

# Toward Scalable Docker-Based Emulations of Blockchain Networks

---

Diego Pennino and Maurizio Pizzonia

26 May 2023, Bologna

5<sup>th</sup> Distributed Ledger Technology Workshop  
DLT 2023

# Blockchain Network Testing

Performing realistic experiments for blockchain networks is notoriously hard.

# Blockchain Network Testing

Performing realistic experiments for blockchain networks is notoriously hard.

The complexity:

# Blockchain Network Testing

Performing realistic experiments for blockchain networks is notoriously hard.

The complexity:

- Large number of nodes,

# Blockchain Network Testing

Performing realistic experiments for blockchain networks is notoriously hard.

The complexity:

- Large number of nodes,
- Software complexity,

# Blockchain Network Testing

Performing realistic experiments for blockchain networks is notoriously hard.

The complexity:

- Large number of nodes,
- Software complexity,
- Communications among nodes (properties of transport protocols),

# Blockchain Network Testing

Performing realistic experiments for blockchain networks is notoriously hard.

The complexity:

- Large number of nodes,
- Software complexity,
- Communications among nodes (properties of transport protocols),
- Nodes are spread over the internet (delay and packet loss)

# Blockchain Network Testing

Performing realistic experiments for blockchain networks is notoriously hard.

The complexity:

- Large number of nodes,
- Software complexity,
- Communications among nodes (properties of transport protocols),
- Nodes are spread over the internet (delay and packet loss)

**Emulation or Simulation???**



# Blockchain Network Testing

Performing realistic experiments for blockchain networks is notoriously hard.

The complexity:

- Large number of nodes,
- Software complexity,
- Communications among nodes (properties of transport protocols),
- Nodes are spread over the internet (delay and packet loss)

**Simulation**

**Emulation or Simulation???**

# Blockchain Network Testing

Performing realistic experiments for blockchain networks is notoriously hard.

The complexity:

- Large number of nodes, **Simulation**
- Software complexity, **Emulation**
- Communications among nodes (properties of transport protocols), **Emulation**
- Nodes are spread over the internet (delay and packet loss) **Emulation**

**Emulation or Simulation???**

# Blockchain Network Testing

Performing realistic experiments for blockchain networks is notoriously hard.

The complexity:

- Large number of nodes, **Simulation**
- Software complexity, **Emulation**
- Communications among nodes (properties of transport protocols), **Emulation**
- Nodes are spread over the internet (delay and packet loss) **Emulation**

Emulation or **Simulation**???

## Emulation: two choices



# Emulation: two choices

## PRO:

- Simple to handle
- Simple to modify

## CONS:

- Small amount of nodes (few hundreds)
- Fake Networks

## Local



# Emulation: two choices

## PRO:

- Huge amount of nodes
- “Real” Network environment

## CONS:

- Hard to handle
- Slow to modify
- Clusters

**Distributed**



## PRO:

- Simple to handle
- Simple to modify

## CONS:

- Small amount of nodes (few hundreds)
- Fake Networks

**Local**



# Emulation: Our Solution

## PRO:

- Huge amount of nodes
- “Real” Network environment

## CONS:

- Hard to handle
- Slow to modify
- Clusters

Distributed



## PRO:

- Simple to handle
- Simple to modify

## CONS:

- Small amount of nodes (few hundreds)
- Fake Networks

Local



## Boosted Hardware on premises



### PRO:

- Huge amount of nodes
- “Real” Network environment

### CONS:

- ~~• Hard to handle~~
- ~~• Slow to modify~~
- ~~• Clusters~~

### PRO:

- Simple to handle
- Simple to modify

### CONS:

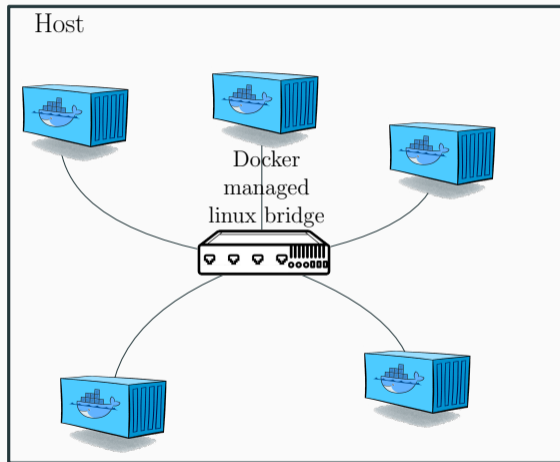
- ~~• Small amount of nodes  
(few hundreds)~~
- ~~• Fake Networks~~



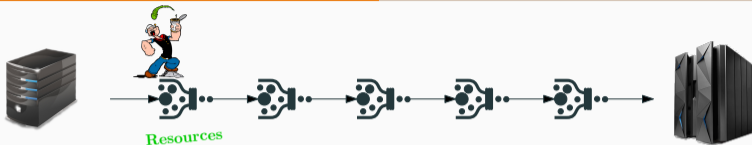
# Roadmap



# The typical Docker scenario



# 1<sup>st</sup> bottleneck: Resources usage



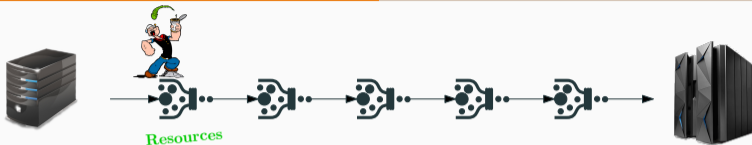
Security limits

`/etc/security/limits.conf`

Kernel parameters

`/etc/sysctl.conf`

# 1<sup>st</sup> bottleneck: Resources usage



Security limits

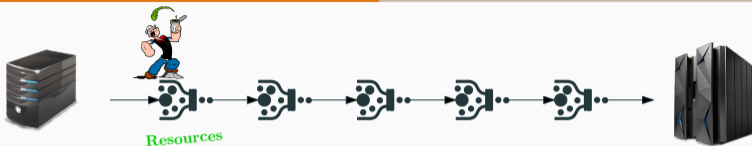
`/etc/security/limits.conf`

{ **nfile** number of open files  
**nproc** maximum number of processes

Kernel parameters

`/etc/sysctl.conf`

# 1<sup>st</sup> bottleneck: Resources usage



Security limits

`/etc/security/limits.conf`

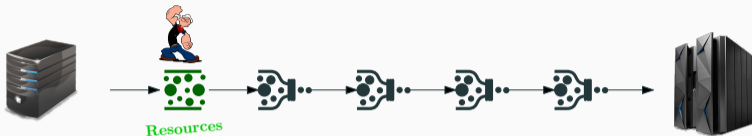
{ **nofile** number of open files  
**nproc** maximum number of processes

Kernel parameters

`/etc/sysctl.conf`

{ **pty** maximum number of pseudo-terminal def:4096  
**gc\_thresh1** garbage collector ARP entries def:128  
**gc\_thresh2** garbage collector ARP entries def:512  
**gc\_thresh3** garbage collector ARP entries def:1024

# 1<sup>st</sup> bottleneck: Resources usage



Security limits

`/etc/security/limits.conf`

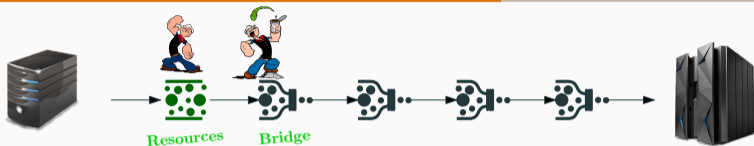
{ **nofile** number of open files  
**nproc** maximum number of processes

Kernel parameters

`/etc/sysctl.conf`

{ **pty** maximum number of pseudo-terminal def:4096  
**gc\_thresh1** garbage collector ARP entries def:128  
**gc\_thresh2** garbage collector ARP entries def:512  
**gc\_thresh3** garbage collector ARP entries def:1024

## 2<sup>nd</sup> bottleneck: The Bridge

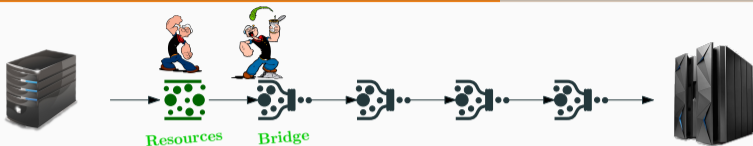


Default Linux Bridge



$2^{10} = 1024$  ports

## 2<sup>nd</sup> bottleneck: The Bridge

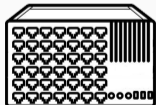


Default Linux Bridge



$$2^{10} = 1024 \text{ ports}$$

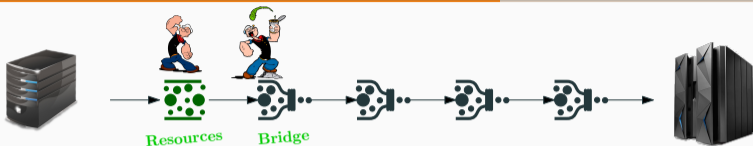
Our Linux Bridge



$$2^{17} = 131,072 \text{ ports}$$



## 2<sup>nd</sup> bottleneck: The Bridge

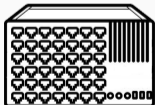


Default Linux Bridge



$2^{10} = 1024$  ports

Our Linux Bridge

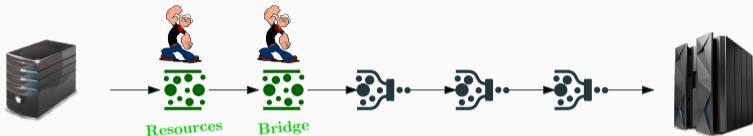


$2^{17} = 131,072$  ports

Kernel

```
/ net / bridge / br_private.h
23 #define BR_HASH_SIZE (1 << BR_HASH_BITS)
24
25 #define BR_HOLD_TIME (1*HZ)
26
27 #define BR_PORT_BITS 10
28 #define BR_MAX_PORTS (1<<BR_PORT_BITS)
29
```

## 2<sup>nd</sup> bottleneck: The Bridge

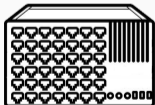


Default Linux Bridge



$2^{10} = 1024$  ports

Our Linux Bridge



$2^{17} = 131,072$  ports

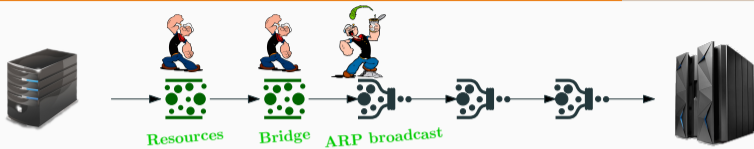
Kernel

```
/ net / bridge / br_private.h
23 #define BR_HASH_SIZE (1 << BR_HASH_BITS)
24
25 #define BR_HOLD_TIME (1*HZ)
26
27 #define BR_PORT_BITS 10
28 #define BR_MAX_PORTS (1<<BR_PORT_BITS)
29
```

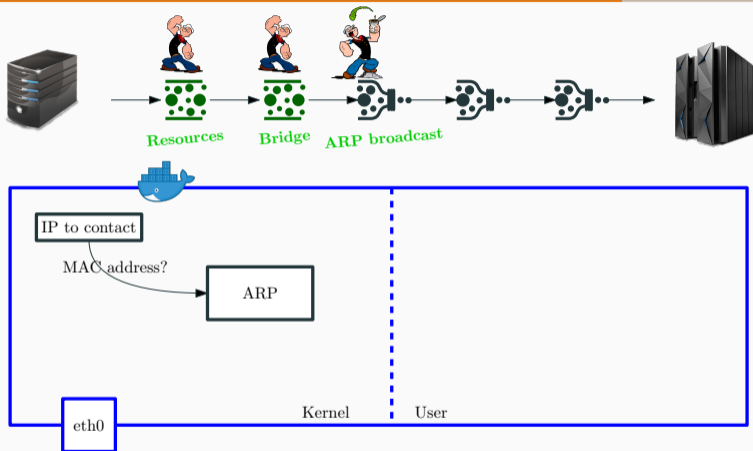
# 3<sup>rd</sup> bottleneck: ARP Broadcast



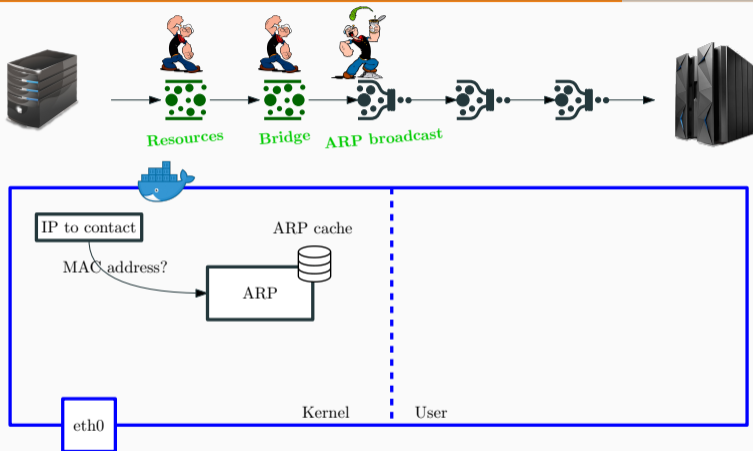
# 3<sup>rd</sup> bottleneck: ARP Broadcast



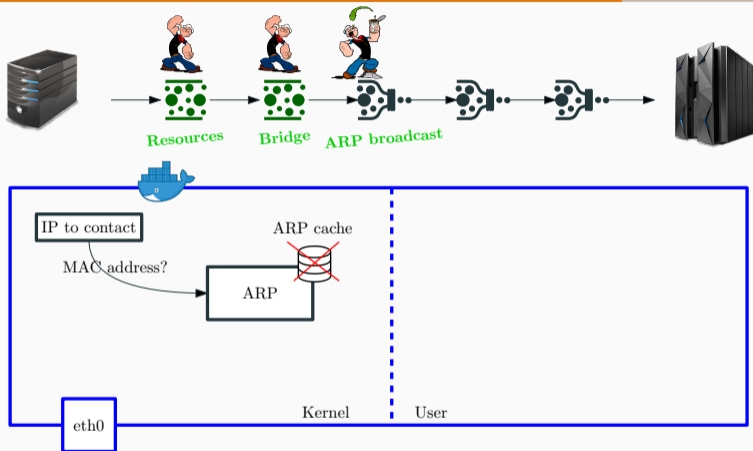
# 3<sup>rd</sup> bottleneck: ARP Broadcast



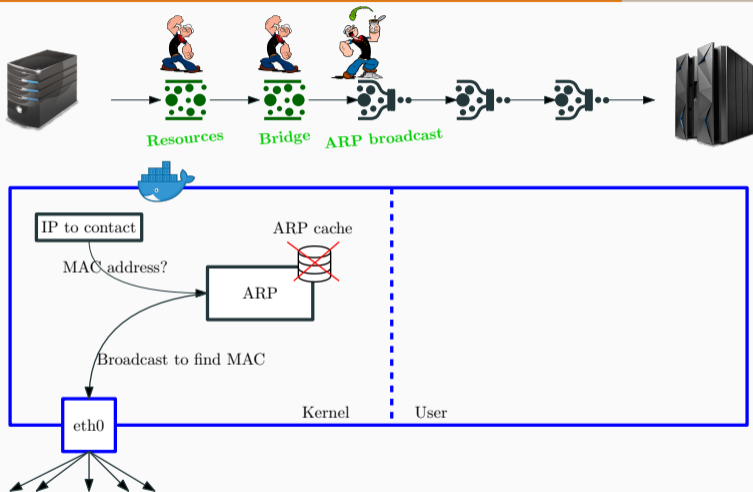
# 3<sup>rd</sup> bottleneck: ARP Broadcast



# 3<sup>rd</sup> bottleneck: ARP Broadcast

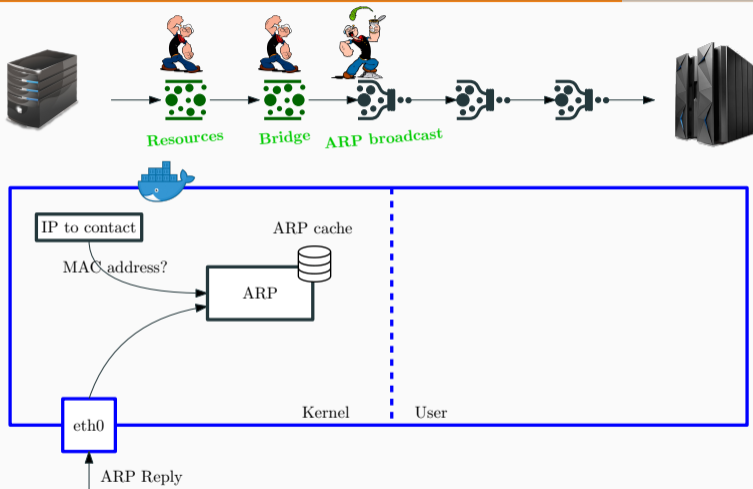


# 3<sup>rd</sup> bottleneck: ARP Broadcast

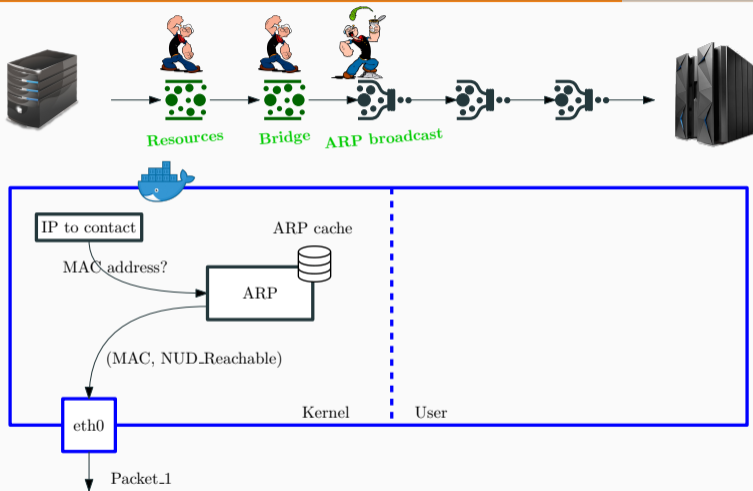




# 3<sup>rd</sup> bottleneck: ARP Broadcast

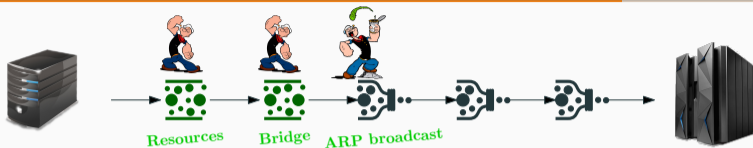


# 3<sup>rd</sup> bottleneck: ARP Broadcast

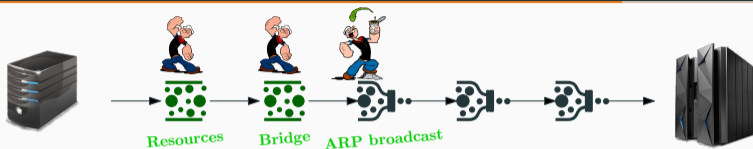


**Neighbour Unreachability Detection(NUD):** *Reachable*= Valid entry recently used

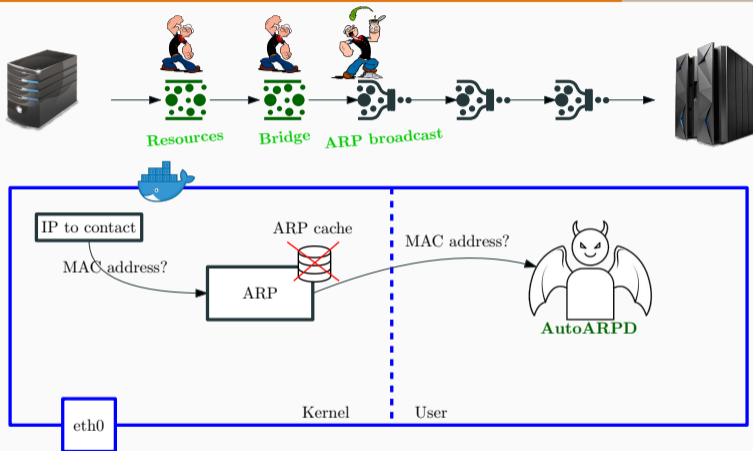
## 3<sup>rd</sup> bottleneck: ARP Broadcast



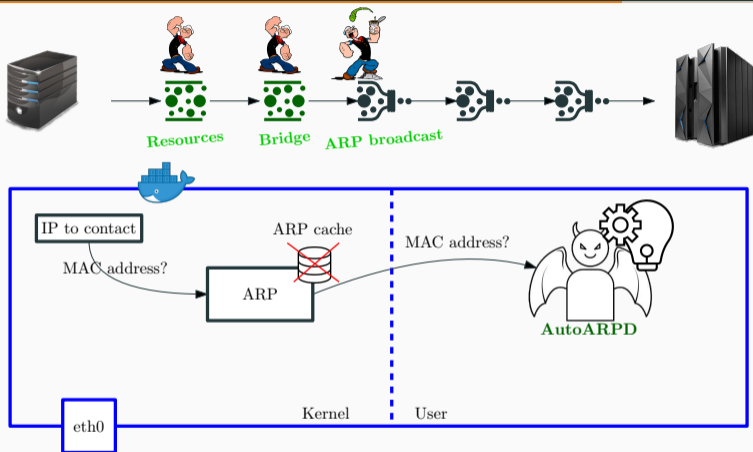
## 3<sup>rd</sup> bottleneck: ARP Broadcast



# 3<sup>rd</sup> bottleneck: ARP Broadcast



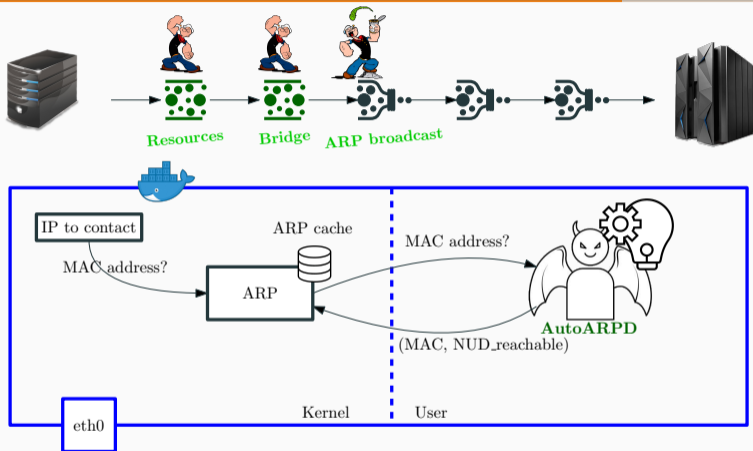
## 3<sup>rd</sup> bottleneck: ARP Broadcast



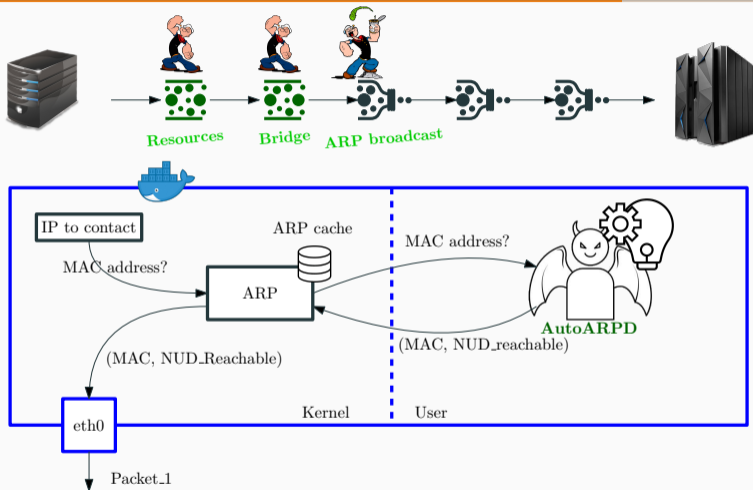
**Usage:** `autoarpd <RULE> [ interfaces ]` **E.g:** `autoarpd 02:42:ip1:ip2:ip3:ip4 eth0`

**AutoArpd.** <https://gitlab.com/uniroma3/compunet/networks/AutoARPD>

## 3<sup>rd</sup> bottleneck: ARP Broadcast

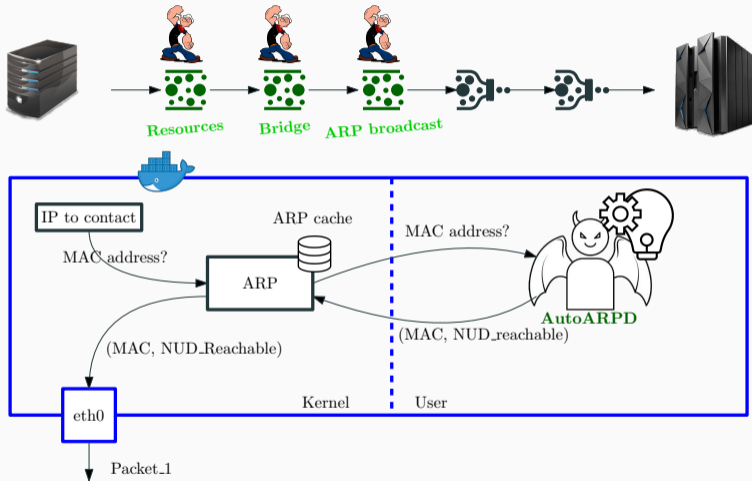


# 3<sup>rd</sup> bottleneck: ARP Broadcast

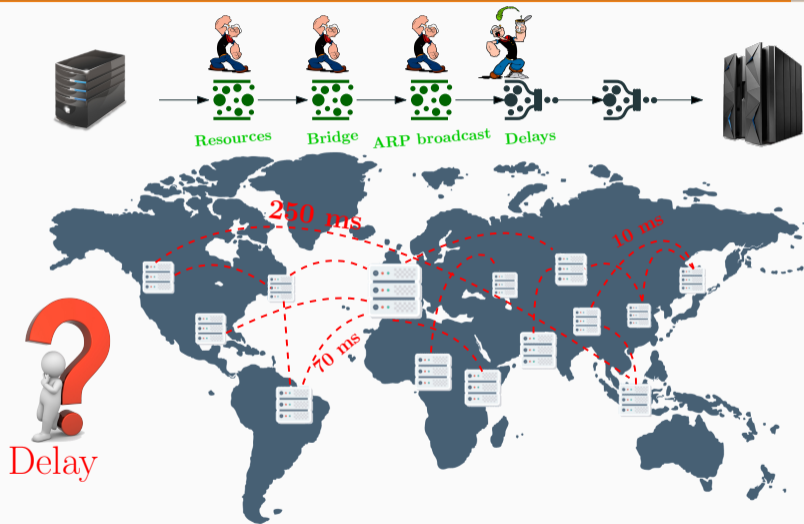




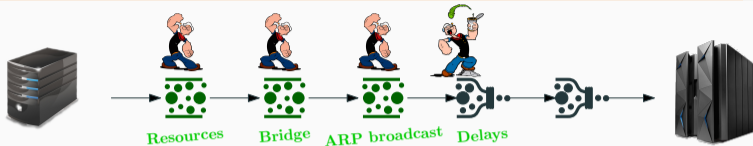
# 3<sup>rd</sup> bottleneck: ARP Broadcast



# 4<sup>th</sup> bottleneck: Emulating Realistic Internet Delays



# 4<sup>th</sup> bottleneck: Emulating Realistic Internet Delays



Measurement-Based Analysis, Modeling, and Synthesis of the Internet Delay Space

Yu Zhang, Tai Sing Eugene Ng, Anirban Nandi, Rishabh R. Rast, Peter Doolittle, and Gabriel Wang

Abstract—Understanding the characteristics of the Internet delay space (i.e., the all-pairs set of state-merged propagation delays across edge networks in the Internet) is important for the design of global-scale distributed systems. The authors characterize the delay space by using a wide variety of validation of the bridge topology and to the graph properties within the Internet delay space. New characterizations of distributed systems allow for an analysis and validation to reach change alternatives. They use a realistic model of the Internet delay space. In this paper, the authors present their specific contributions of distributed systems and quantify key properties that are important for distributed system design. Our results show that existing delay space models do not adequately capture these important properties of the Internet delay space. Furthermore, we derive a simple model of the Internet delay space based on our analytical findings. This model preserves the relevant metrics for better delay modeling, offers for a compact representation, and can be used to synthesize delay data for simulation and estimation of a wide range of distributed systems and applications. We present the design of a publicly available delay space synthesizer and offer D2P and experimental evaluations.

Index Terms—network distributed system, Internet delay space, measurement, modeling, simulation, synthesis.

### 1. INTRODUCTION

DESIGNERS of large-scale distributed systems rely on simulation and network emulation to study design alternatives and evaluate prototype systems at scale and prior to deployment. To obtain accurate results, such simulation or emulation must include an adequate model of the Internet delay space, the all-pairs set of state-merged propagation delays among edge networks. Such a model must accurately reflect these characteristics of real Internet delays that influence system performance. For example, having realistic clustering properties is important because they can influence the load balance of delay-sensitive middle networks and the effectiveness of server placement policies and caching strategies. Having realistic graph characteristics [10] is equally important because the effectiveness of certain distributed algorithms depends on them. Many distributed systems are also sensitive to the multiplicity of IP routing with respect to delay, which influences load-balancing (multiplicity reduction in the delay space and must be reflected in emulated work.

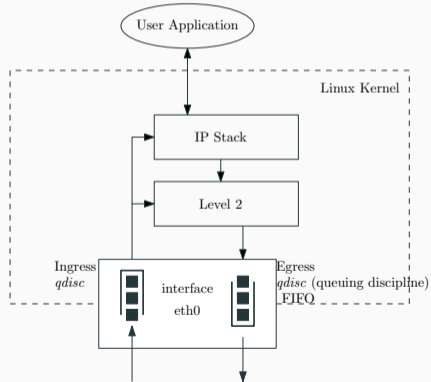
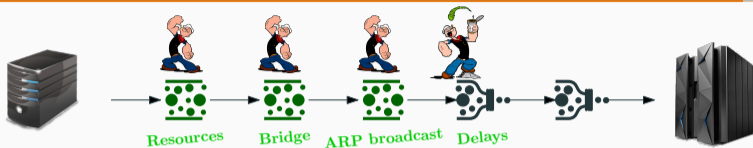
Currently, two approaches are used to obtain a delay model. The first approach, adopted by researchers in the distributed-systems (DS) to collect real delay measurements using a tool such as King [14]. However, due to limitations of the measurement methodology and the quality-time requirement for measuring a delay matrix, measured data tend to be incomplete and their size is the size of a delay matrix that can be captured in practice. To create a delay matrix of size  $1.5 \times 1.5$  million IP space matrix, which is not a trivial amount of data to obtain.

The second approach is to start with a statistical network topology model (e.g., [18], [19], [9], [4], and [17]) and assign individual link delays to the network. The delay space is then modeled by the all-pairs shortest path delays within the topology. The properties of such delay models, however, tend to differ substantially from the actual Internet delay space. This is because these models do not adequately capture rich features in the Internet delay space, such as those caused by geographic constraints, variations in node concentrations, and routing inefficiencies.

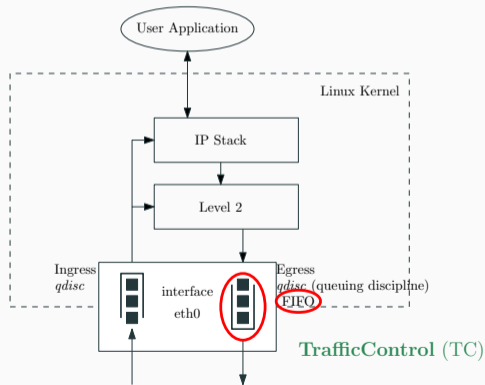
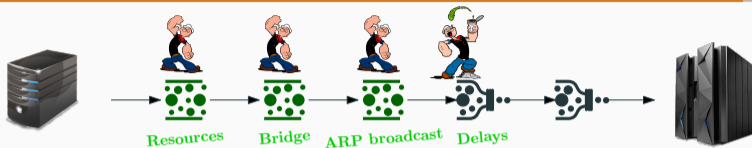
A delay space model suitable for large-scale simulations must adequately capture the relevant characteristics of the Internet delay space. In addition, the model must have a compact representation since large-scale simulation tends to be memory-intensive. The state-of-the-art approach of storing the delay values that consist of a 100 K node network, for instance, requires 100 MB of memory. It is not easy to update the delay values.



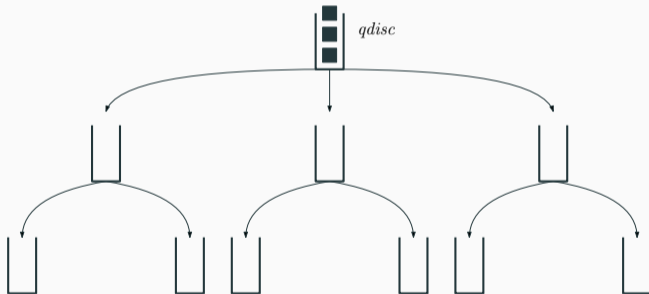
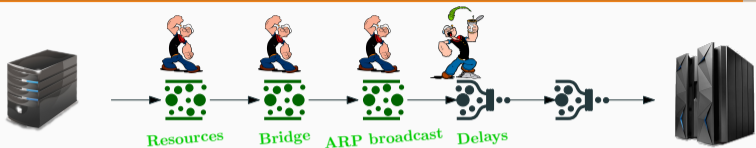
# 4<sup>th</sup> bottleneck: Emulating Realistic Internet Delays



# 4<sup>th</sup> bottleneck: Emulating Realistic Internet Delays

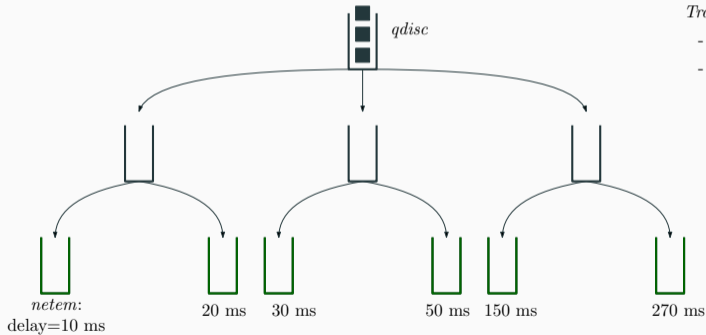
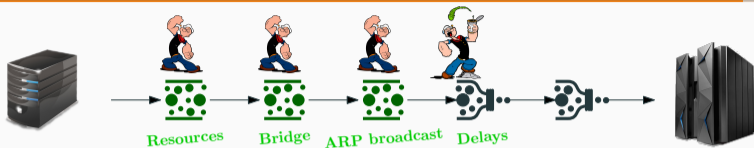


# 4<sup>th</sup> bottleneck: Emulating Realistic Internet Delays



*TrafficControl (TC)*  
- qdiscs structure

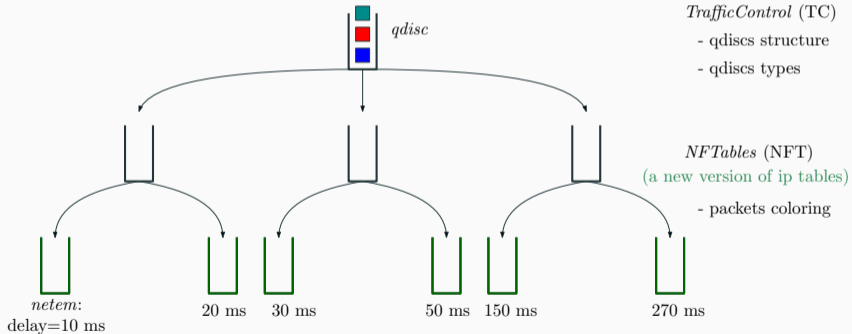
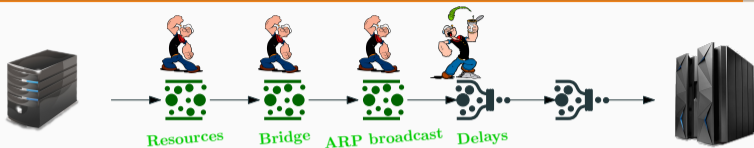
# 4<sup>th</sup> bottleneck: Emulating Realistic Internet Delays



*TrafficControl (TC)*

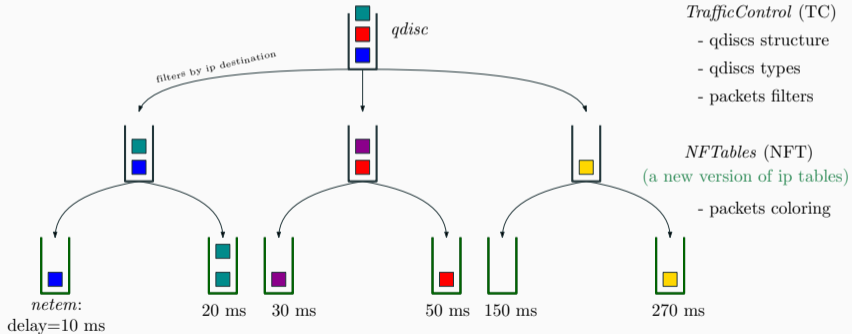
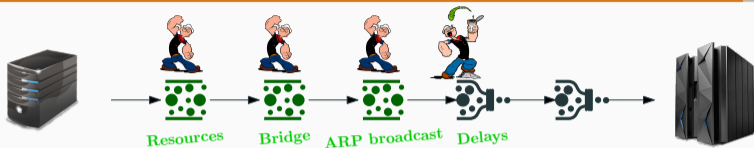
- qdiscs structure
- qdiscs types

# 4<sup>th</sup> bottleneck: Emulating Realistic Internet Delays

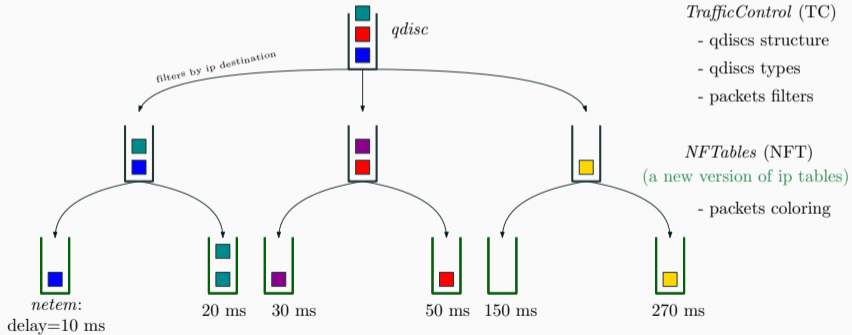




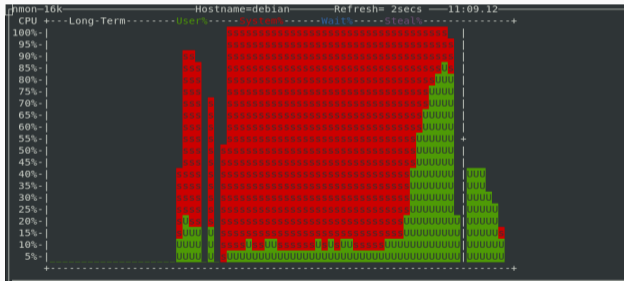
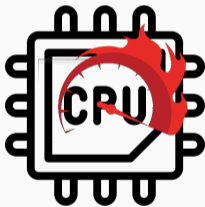
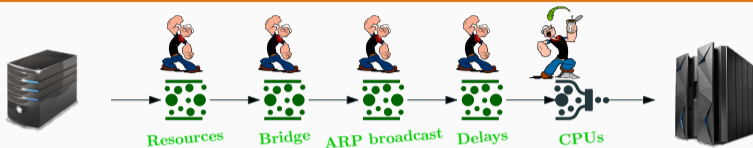
# 4<sup>th</sup> bottleneck: Emulating Realistic Internet Delays



# 4<sup>th</sup> bottleneck: Emulating Realistic Internet Delays

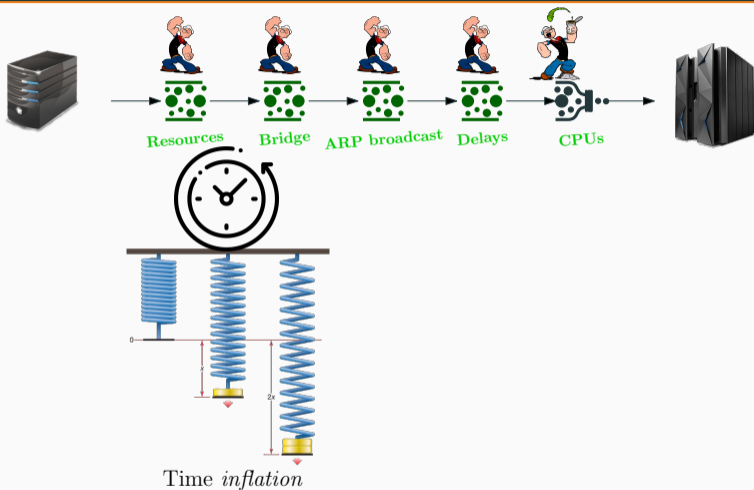


# 5<sup>th</sup> bottleneck: CPUs workload

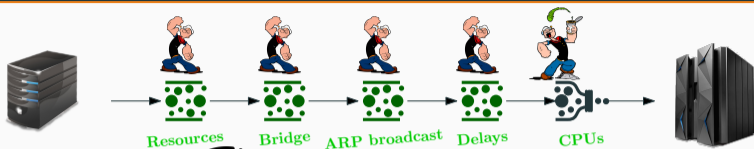


CPUs workload

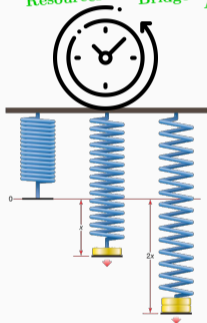
## 5<sup>th</sup> bottleneck: CPUs workload



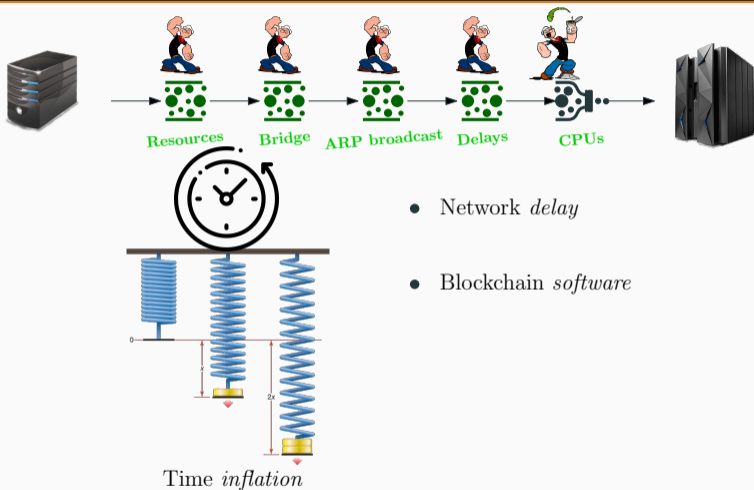
## 5<sup>th</sup> bottleneck: CPUs workload



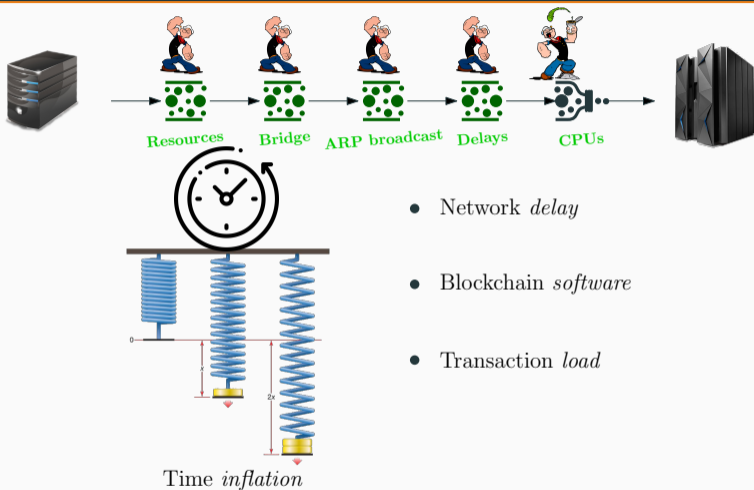
- Network *delay*



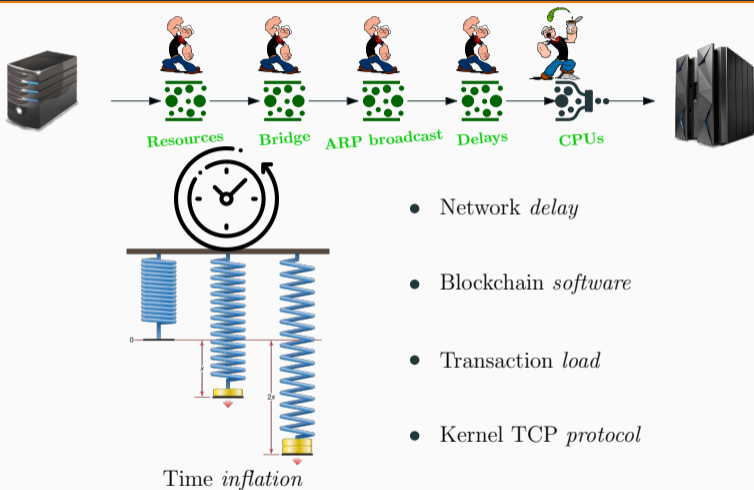
## 5<sup>th</sup> bottleneck: CPUs workload



## 5<sup>th</sup> bottleneck: CPUs workload

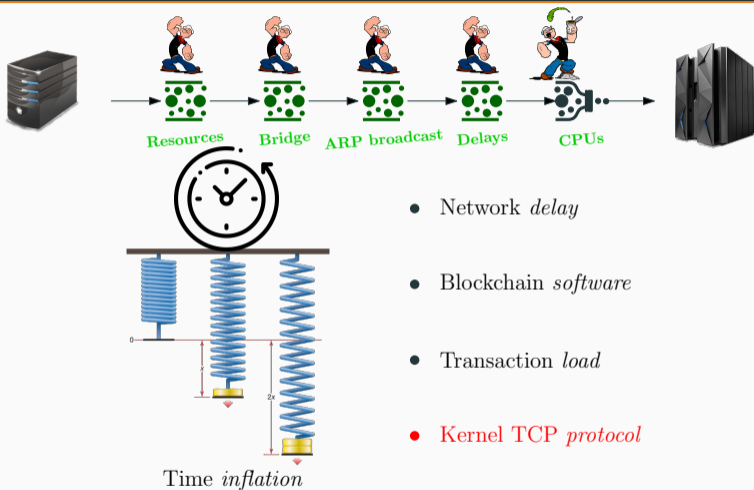


## 5<sup>th</sup> bottleneck: CPUs workload

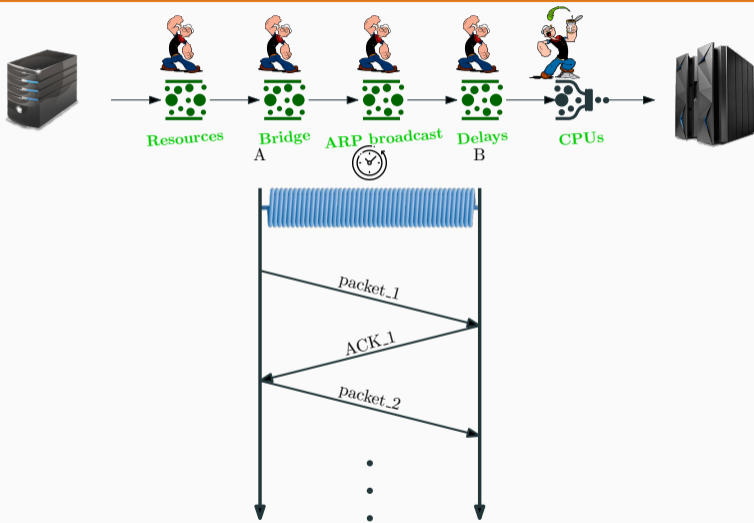




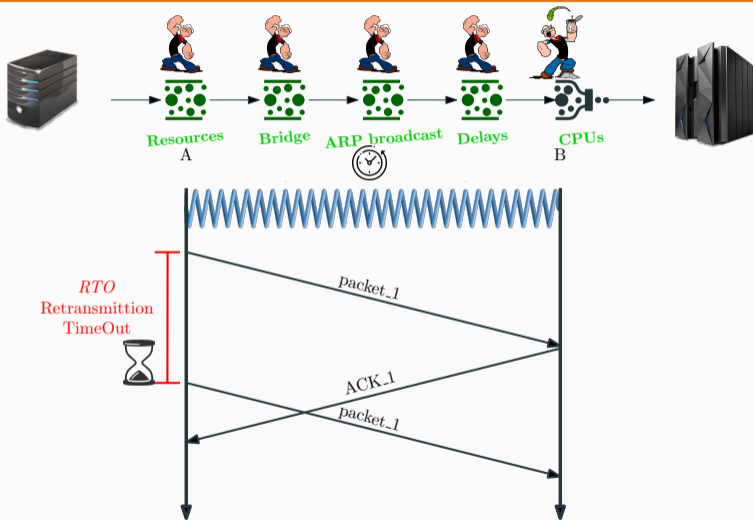
## 5<sup>th</sup> bottleneck: CPUs workload



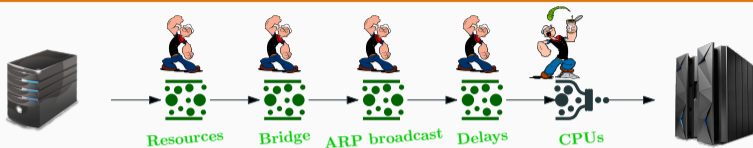
# 5<sup>th</sup> bottleneck: CPUs workload



# 5<sup>th</sup> bottleneck: CPUs workload



## 5<sup>th</sup> bottleneck: CPUs workload

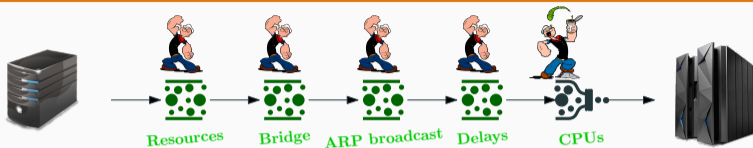


## *Berkeley Packet Filter (BPF) code*

kernel

```
/ include / net / tcp.h
2173 static inline u32 tcp_timeout_init(struct sock *sk)
2174 {
2175     int timeout;
2176
2177     timeout = tcp_call_bpf(sk, BPF SOCK OPS TIMEOUT_INIT, 0, NULL);
2178
2179     if (timeout <= 0)
2180         timeout = TCP_TIMEOUT_INIT;
2181     return timeout;
2182 }
```

## 5<sup>th</sup> bottleneck: CPUs workload

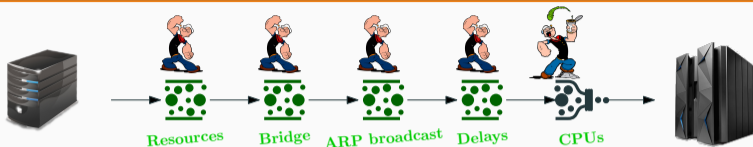


## *Berkeley Packet Filter (BPF) code*

kernel

```
 / include / net / tcp.h
2173 static inline u32 tcp_timeout_init(struct sock *sk)
2174 {
2175     int timeout;
2176
2177     timeout = tcp_call_bpf(sk, BPF SOCK OPS TIMEOUT_INIT, 0, NULL);
2178     -----
2179     if (timeout <= 0)
2180         timeout = TCP_TIMEOUT_INIT;
2181     return timeout;
2182 }
```

## 5<sup>th</sup> bottleneck: CPUs workload



## *Berkeley Packet Filter (BPF) code*

kernel

```
/ include / net / tcp.h
2173 static inline u32 tcp_timeout_init(struct sock *sk)
2174 {
2175     int timeout;
2176     timeout = tcp_call_bpf(sk, BPF_SOCK_OPS_TIMEOUT_INIT, 0, NULL);
2177     if (timeout <= 0)
2178         timeout = TCP_TIMEOUT_INIT;
2179     return timeout;
2180 }
2181
2182 }
```

```
9 __section( "sockops" )
10 int set_initial_rto(struct bpf_sock_ops *skops)
11 {
12     const int timeout = 3; // initial RTO timeout in seconds
13     const int hz = 250; // this value has to match the HZ value of the system
14
15     int op = (int) skops->op;
16     if (op == BPF_SOCK_OPS_TIMEOUT_INIT) {
17         skops->reply = timeout * hz;
18         return 1;
19     }
20
21     return 1;
22 }
```

## 5<sup>th</sup> bottleneck: CPUs workload



## *Berkeley Packet Filter (BPF) code*

kernel

```
/ include / net / tcp.h
2173 static inline u32 tcp_timeout_init(struct sock *sk)
2174 {
2175     int timeout;
2176
2177     timeout = tcp_call_bpf(sk, BPF_SOCK_OPS_TIMEOUT_INIT, 0, NULL);
2178     if (timeout <= 0)
2179         return timeout;
2180     if (timeout <= 0)
2181         timeout = TCP_TIMEOUT_INIT;
2182     return timeout;
2183 }
```

```
9  __section( "sockops" )
10  int set_initial_rto(struct bpf_sock_ops *skops)
11  {
12      const int timeout = 3; // initial RTO timeout in seconds
13      const int hz = 250; // this value has to match the HZ value of the system
14
15      int op = (int) skops->op;
16      if (op == BPF_SOCK_OPS_TIMEOUT_INIT) {
17          skops->reply = timeout * hz;
18          return 1;
19      }
20
21      return 1;
22 }
```

# Our Solution





# Our Solution



## OUR Emulation:

- Huge amount of nodes → 3500 containers in 400GB RAM
- “Real” Network env. → end-to-end realistic internet delays, 8000 TCP-based and 64000 UDP-based connections
- Simple to handle → Makefile and Python scripts
- Simple to modify → Python scripts

# Our Solution



## Future works:

- ★ Simplify the setup
- ★ Multiple host (kubernetes)
- ★ Real software of a blockchain node
- ★ Create a library to create transaction load
- ★ Create a library to support data gathering

**The End**  
**Thank you!**

[diego.pennino@unitus.it](mailto:diego.pennino@unitus.it), [diego.pennino@uniroma3.it](mailto:diego.pennino@uniroma3.it), [pizzonia@ing.uniroma3.it](mailto:pizzonia@ing.uniroma3.it)