

# Evaluation of Message Delay Correlation in Distributed Systems

**Daniel Albeseder**

da@ecs.tuwien.ac.at

Embedded Computing Systems Group

TU Wien, 20.05.2005

## Overview

- Motivation for the  $\Theta$ -Model
- Introduction of the  $\Theta$ -Model
- Evaluation-Setup
- Results
- Conclusion

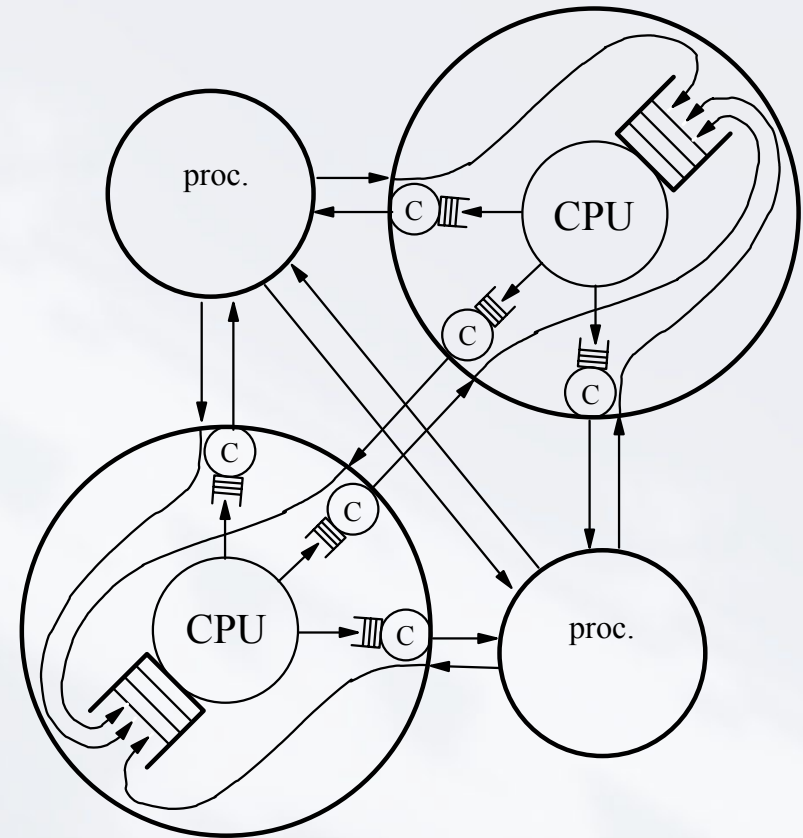
## Timed Algorithms

- Most FT algorithms for distributed RTS have explicit time values (unit „seconds“) in their code / variables
- Toy example: Local real-time clock for timing out a crashed process

```
msg_pong do_roundtrip(msg_ping, p) {  
    send msg_ping to p  
    TIMEOUT :=  $C(t) + 2\tau^+$  /* max. end-to-end delay  $\tau^+$  (sec) */  
    while  $C(t) < \text{TIMEOUT}$  do nothing  
    if msg_pong did not arrive then msg_pong := NIL  
    return msg_pong  
}
```

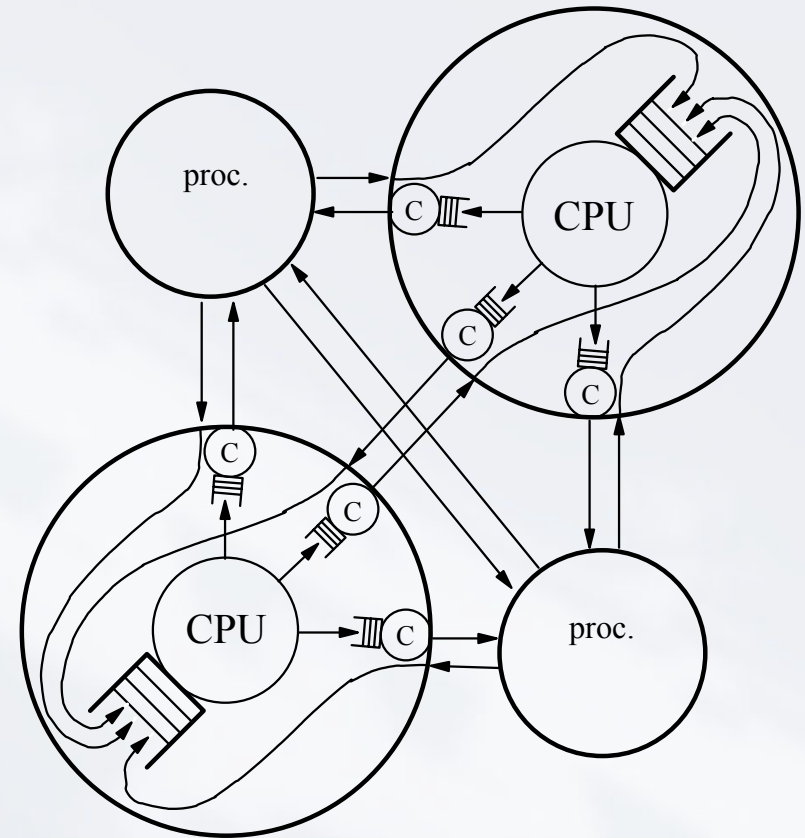
## Why is determining $\tau^+$ difficult ?

- Queuing phenomena:
  - Simultaneous messages from different peers (CPU)
  - Multiple processes (CPU)
  - Multiple messages (Link)

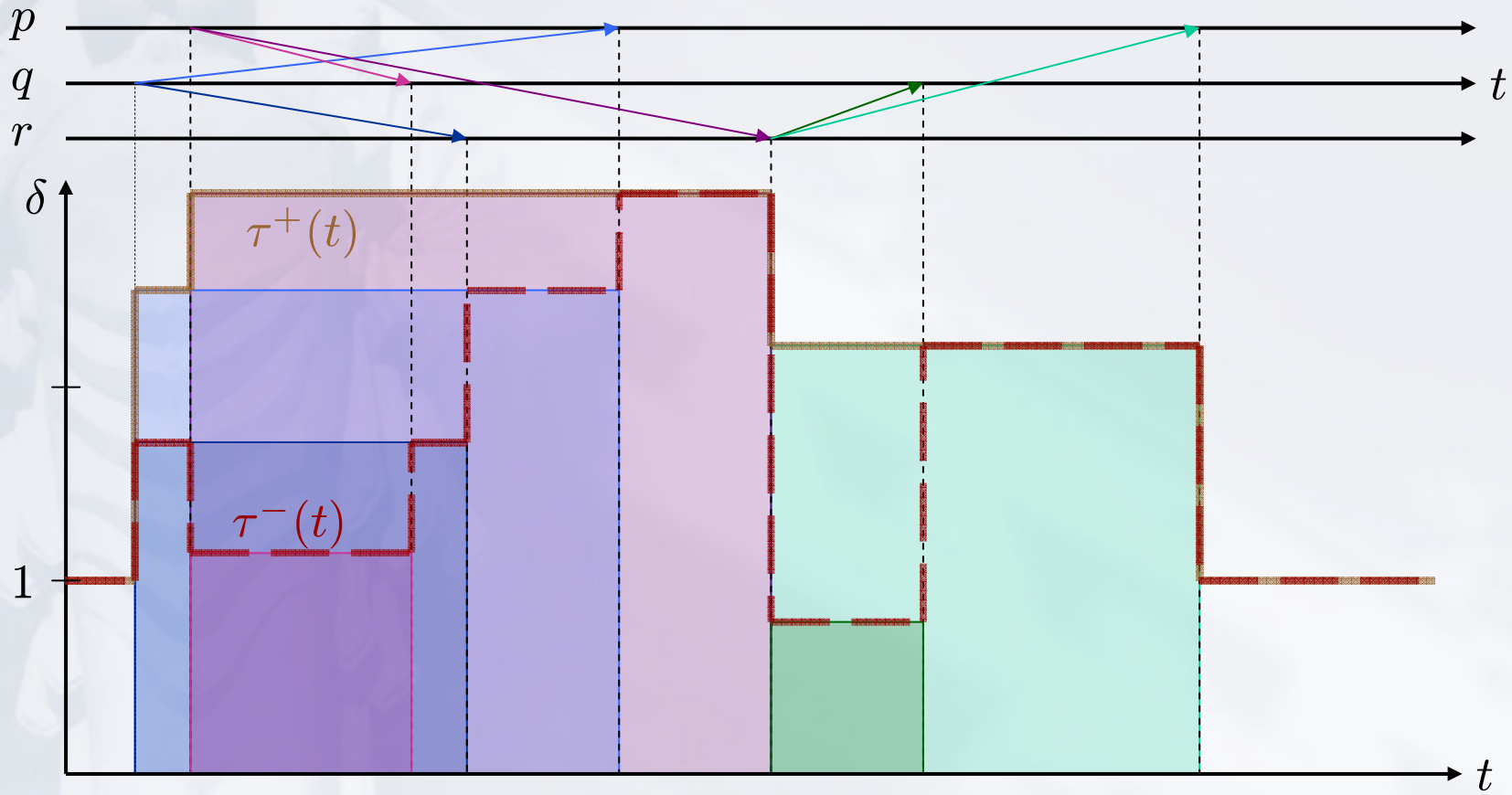


## Why is determining $\tau^+$ difficult ?

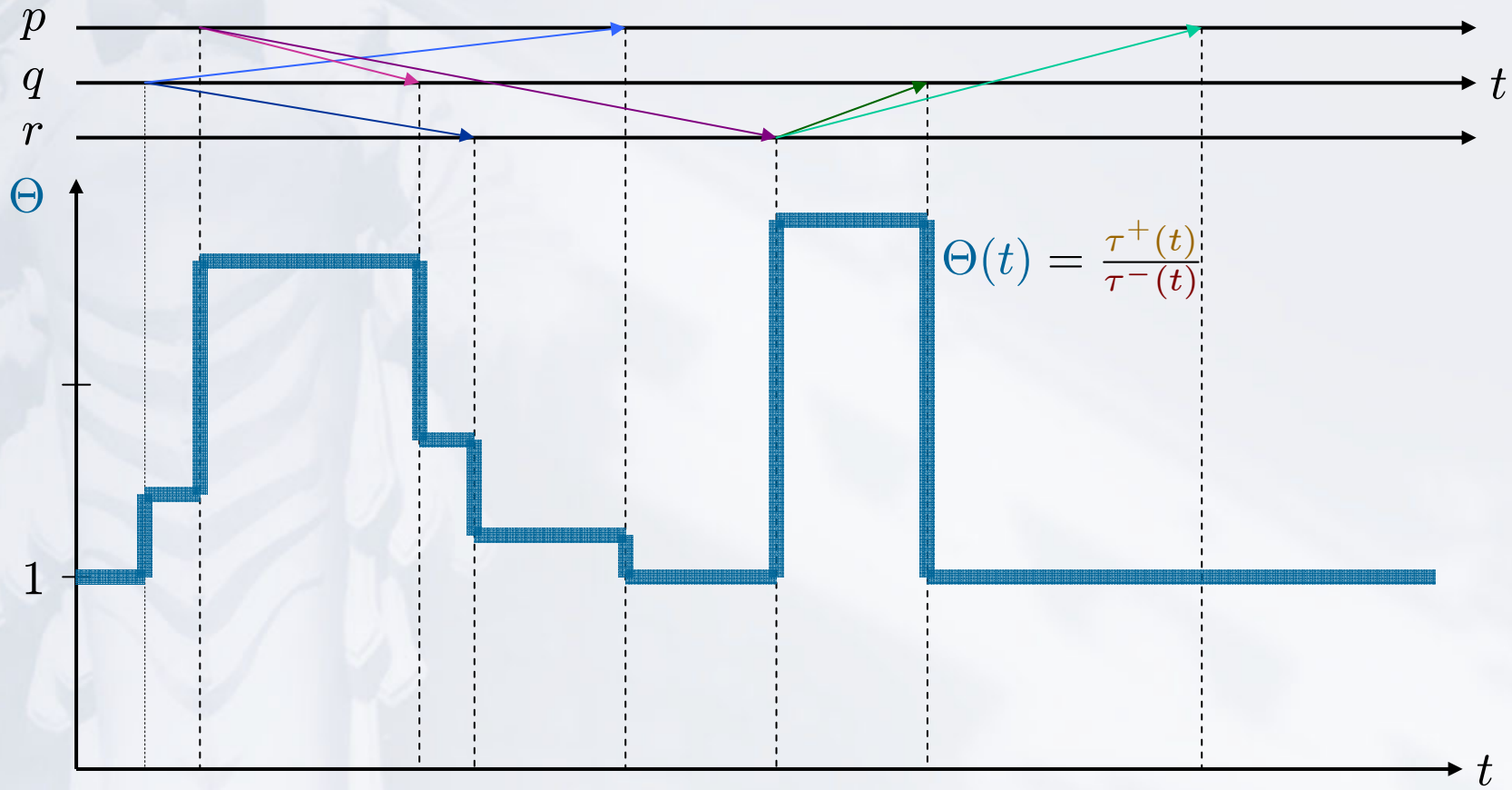
- End-to-end delays hence depend upon
  - message & computational complexity of algorithms
  - interaction („blocking factors“)
  - load conditions
  - **scheduling disciplines**



## The $\Theta$ -Model



## The $\Theta$ -Model

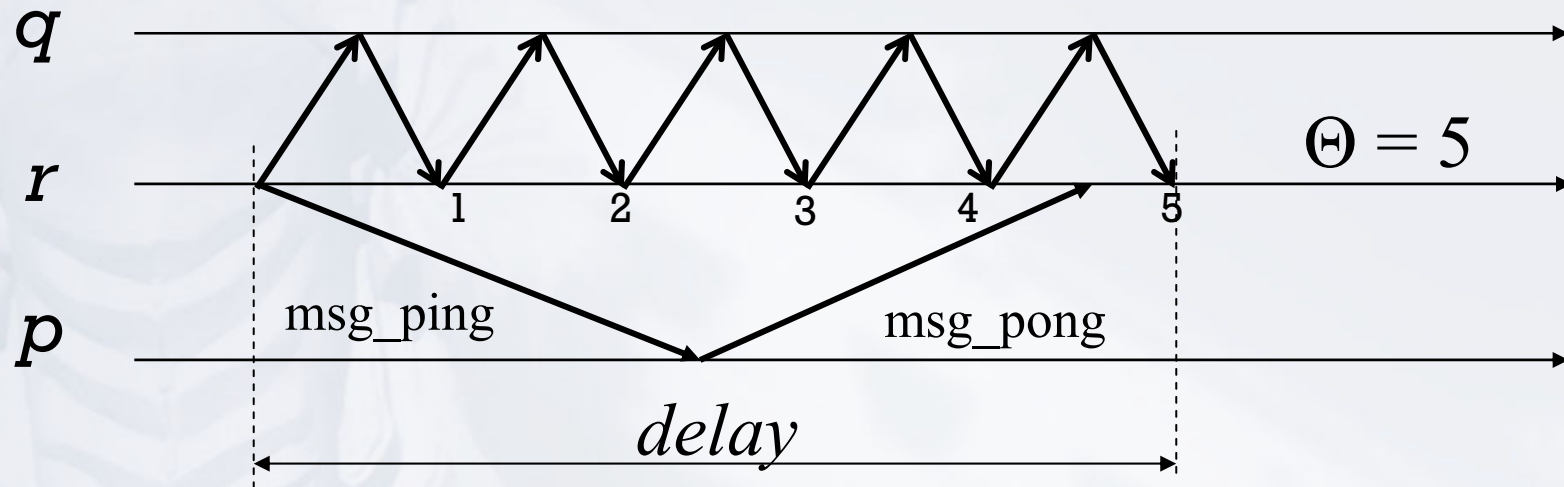


## Time-free implementation of do\_roundtrip(.)

```
send msg_ping to p
for i=1 to  $\Theta$  do      /*  $\Theta$  is dimensionless ! */
begin                  /* do additional roundtrips for waiting */
    send delay_ping(i) to process q
    wait for delay_pong(i) from process q
end
if msg_pong did not arrive then msg_pong := NIL
return msg_pong
```



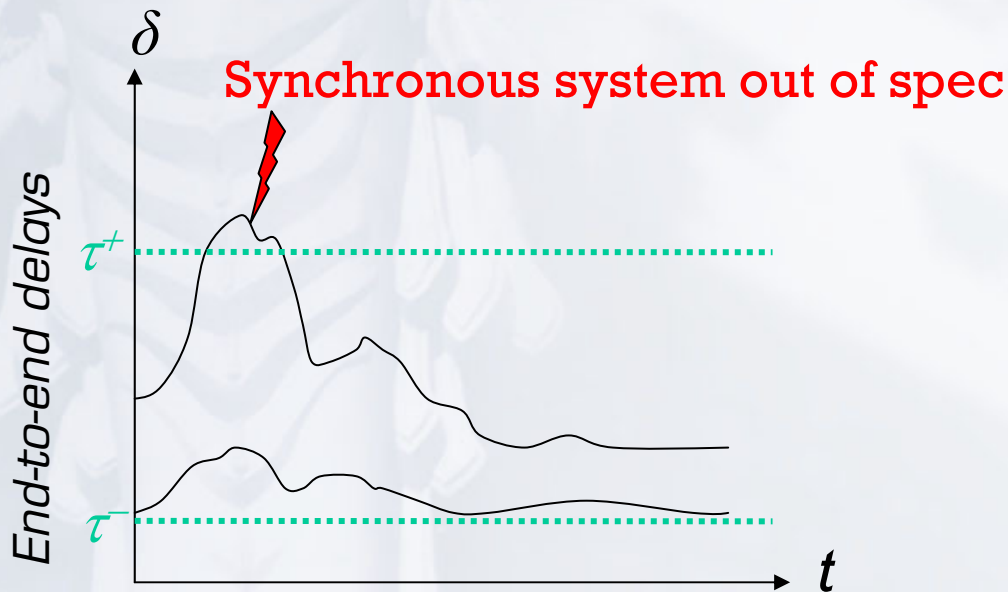
# Time-free implementation of `do_roundtrip(.)`



- Timing behavior solely emerges from the underlying system [*delay* adapts automatically to actual speed]

## Correlation $\rightarrow$ Coverage Expansion

- Given some bound  $\tau^+$  and  $\tau^-$  assumed during system design, compute  $\Theta = \tau^+ / \tau^-$
- Unanticipated overload:  $\tau^+(t) > \tau^+$



if  $\tau^+(t) \leq \Theta \tau^-(t)$   
 however,  
 **$\Theta$ -system still OK**

## Goals

- Is there a correlation between maximum and minimum end-to-end delay?
- Is the minimum end-to-end delay increasing for higher load settings?
- How big is  $\Theta$  in a real system?

A large, faded, light-colored image of a classical statue, possibly a figure with a large, ornate, multi-tiered headdress or wings, serving as a background for the slide.

# Evaluation

## Clock Synchronization Algorithm

```
VAR k : integer := 1;
```

```
send (init, 1) to all [once];           /* initialization */
```

```
if received (init, j) from at least  $f+1$  distinct processes with  $j \geq k$ 
```

```
    k := j;                               /* catch-up with new round */
```

```
    send (init, k) to all [once];
```

```
fi
```

```
if received (init, k) from at least  $n-f$  distinct processes
```

```
    delay(D);
```

```
    k := k + 1;                             /* advance to next round */
```

```
    send (init, k) to all [once];
```

```
fi
```

## Significant Messages

- Not all message are significant for  $\Theta$ -assumption
- Actually only the messages which contribute to a round-switch (RS) are relevant
- Therefore, we define:

$\tau_r^-(t)$  ... n-f fastest message contributing to a RS

$\Theta = \max_t \frac{\tau^+(t)}{\tau_r^-(t)}$  ... Significant uncertainty ratio

## Evaluation Setup

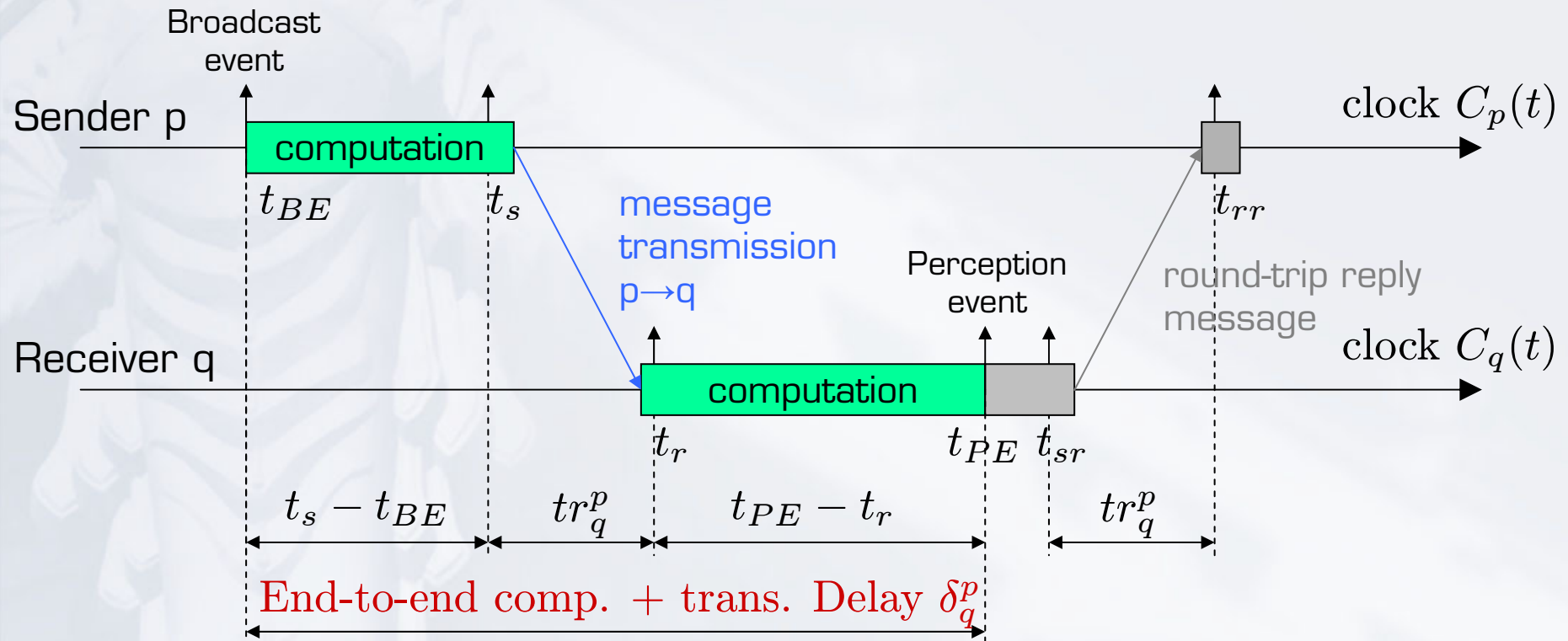
- Pentium4 workstations (2,4GHz FSB533)
- Fully switched Fast-Ethernet over two Cisco Catalyst 2950 switches (connected over fiber Gigabit-Ethernet backbone)
- Red Hat Linux 7.2 with 2.4.20 kernel, patched with High-Resolution-Timers and Kernel-Preemption

## Measurements

- Number of processors  $n=4$  and  $n=7$
  - Failure settings  $f=1$  and  $f=2$   
(without failure injection)
  - Interround-Delays  $D=1\text{ms}$ ,  $D=100\mu\text{s}$  and omitted
  - Various load cases:
    - Constant load,
    - increasing load and
    - load jumps
- } 1. Processor load  
2. Network load  
3. Combined load



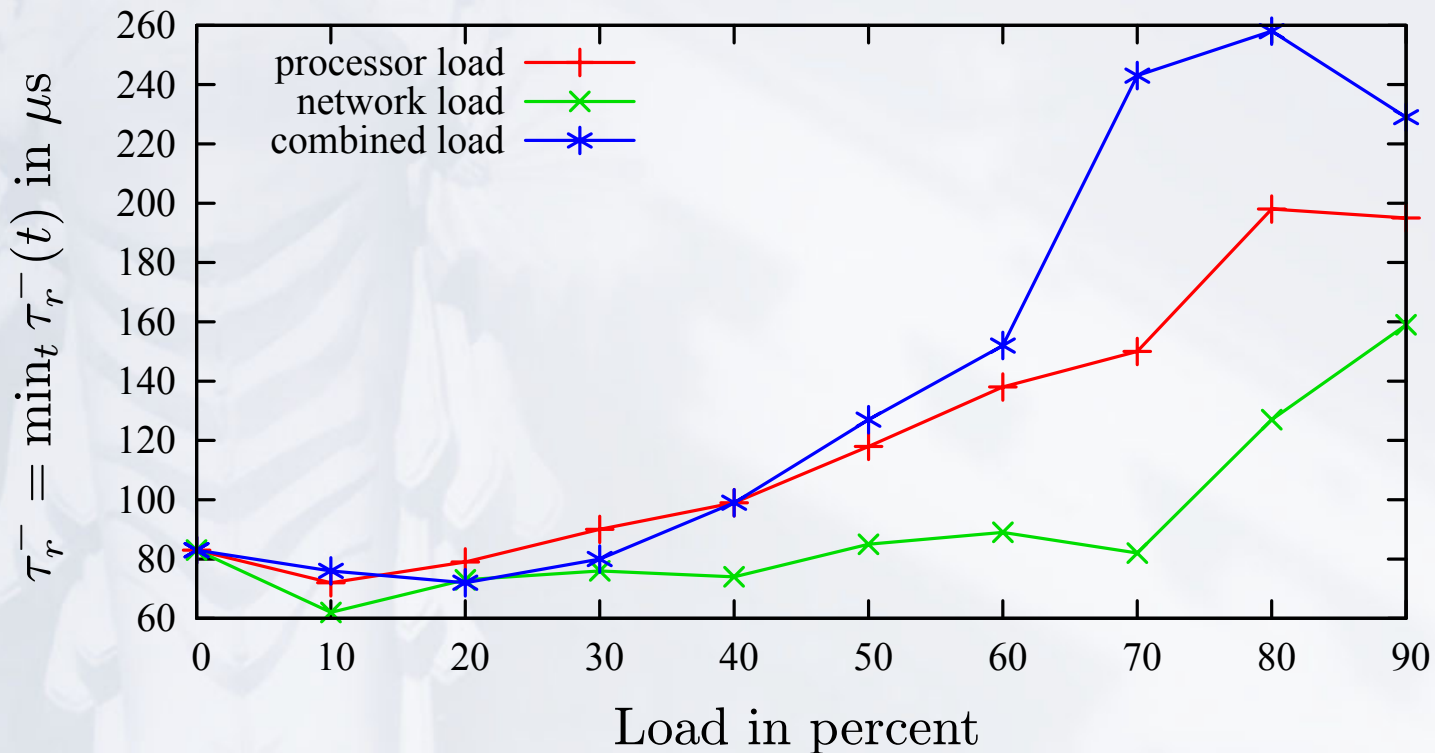
## Round-Trip Measurement



A faded, light-colored background image of a classical building facade. On the left side, there is a large, detailed sculpture of an eagle with its wings spread, perched on a pedestal. The rest of the image shows architectural details like columns and arches, all rendered in a very light, almost white tone.

# Results

## Minimal End-To-End Delay

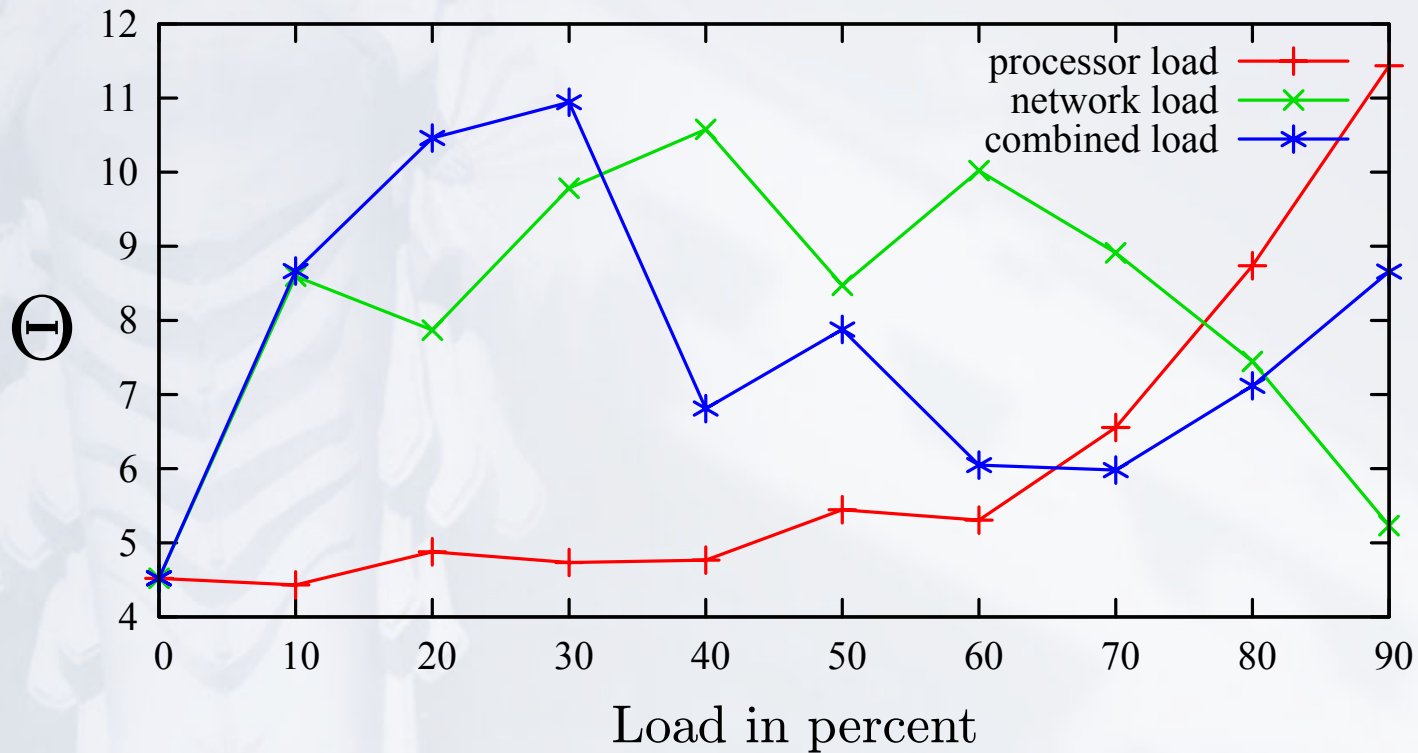


$n = 7$

$f = 2$

$D = 1\text{ms}$

## Uncertainty Ratio $\Theta$



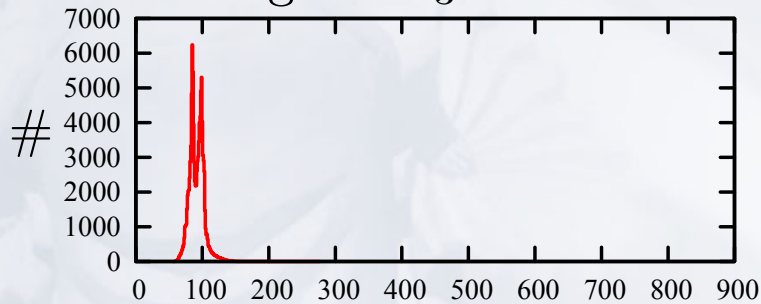
$n = 7$

$f = 2$

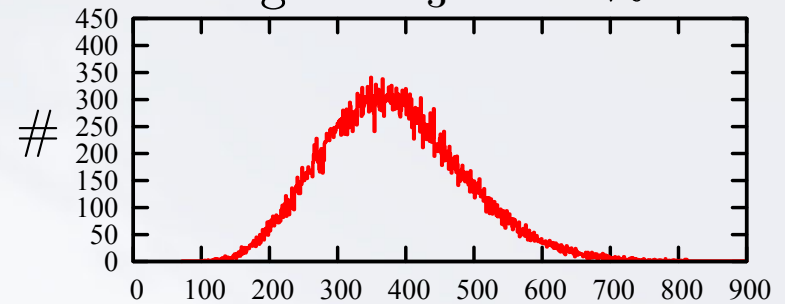
$D = 1m7s$

## Histograms

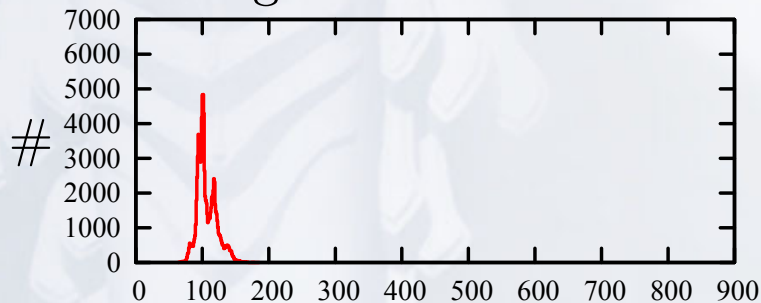
Histogram  $h_3$  for no load



Histogram  $h_3$  for 70% load

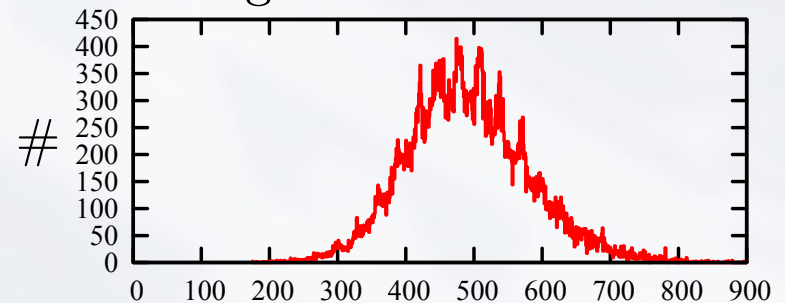


Histogram of  $\tau^+$  for no load



End-to-end delay in  $\mu\text{s}$

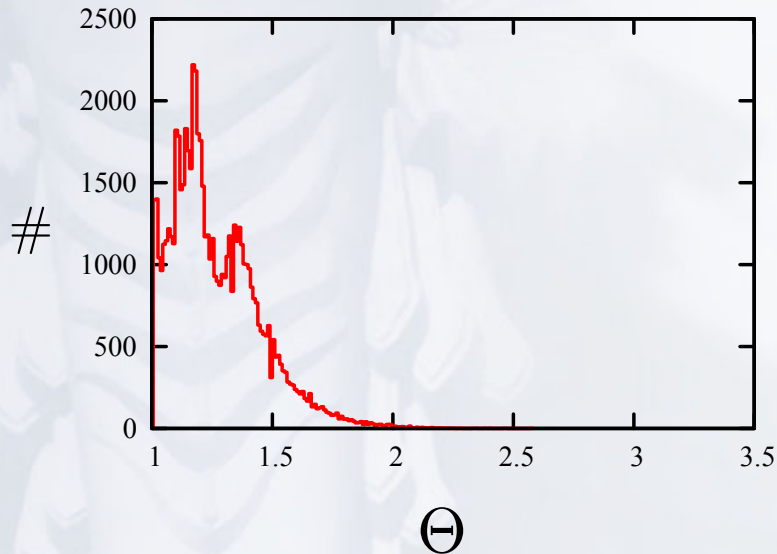
Histogram of  $\tau^+$  for 70% load



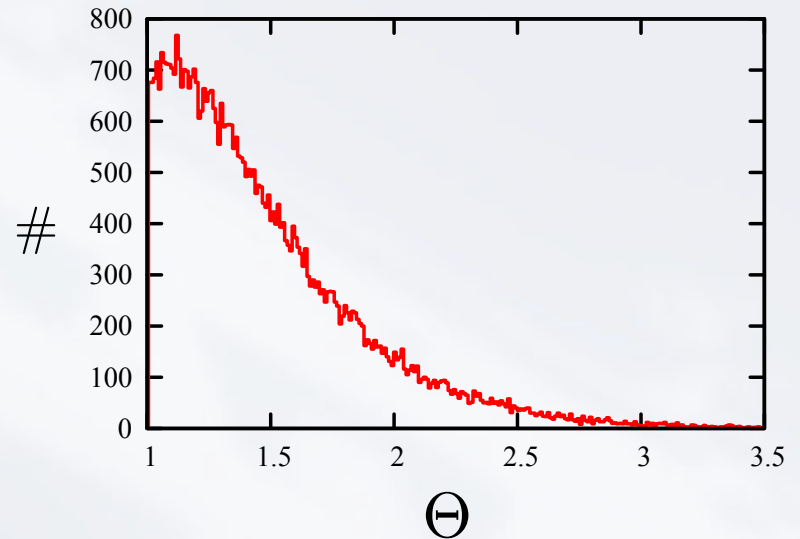
End-to-end delay in  $\mu\text{s}$

## Histograms of $\Theta > 1$

no load



70% load



## Conclusion

- Experimental evaluation of a clock synchronization algorithm
- Common OS (Linux) with enhancements used
- $\Theta$ -Model is valid for our assumptions
- This evaluation is only the beginning of a more comprehensive study

A large, faded, light-colored image of a classical statue, possibly a figure with a large, ornate, multi-tiered headdress or wings, serving as a background for the slide.

**Thank you!**