# Kollaps/Thunderstorm: Reproducible Evaluation of Distributed Systems

Miguel Matos

U. Lisbon & INESC-ID Portugal

Tutorial @ DISCOTEC/DAIS 2020

# OUTLINE

- Overview of Kollaps & Thunderstorm

- Hands-on tutorial: basics

- Hands-on tutorial: advanced features

# MOTIVATION

Amazon Found Every 100ms of Latency Cost them 1% in Sales

Nov 10, 2016, 11:43am EST

**Why Brands Are Fighting Over Milliseconds**

Video news

**Buffering reduces video watch time by ~40%, according to research**

September 14, 2016 (4 years ago)

**Post Mortem: What Yesterday's Network Outage Looked Like**

Zalando saw a 0.7% increase in revenue when they shaved 100ms off their load time.

# MOTIVATION

- Performance depends heavily on underlying network
- **Variability** and **Failures** are the norm

- Need for tools for systematic evaluation of distributed applications

- Ability to answer key questions:
  - What is the impact of halving the network latency in application throughput?
  - What is the effect of packet loss?
  - What if …

ANGAINOR

inesc id
lisboa

| Name | Year | Mode | HW ind. | Orchestration | Concurrent deployments | Path congestion | Link-Level emulation capabilities | | | | Any Language | Topology dynamics | Unit |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | Bandwidth | Delay | Packet loss | Jitter | | | |
| DelayLine [47] | 1994 | User | ✓ | Centralized | ✗ | ✗ | ✗ | ✓ | ✓ | ✗ | ✓ | ✗ | P |
| ModelNet [81] | 2002 | Kernel | ✓ | Centralized | ✗ | ✓ | ✓⚡ | ✓⚡ | ✓⚡ | ✗ | ✓ | ✓ | P |
| Nist NET [33] | 2003 | Kernel | ✓ | Centralized | ✗ | ✗ | ✓⚡ | ✓⚡ | ✓⚡ | ✓⚡ | ✓ | ✗ | P |
| NetEm [45] | 2005 | Kernel | ✓ | (N/A: single link emulation only) | | | ✗ | ✓ | ✓ | ✓ | ✓ | ✗ | P |
| Trickle [39] | 2005 | User | ✓ | (N/A: single link emulation only) | | | ✓ | ✓ | ✗ | ✗ | ✓ | ✗ | P |
| EmuSocket [23] | 2006 | User | ✓ | (N/A: single link emulation only) | | | ✓⚡ | ✗ | ✗ | ✗ | ✓ | ✗ | P |
| ACIM/FlexLab [71] | 2007 | Kernel | ✓ | Centralized | | | | ✓⚡ | ✓⚡ | ✓⚡ | ✓ | ✓ | V |
| NCTUns [85] | 2007 | Kernel | ✓ | Centralized | | | | ✓ | ✓ | ✓ | ✓ | ✗ | P |
| Emulab [46, 88] | 2008 | Kernel | ✗ | Centralized | | | | ✓⚡ | ✓⚡ | ✗ | ✓ | ✗ | V |
| IMUNES [70] | 2008 | Kernel | ✗ | Centralized | | | | ✓ | ✓ | ✗ | ✓ | ✗ | P |
| MyP2P-World [75] | 2008 | User | ✓ | Centralized | | | | ✓ | ✓ | ✗ | ✗ | ✗ | P |
| P2PLab [61] | 2008 | Kernel | ✓ | Centralized | | | | ✓ | ✓ | ✗ | ✗ | ✗ | P |
| Netkit [67] | 2008 | Kernel | ✓ | Centralized | | | | ✓⚡ | ✓⚡ | ✗ | ✓ | ✗ | V |
| DFS [79] | 2009 | User | ✓ | Centralized | | | | ✓⚡ | ✓ | ✓ | ✗ | ✓ | P |
| Dummynet [32] | 2010 | Kernel | ✓ | Centralized | | | | ✓⚡ | ✓⚡ | ✗ | ✓ | ✗ | P |
| Mininet [53] | 2010 | Kernel | ✓ | Centralized | | | | ✓⚡ | ✓⚡ | ✓⚡ | ✓ | ✓ | P |
| SliceTime [86] | 2011 | Kernel | ✗ | Centralized | | | | ✓ | ✗ | ✗ | ✓ | ✓ | V |
| Mininet-HiFi [44] | 2012 | Kernel | ✓ | Centralized | | | | ✓⚡ | ✓⚡ | ✓⚡ | ✓ | ✓ | C |
| SplayNet [76] | 2013 | User | ✓ | Decentralized | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✓ | P |
| MaxiNet [87] | 2014 | Kernel | ✓ | Centralized | ✗ | ✓ | ✓⚡ | ✓⚡ | ✓⚡ | ✓⚡ | ✓ | ✓ | P |
| Dockemu [80] | 2015 | User | ✓ | Centralized | ✗ | ✗ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | C |
| EvalBox [77] | 2015 | Kernel | ✓ | Centralized | ✗ | ✗ | ✓⚡ | ✓⚡ | ✓⚡ | ✓⚡ | ✓ | ✓ | P |
| ContainerNet [65] | 2016 | Kernel | ✓ | Centralized | ✗ | ✓ | ✓⚡ | ✓⚡ | ✓⚡ | ✓⚡ | ✓ | ✓ | C,V |
| Kathará [30] | 2018 | Kernel | ✓ | Centralized | ✗ | ✓ | ✓⚡ | ✓⚡ | ✓⚡ | ✗ | ✓ | ✗ | C |
| **KOLLAPS** | **2020** | **Kernel** | ✓ | **Decentralized** | ✓ | ✓ | ✓⚡ | ✓⚡ | ✓⚡ | ✓⚡ | ✓ | ✓ | C,V |

Main limitations:
- scalability
- accuracy
- dynamics

ANGAINOR

inesc id lisboa

# KOLLAPS IN A NUTSHELL

- Applications are concerned about the end-to-end properties:
  - latency, jitter, bandwidth, packet loss
- Rather than the internal network state leading to those properties

# KOLLAPS IN A NUTSHELL

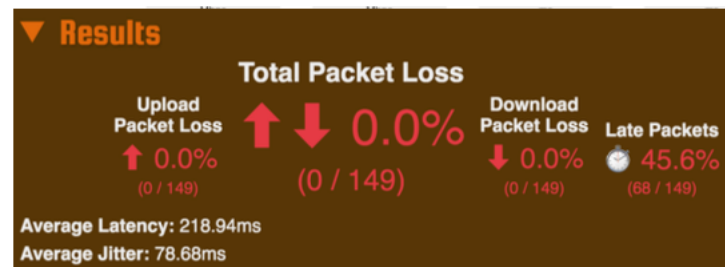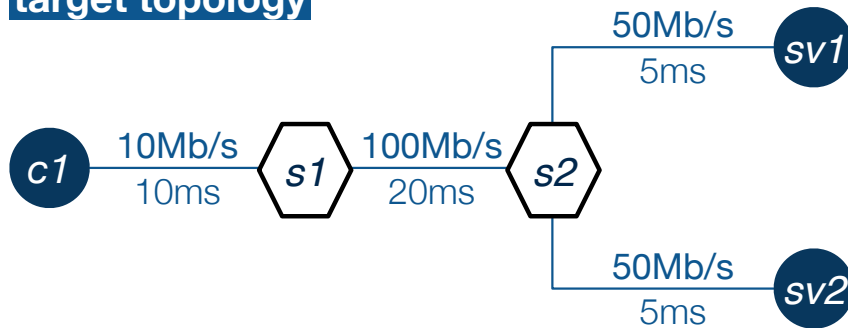- Applications are concerned about the end-to-end properties:
  - latency, jitter, bandwidth, packet loss
- Rather than the internal network state leading to those properties



- Emulate **only** end-to-end properties
- Allows decentralized highly scalable emulation

# NETWORK COLLAPSING

**target topology**

# NETWORK COLLAPSING

**target topology**

c1 — 10Mb/s / 10ms — s1 — 100Mb/s / 20ms — s2 — 50Mb/s / 5ms — sv1

s2 — 50Mb/s / 5ms — sv2

**collapsed topology**

- Node
- Router
- Throughput / Latency

c1 — 10Mb/s / 35ms — 

10Mb/s / 35ms

sv1 — 50Mb/s / 10ms — sv2

# NETWORK COLLAPSING

**target topology**

c1 — 10Mb/s / 10ms — s1 — 100Mb/s / 20ms — s2 — 50Mb/s / 5ms — sv1

s2 — 50Mb/s / 5ms — sv2

**collapsed topology**

- ● Node
- ⬡ Router

Throughput / Latency

c1 — 10Mb/s / 35ms — sv1

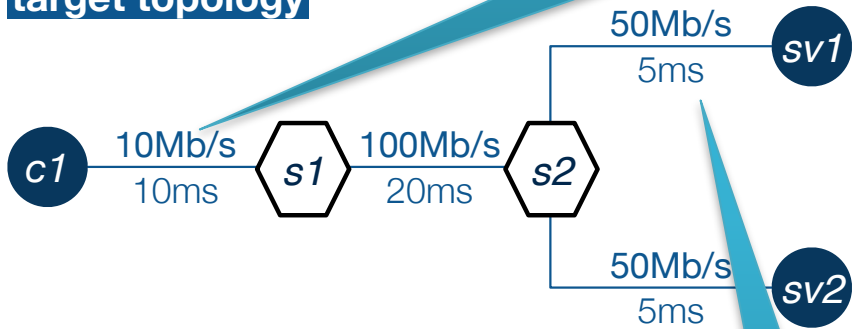c1 — 10Mb/s / 35ms — sv2

sv1 — 50Mb/s / 10ms — sv2

ANGAINOR

inesc id lisboa

# NETWORK COLLAPSING

Minimum bandwidth of all links

**target topology**

**collapsed topology**

c1 —10Mb/s 10ms— s1 —100Mb/s 20ms— s2

50Mb/s 5ms — sv1

50Mb/s 5ms — sv2

Node

Router

Throughput / Latency

10Mb/s 35ms — c1 — 10Mb/s 35ms

sv1 — 50Mb/s 10ms — sv2

Pre-computation of static properties

Sum of latencies of all links
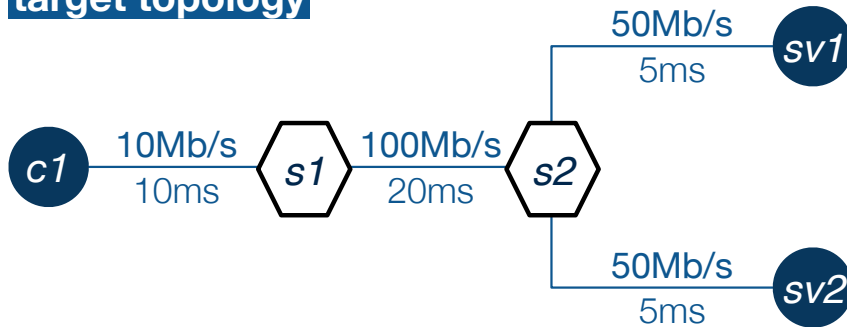
ANGAINOR

inesc id lisboa
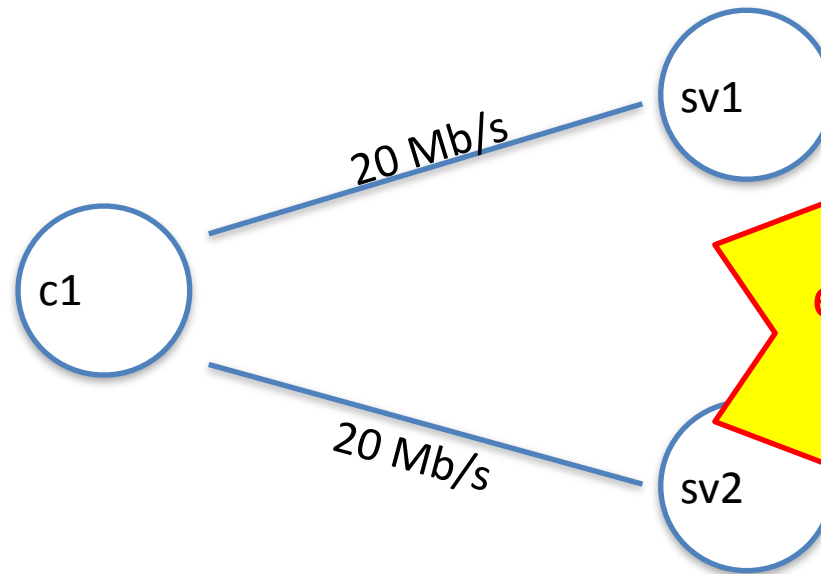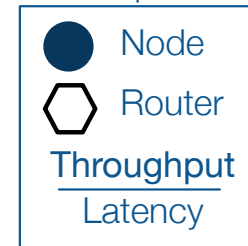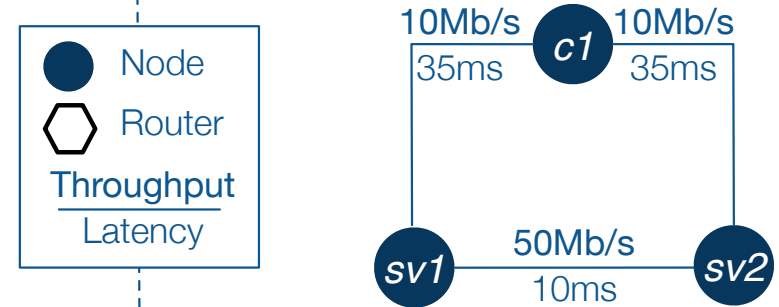
# ARCHITECTURE

# EMULATION MANAGER (EM)

- One instance per physical machine

- Enforces topology properties
  - static properties
  - dynamic properties

# EM: DYNAMIC PROPERTIES

**target topology**

50Mb/s
5ms
sv1

c1 — 10Mb/s / 10ms — s1 — 100Mb/s / 20ms — s2

50Mb/s
5ms
sv2

**collapsed topology**

● Node
⬡ Router

**Throughput**
Latency

10Mb/s / 35ms — c1 — 10Mb/s / 35ms

sv1 — 50Mb/s / 10ms — sv2

sv1

20 Mb/s

c1

20 Mb/s

sv2

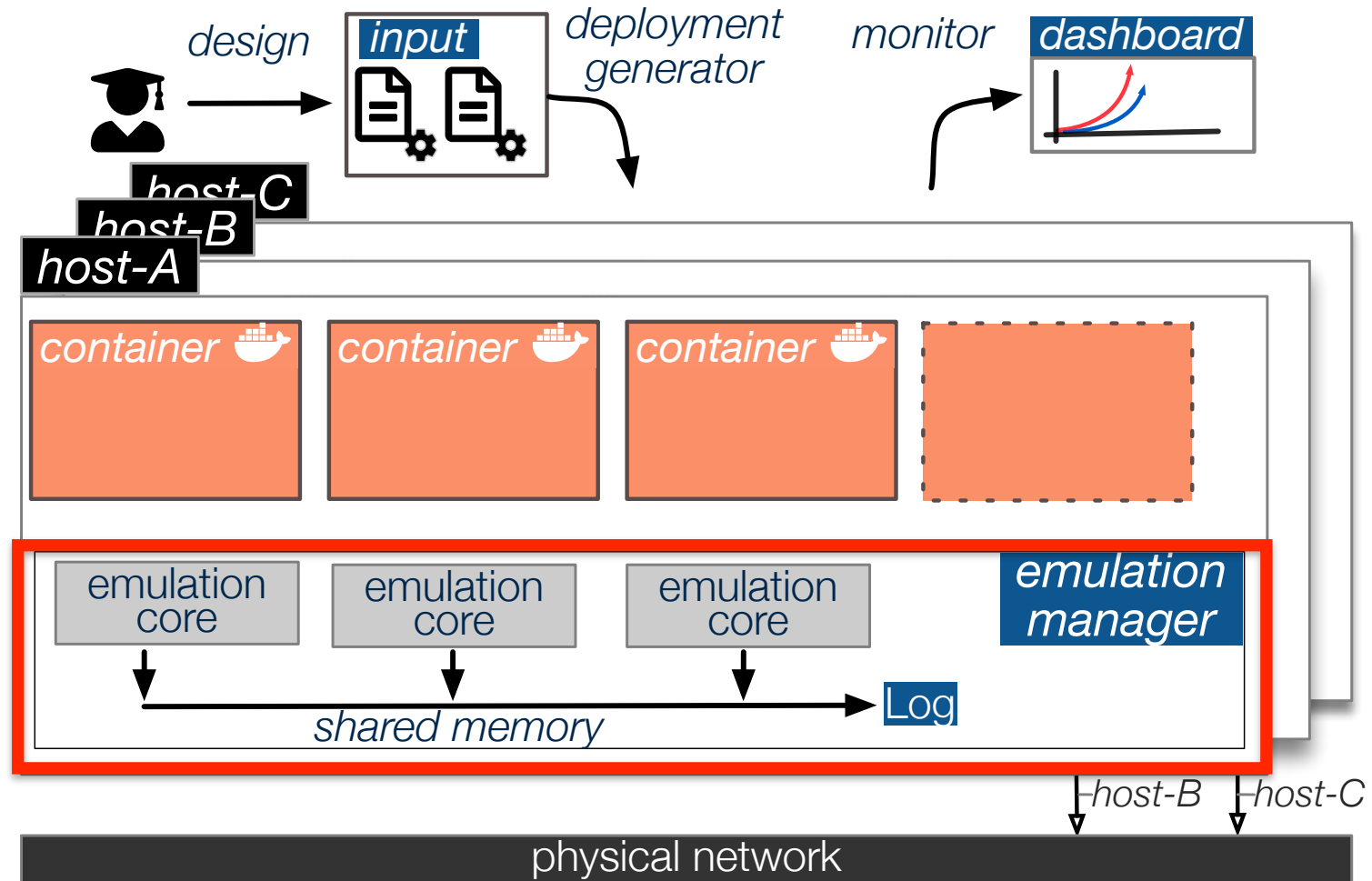How to share enforce bandwidth sharing under congestion?

# EM: DYNAMIC PROPERTIES

- RTT-Aware Min-Max model:

$$Share(f) = \left( RTT(f) \sum_{i=1}^{n} \frac{1}{RTT(f_i)} \right)^{-1}$$

- Intuition
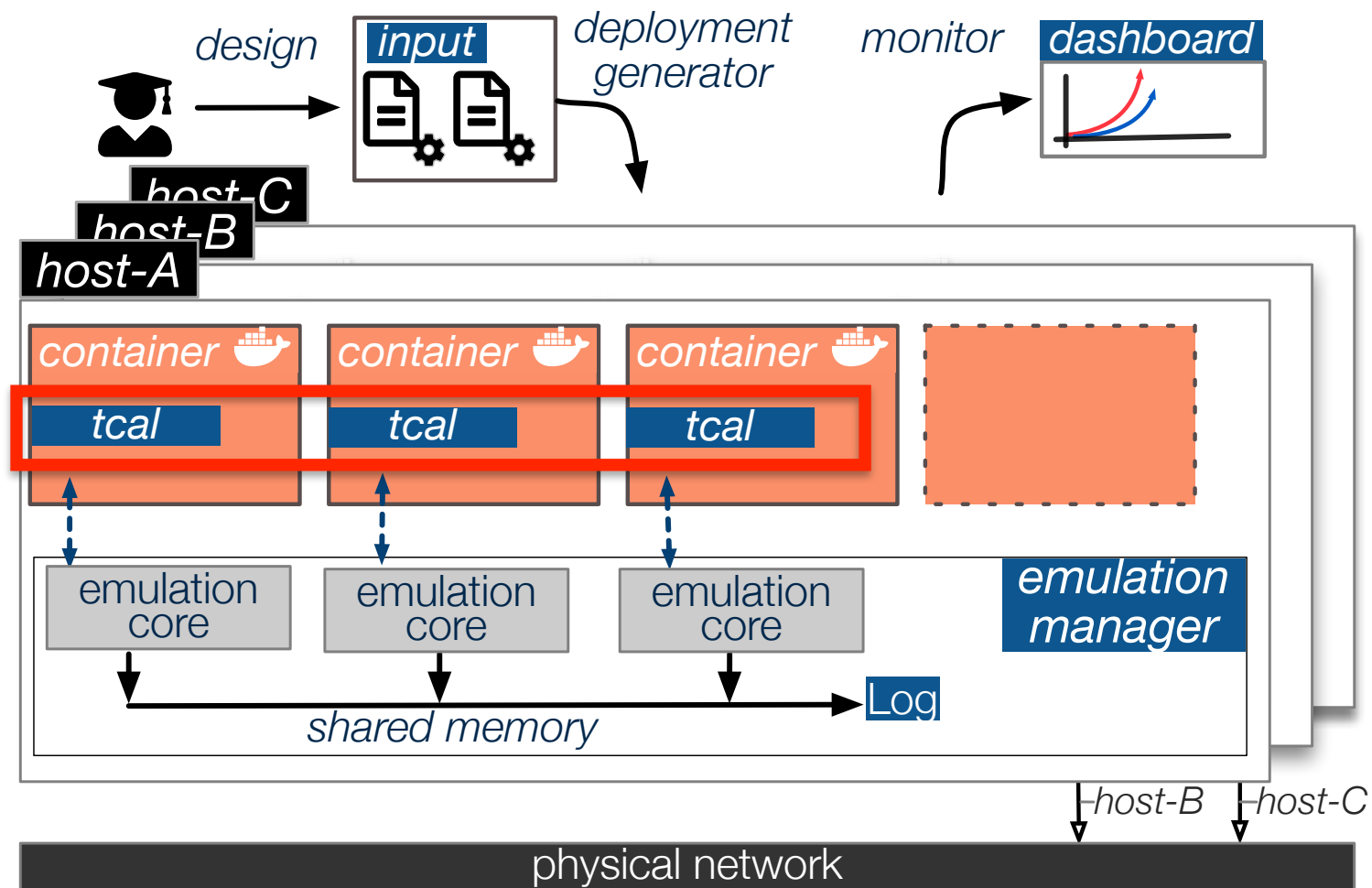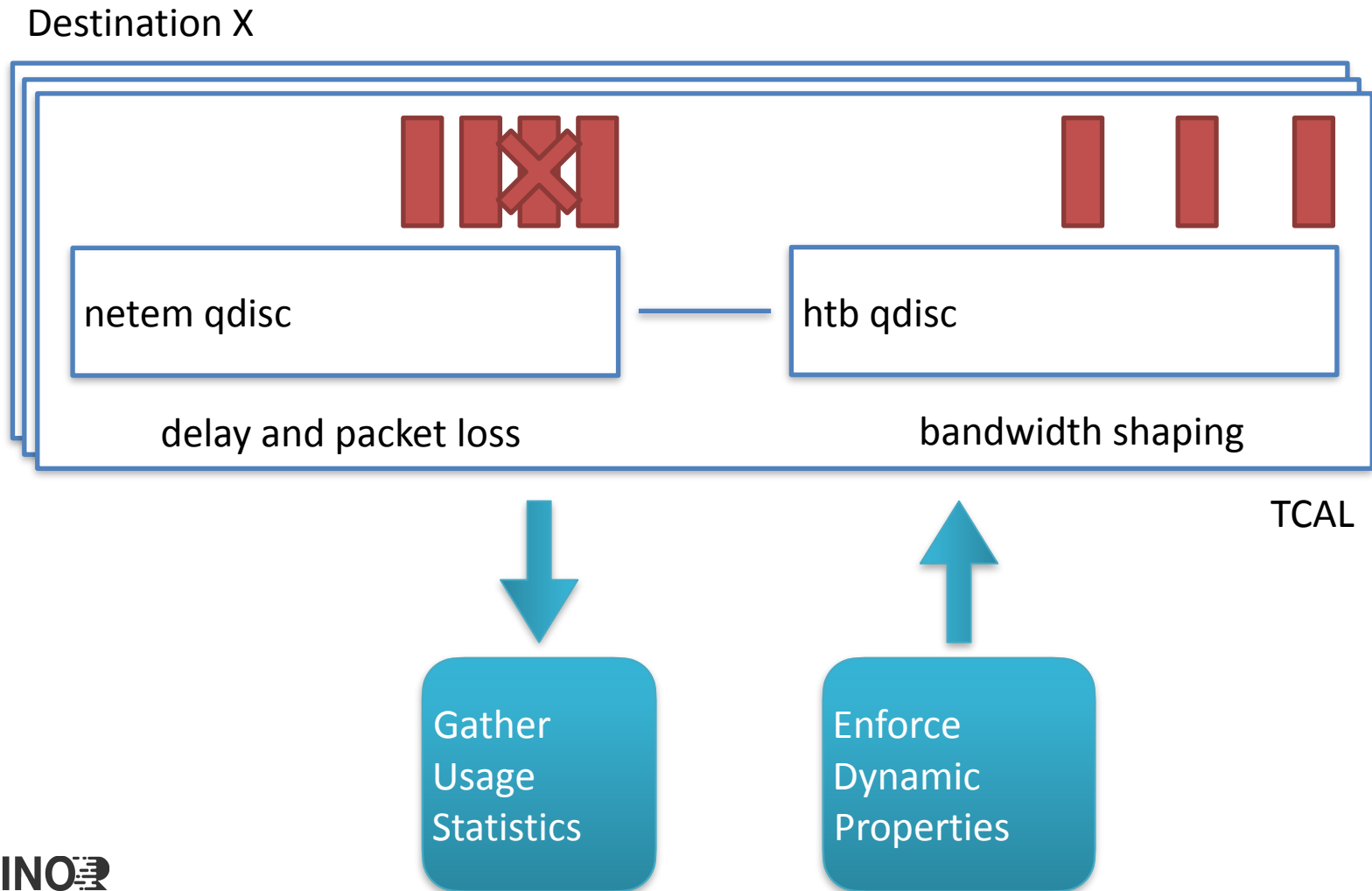  - Available bandwidth is inversely proportional to the RTT

# ARCHITECTURE

# EMULATION CORE

- Spawned by the Emulation Manager

- One instance per container

- Collect's container's usage

- Exchanges metadata with EC through shared memory
  - no bandwidth overhead for local containers

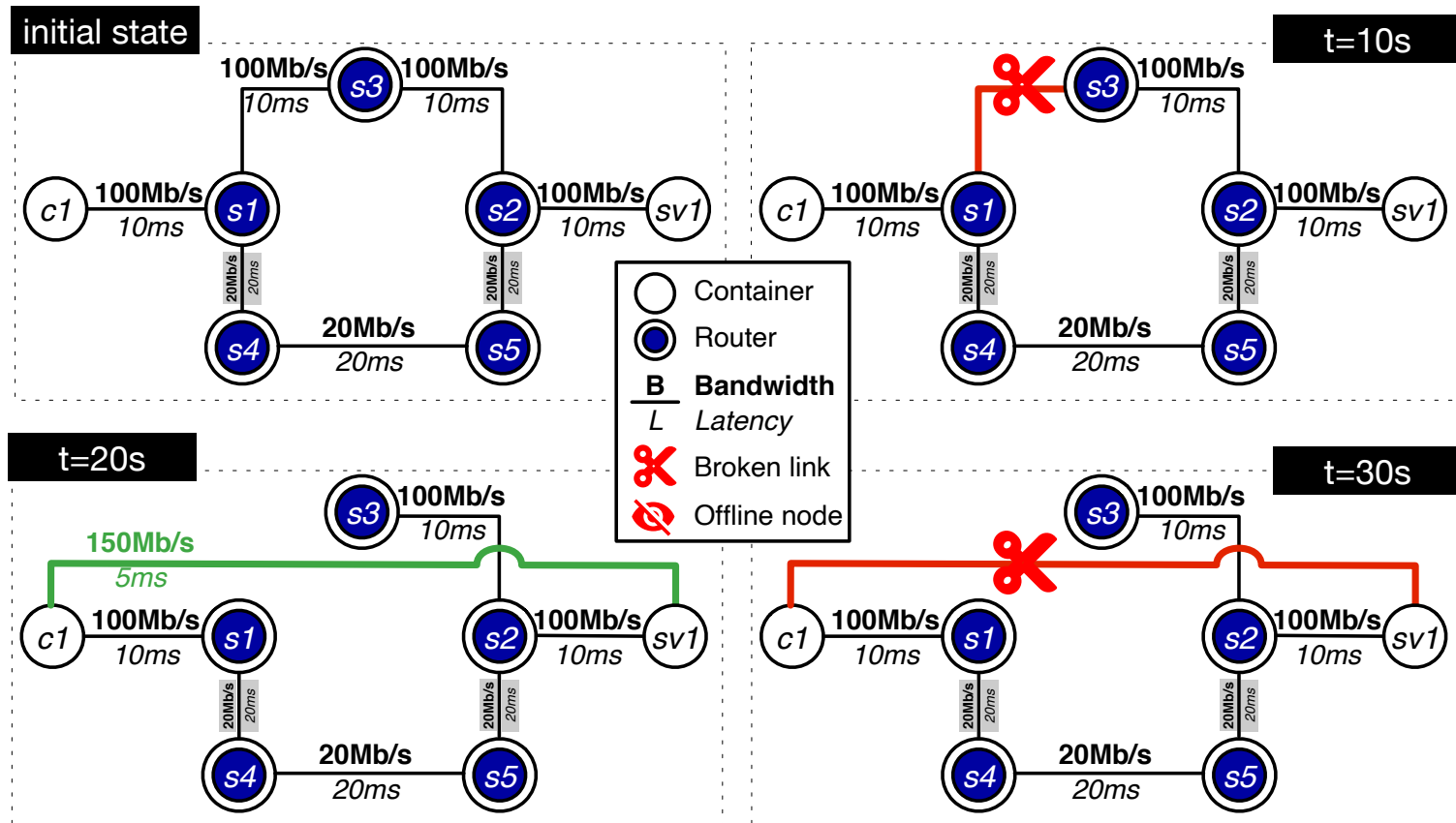# ARCHITECTURE

# LINUX TC ABSTRACTION LAYER

Destination X

netem qdisc ——— htb qdisc

delay and packet loss    bandwidth shaping

TCAL

Gather Usage Statistics

Enforce Dynamic Properties

# DESCRIBE DYNAMIC EXPERIMENTS

- Emulation of network dynamics

# THUNDERSTORM DESCRIPTION LANGUAGE

```
 1  experiment:
 2   services:
 3    name: c1
 4     image: "iperf"
 5    name: sv
 6     image: "nginx"
 7      replicas: 2
 8   bridges:
 9    name: s1
10    name: s2
11   links:
12    orig: c1
13     dest: s1
14     latency: 10
15     up: 10Mbps
16     down: 10Mbps
17     jitter: 0.25
```

```
19  dynamic:
20   orig: c1
21    dest: s1
22    jitter: 0.5
23    time: 120
24   action: leave
25    name: s1
26    time: 200
27   action: join
28    orig: c1
29    dest: s2
30    up: 100Mbps
31    down: 100Mbps
32    latency: 10
33    time: 210
34   action: leave
35    name: sv
```

# EVALUATION

- Link-level emulation
- Scalability and metadata overhead
- Short- and long-lived connections
- Cubic and Reno congestion control algorithms
- Dynamic behavior
- Large-scale topologies
- Reproducing published results
- Geo-replicated Systems
- What-if use cases

# EVALUATION

- Link-level emulation
- Scalability and metadata overhead
- Short- and long-lived connections
- Cubic and Reno congestion control algorithms
- Dynamic behavior
- **Large-scale topologies**
- Reproducing published results
- **Geo-replicated Systems**
- **What-if use cases**

ANGAINOR

# LARGE-SCALE TOPOLOGIES

- Scale-free networks with random ping requests
- Mean-square error w.r.t. theoretical RTT:

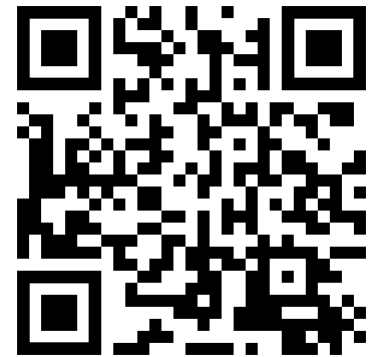| Size (# nodes + # switches) | KOLLAPS | Mininet | Maxinet |
|---|---|---|---|
| 1000 | 0.0261 | 0.0079 | 28.0779 |
| 2000 | 0.0384 | N/A | 347.5303 |
| 4000 | 0.0721 | N/A | N/A |

# GEO-REPLICATED SYSTEM

- Cassandra on EC2 (replication factor: 2)
  - 4 replicas in Frankfurt
  - 4 replicas in Sydney
  - 4 YCSB clients in Frankfurt

# WHAT-IF SCENARIO

- Cassandra on EC2 (replication factor: 2)
  - 4 replicas in Frankfurt
  - 4 replicas in ~~Sydney~~ (Seoul)
  - 4 YCSB clients in Frankfurt

# CONCLUSION AND FUTURE WORK

- KOLLAPS: a decentralized topology emulator

  - Focuses on end-to-end properties

  - Relies on network collapsing techniques

  - Leverages Linux tc and Container Technologies

- Thunderstorm

  - Language to concisely write dynamic experiments

  - Precise description of experiments

    - Key to reproducibility

- Team

  - Miguel Matos, Valerio Schiavoni, Shady Issa, Paulo Gouveia, João Neves, Carlos Segarra, Luca Liechti
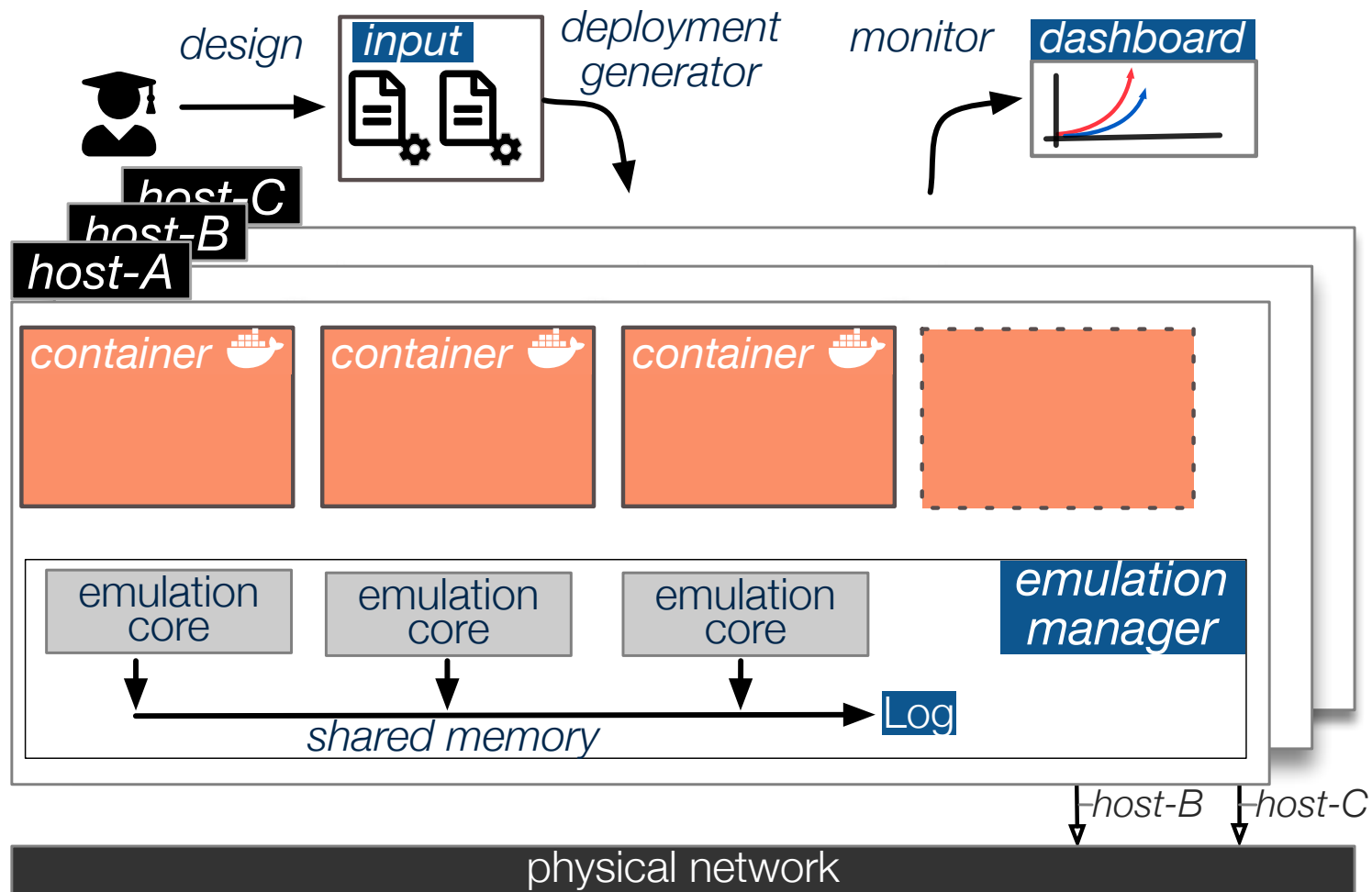
https://github.com/miguelammatos/Kollaps

# PART II: HANDS-ON TUTORIAL

- Overview of Kollaps & Thunderstorm

- **Hands-on tutorial: basics**

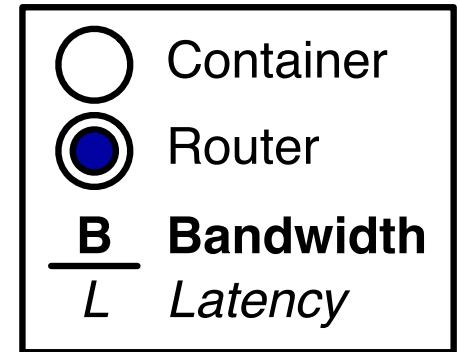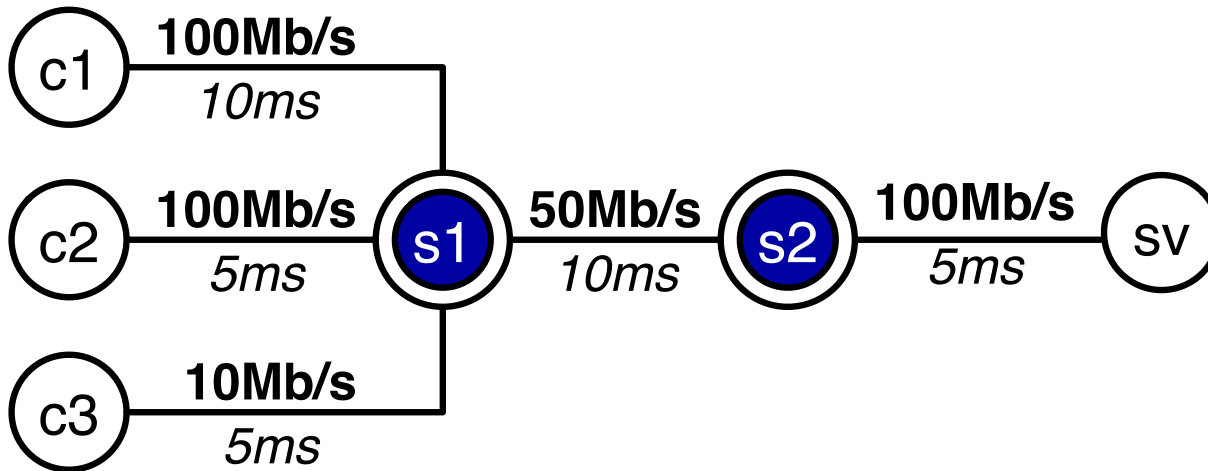- Hands-on tutorial: advanced features

# PART II: HANDS-ON TUTORIAL

- Goal
  - Install Kollaps/Thunderstorm
    - Assumptions
      - *Linux*
      - *Docker and Docker Swarm*
  - Run a simple experiment
    - iPerf3 server
    - iPerf3 client
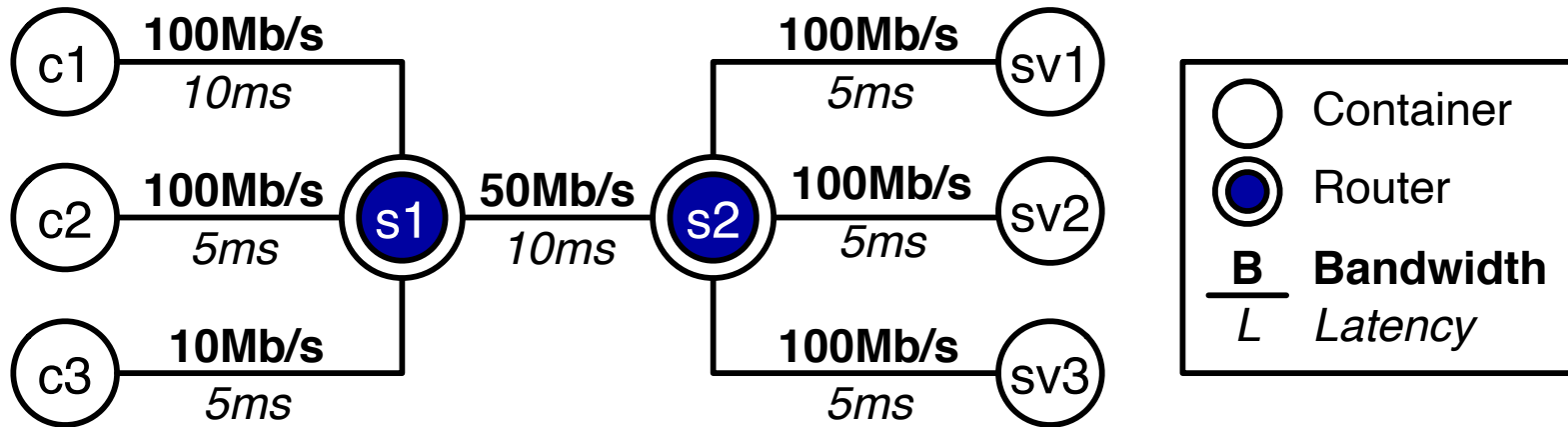    - measure bandwidth
    - measure ping

# ARCHITECTURE

# IPERF3 TOPOLOGY

# IPERF3 TOPOLOGY

# PART II: HANDS-ON TUTORIAL

- Goal
  - Install Kollaps/Thunderstorm
    - Assumptions
      - *Linux*
      - *Docker and Docker Swarm*
  - Run a simple experiment
    - iPerf3 server
    - iPerf3 client

https://github.com/miguelammatos/Kollaps
https://github.com/miguelammatos/Kollaps/wiki

ANGAINOR

inesc id
lisboa