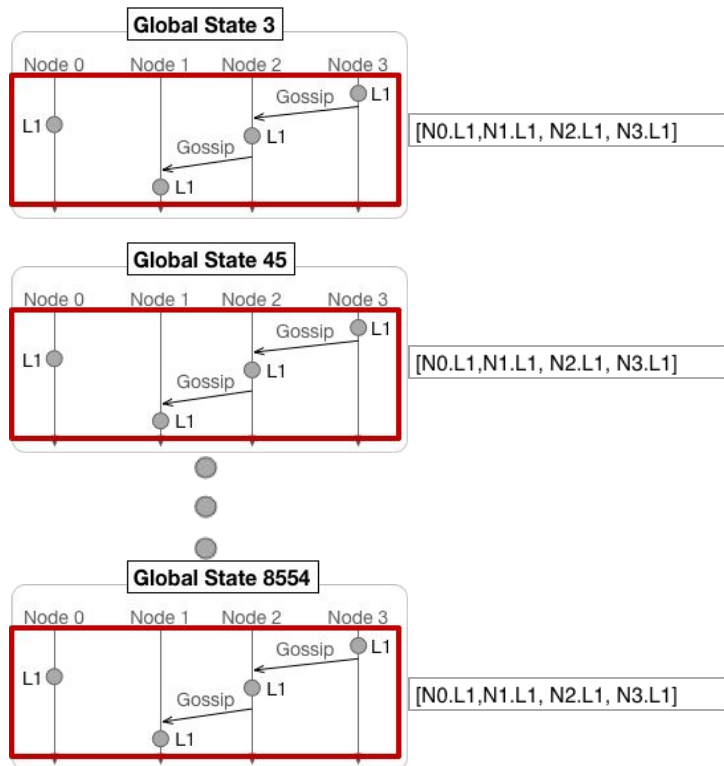
- Distributed States of the same signature are bucketed together
- More states means stronger invariants

Daikon Bucketing



- Distributed States of the same signature are bucketed together
- More states means stronger invariants
- Each new global state provides more evidence that an invariant is true

Daikon Bucketing



- Distributed States of the same signature are bucketed together
- More states means stronger invariants
- Each new global state provides more evidence that an invariant is true
- Ex) $N1_Events == N2_Events == N3_Events$



Distributed Invs

Daikon Template Invariant inference



Daikon systematically tests variables for invariants

Invariants are pre set in templates

Example operators)

`==, >, >=, <, <=,`

`var1 + var2 = var3`

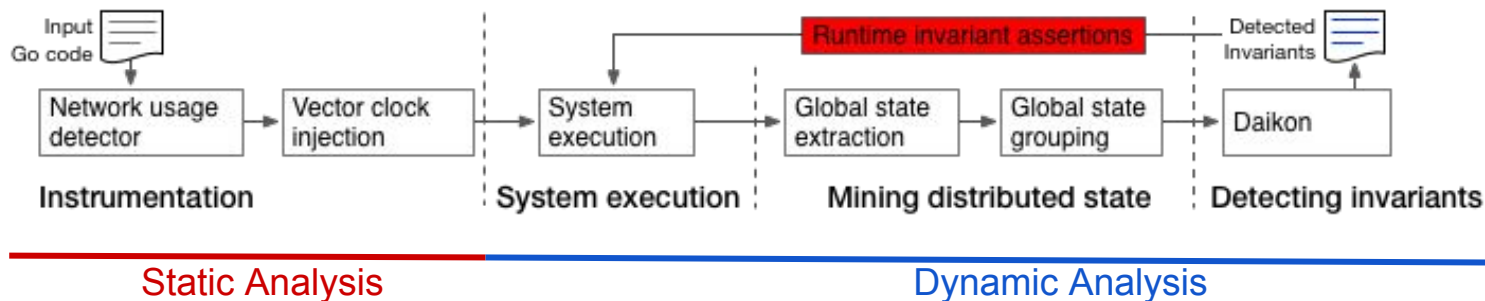
Algorithm - Greater Than

```
For all ints i,j {  
    If i > j {  
        GreaterThan[i][j].Evidence++  
    } Else {  
        GreaterThan[i][j].Invariant = false  
    }  
}
```

Dinv Overview

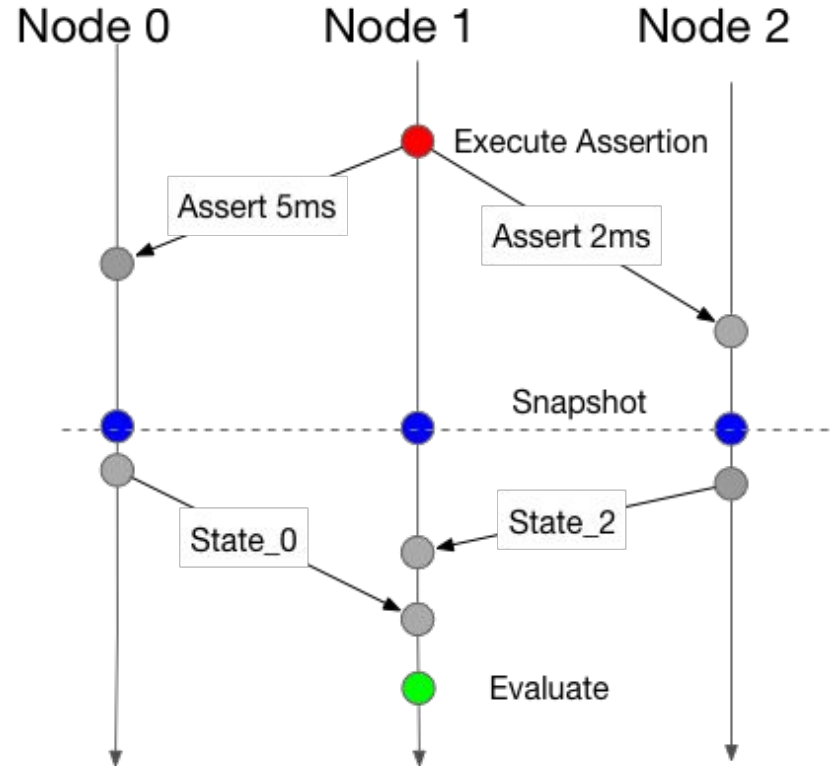
1. Distributed Invariant Inference Challenges

- What state should be logged and when?
- How to infer distributed invariants from logged state?
- How to enforce distributed invariants?**



Assertions

- Dinv has distributed assertions to enforce predicates at runtime
- See me after class if you want an overview!



Dinv Limitations

- Dinv's dynamic analysis is incomplete
- Ground state sampling is poor on loosely coupled systems
- Temporal invariants are currently out of scope



Final take-aways

- Introduced dynamic and static analysis
- Discussed consistent+inconsistent cuts, distributed lattice, ground states
- Dinv tool for detecting data invariants in distributed systems + how it works:
 - Static identification of distributed state
 - Automatic static instrumentation
 - Post-execution merging of distributed states

Source code: <https://bitbucket.org/bestchai/dinv>

Demo: <https://www.youtube.com/watch?v=n9fH9ABJ6S4>