

PTPv2 clock synchronization for the financial sector

Pedro V. Estrela, PhD
Performance Engineer
13-April-2016



- Part 1: Financial markets overview
 - How electronic markets work
 - A brief history of low-latency trading
- Part 2: Network monitoring
 - Bandwidth measurements
 - Latency measurements
- Part 3: Clock synchronization
 - PTPv2 recap
 - State-of-the-art robustness issues



About the presenter

- Pedro V. Estrela

- PhD in Computer Science (2007 IST-UL)
- Found Financial industry by luck 😊

- Performance Engineer

- “Mechanic” of driver-less Formula 1
- Measure + Remove latency bottlenecks



About IMC

- Think of a currency house, but for:
 - Options / Futures / Stocks / Bonds / ETFs / FX
- Some numbers about IMC
 - All major worldwide Markets, All Timezones, 6 offices, ~500px
 - ~60 datacenters, ~200 links, 10000s equipments
 - ~2000 SW deployments
- Teams' responsibilities
 - Trading team = Find the Price
 - Technology team = Adjust orders Quickly



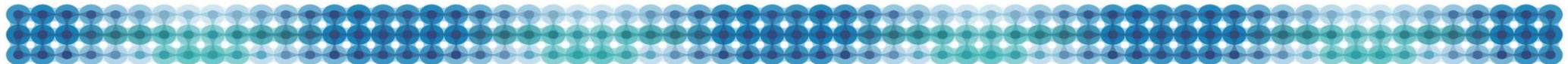
- Competitive landscape



- Relative latency

- $$\text{Total1} = A + \mathbf{B} + C + D + E$$
- $$\text{Total2} = A + B + C + D + E$$
- Trend is clearly: Faster / Raw Hardware / More expensive

Financial Markets overview



WHAT WE DO IN LESS THAN 2 MIN.



<http://www.imc.com/eu/about-us#what-we-do>

• Exchange Price-time priority

- Buyers and Sellers meet at a regulated exchange
- Express their intention to buy / sell
- Orders continuously matched first by price, then time

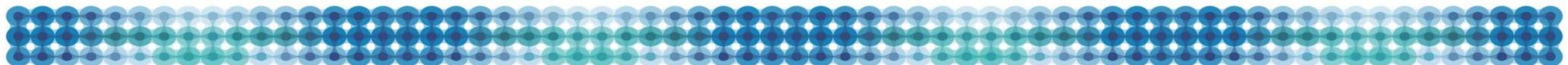


- Imagine this just happened...

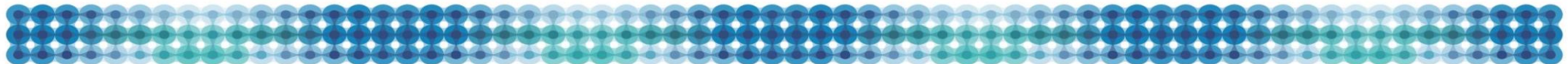
	BUYERS				SELLERS		
Chicago	9.98	9.99	10.00		10.01	10.02	10.03
New York	9.96	9.97	9.98		9.99	10.00	10.01

Questions

- Q1: do you see a trading opportunity here?
- Q2: what should the market maker do here?



Low-Latency



A brief history of low-latency



- 1815: pigeons used to learn of Napoleon's loss
- 1830: Optical-telegraph sends financial messages
- 1865: Fast vessels outrun Mail ships with Union victory news
- 2010: Spread networks drills mountains on NY-Chicago
- 2012: McKay jumps mountains using radio
- 2015: Hibernia builds new straighter Atlantic cable
- 201X:

Days

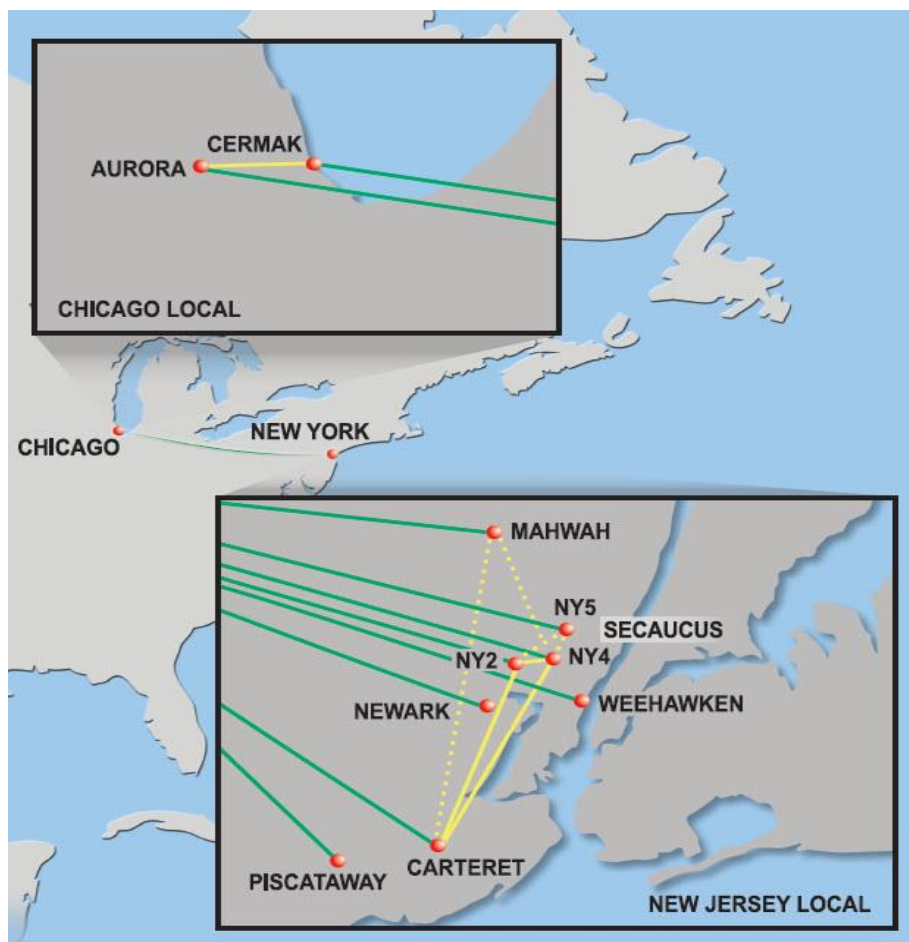
Milliseconds

Source: <http://www.forbes.com/forbes/2010/0927/outfront-netscape-jim-barksdale-daniel-spivey-wall-street-speed-war.html>

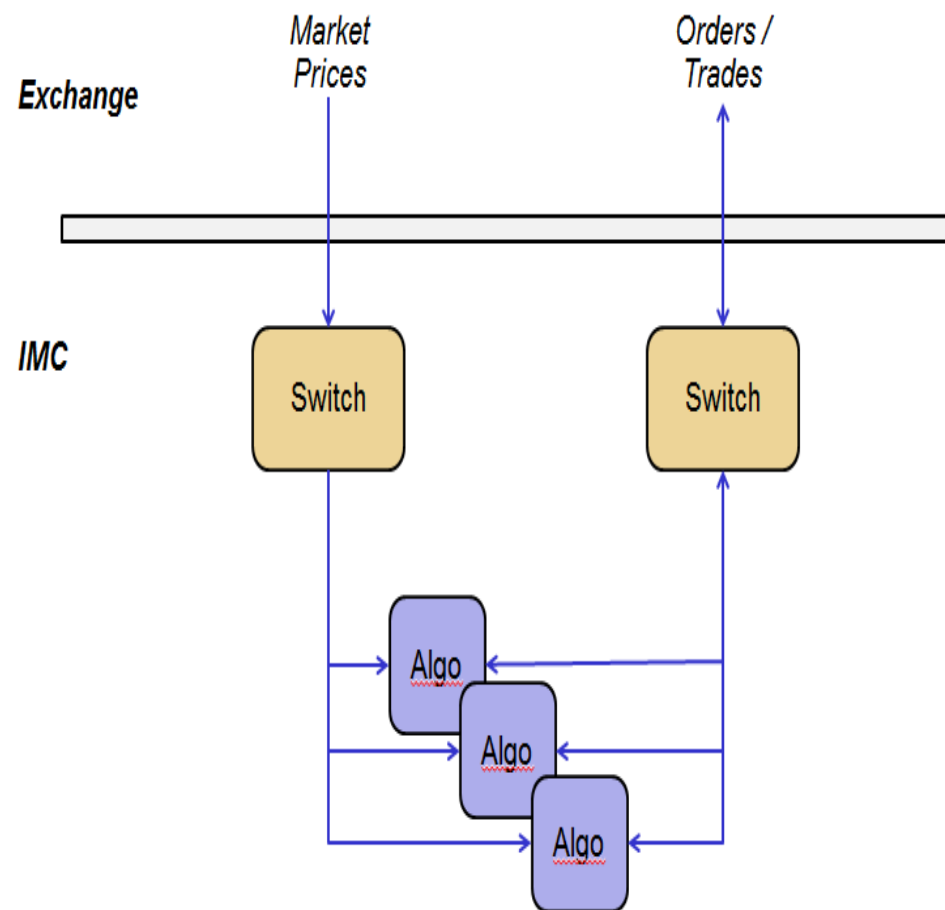
Source: [http://www.moaf.org/publications-collections/financial-history-magazine/111/res/id=sa_File1/Plundered by Harpies.pdf](http://www.moaf.org/publications-collections/financial-history-magazine/111/res/id=sa_File1/Plundered%20by%20Harpies.pdf)



Wide Area

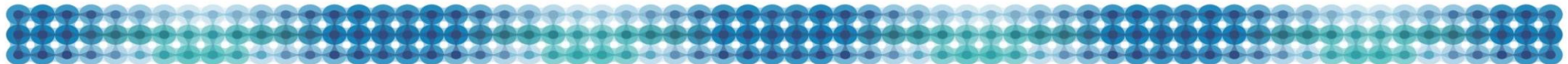
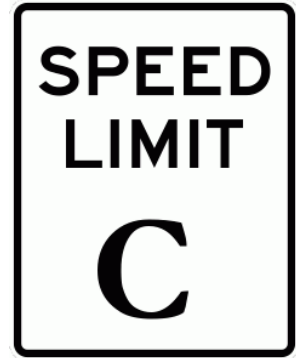


Local Area

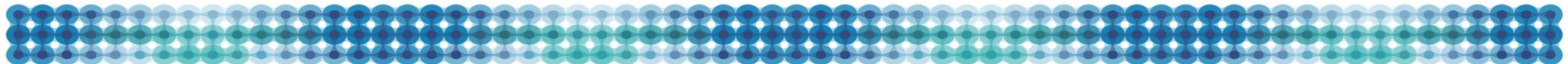


How long is...

- millisecond (ms)
 - A camera flash illuminates for 1 millisecond
 - Distance between countries
- microsecond (μ s)
 - 3 microseconds – Light to travel one Kilometer (1 billion km/h)
 - In and Out a machine, including all processing
- nanosecond (ns)
 - 1 nanosecond – Light to travel one foot
 - 350ns packet forward in a switch

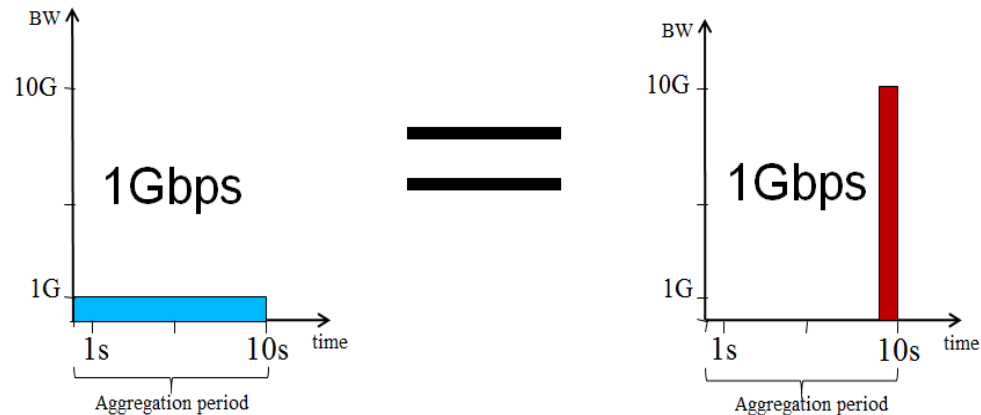


Part 2: Network monitoring



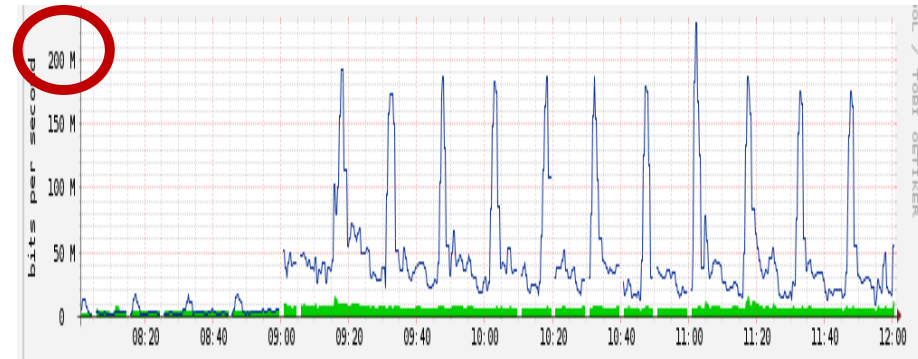
Bandwidth (Micro-Bursts)

$$\text{Rate (Gbps)} = \frac{\text{Volume (bits)}}{\text{Time (seconds)}}$$

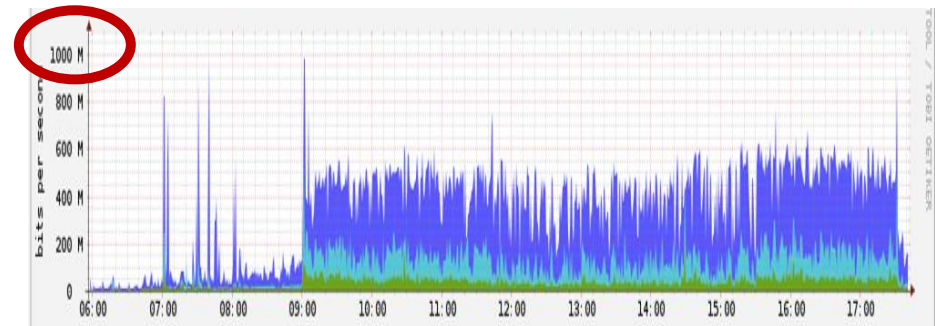


AVG Rate is 1Gbps in both cases!!

- 30 seconds windows:



- 0.010 seconds windows:

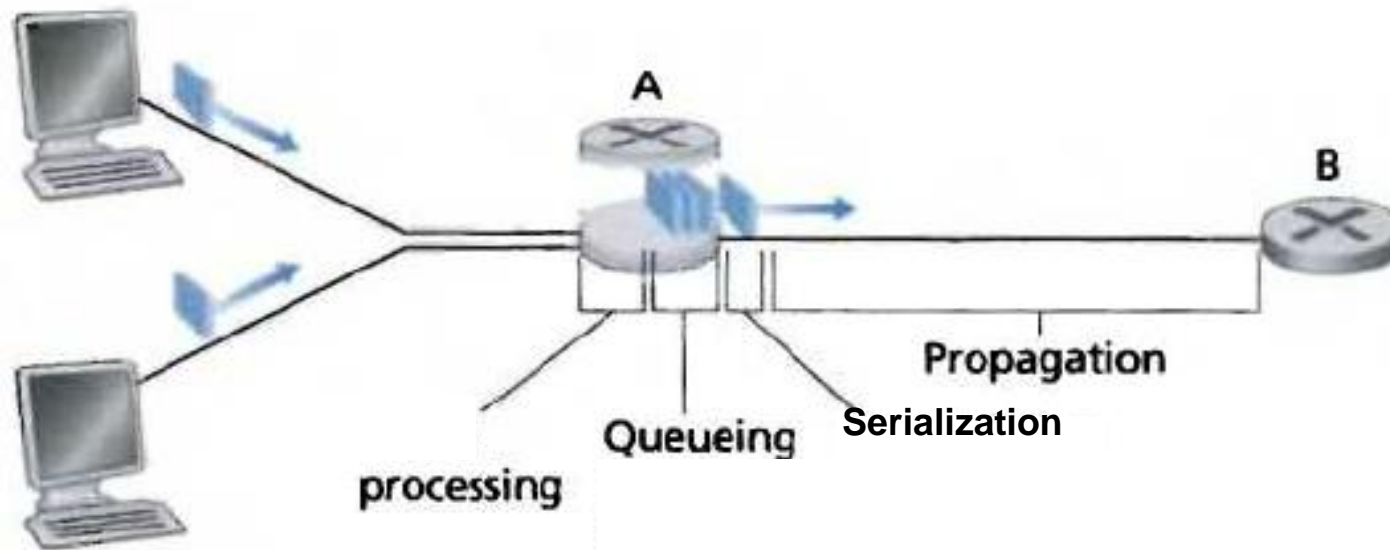


Example: TOP

```
root@tecmint:- http://www.tecmint.com
top - 11:36:04 up 1 day, 22:51, 2 users, load average: 0.06, 0.11, 0.09
Tasks: 141 total, 1 running, 139 sleeping, 0 stopped, 1 zombie
Cpu(s): 0.7%us, 0.5%sy, 0.0%ni, 98.8%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st
Mem: 1021108k total, 982904k used, 38204k free, 134576k buffers
Swap: 2046968k total, 0k used, 2046968k free, 599576k cached
```

PID	USER	PR	NI	VIRT	RES	SHR	S	%CPU	%MEM	TIME+	COMMAND
25454	root	20	0	397m	107m	13m	S	2.3	10.8	25:19.18	skype
269	root	20	0	0	0	0	S	0.3	0.0	0:51.24	scsi_eh_1
1	root	20	0	2872	1400	1200	S	0.0	0.1	0:00.89	init
2	root	20	0	0	0	0	S	0.0	0.0	0:00.00	kthreadd
3	root	RT	0	0	0	0	S	0.0	0.0	0:00.16	migration/0
4	root	20	0	0	0	0	S	0.0	0.0	0:51.23	ksoftirqd/0
5	root	RT	0	0	0	0	S	0.0	0.0	0:00.00	migration/0
6	root	RT	0	0	0	0	S	0.0	0.0	0:00.33	watchdog/0
7	root	RT	0	0	0	0	S	0.0	0.0	0:00.10	migration/1
8	root	RT	0	0	0	0	S	0.0	0.0	0:00.00	migration/1
9	root	20	0	0	0	0	S	0.0	0.0	6:44.34	ksoftirqd/1
10	root	RT	0	0	0	0	S	0.0	0.0	0:00.27	watchdog/1
11	root	20	0	0	0	0	S	0.0	0.0	0:04.05	events/0
12	root	20	0	0	0	0	S	0.0	0.0	0:04.87	events/1
13	root	20	0	0	0	0	S	0.0	0.0	0:00.00	cgroup
14	root	20	0	0	0	0	S	0.0	0.0	0:00.00	khelper
15	root	20	0	0	0	0	S	0.0	0.0	0:00.00	netns
16	root	20	0	0	0	0	S	0.0	0.0	0:00.00	async/mgr
17	root	20	0	0	0	0	S	0.0	0.0	0:00.00	pm
18	root	20	0	0	0	0	S	0.0	0.0	0:00.23	sync_supers

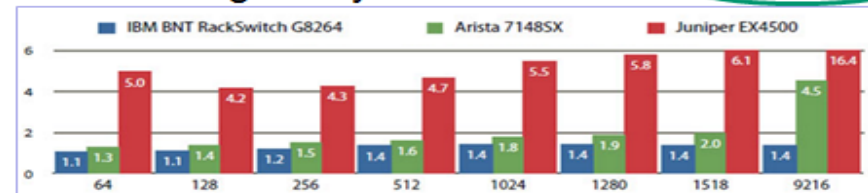
Major latency Factors



Faster switches (\$)

Processing Delay

Constant



Distance (Kilometers)

Propagation Delay

Fixed

Fiber Length	One-Way Delay	Round Trip Delay
1m	5ns	10ns
1km	5µs	10µs
10km	50µs	100µs
100km	500µs	1ms
1,000km	5ms	10ms
10,000km	50ms	100ms

Packet size (bytes)

Serialization Delay

Variable

Packet Size	10Mb/s	100Mb/s	1Gb/s	10Gb/s
64B	51.2µs	5.1µs	0.51µs	0.05µs
512B	0.41ms	41µs	4.1µs	0.41µs
1500B	1.2ms	0.12ms	12µs	1.2µs
9000B	7.2ms	0.72ms	72µs	7.2µs

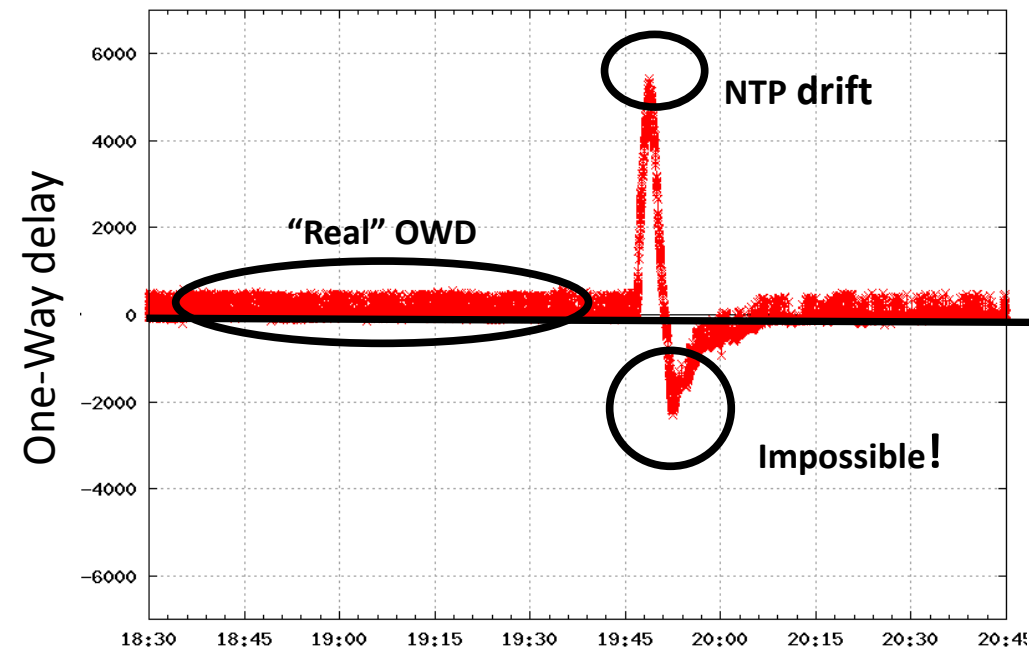
Current Load (Gbps)

Queueing Delay

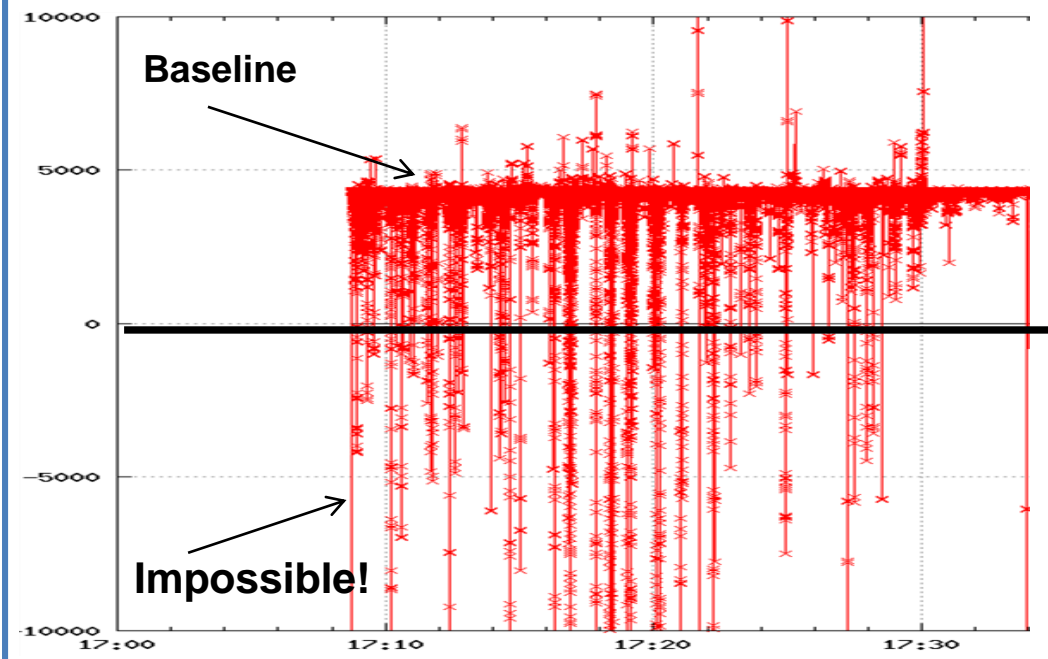
Highly Variable

Queueing Size (IMIX Pkts)	10Mb/s	100Mb/s	1Gb/s	10Gb/s
64 Pkts	17ms	1.7ms	0.17ms	17µs
128 Pkts	33.8ms	3.38ms	0.33ms	33µs
512 Pkts	135ms	13.5ms	1.35ms	0.13ms
1024 Pkts	270ms	27ms	2.7ms	0.27ms

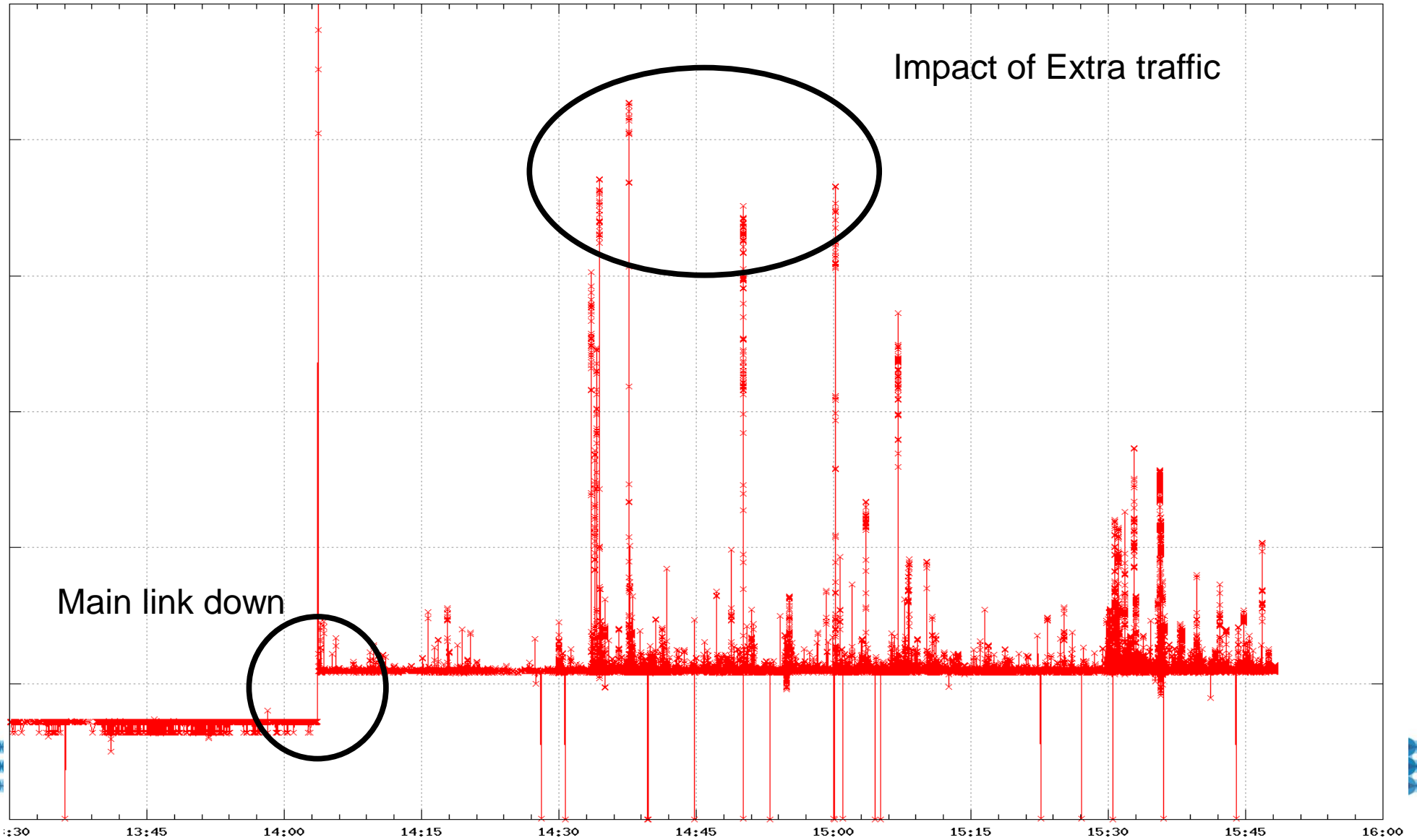
Bad Clock Sync



Bad Timestamps

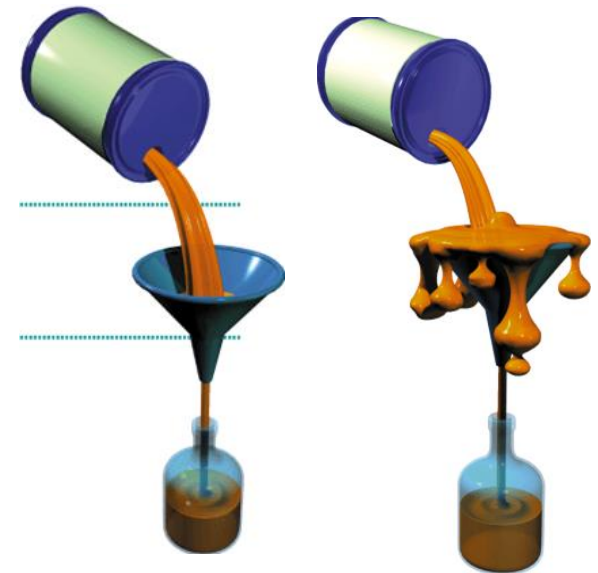
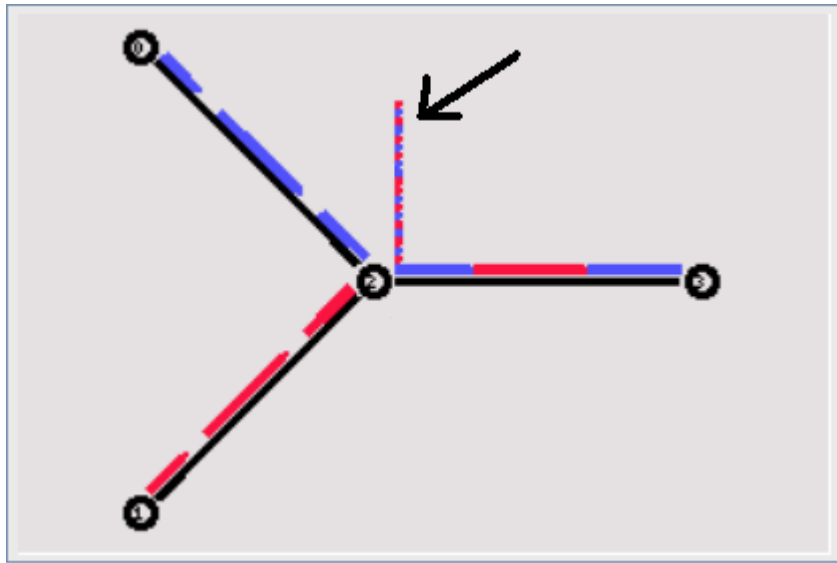


Latency example



Queuing example

- A) Extreme example: <http://youtu.be/D3sEu4CT-nw>
B) Typical example: <http://youtu.be/h-t5XTYNDnU>



Part 3: Increasing robustness of PTPv2 Financial networks



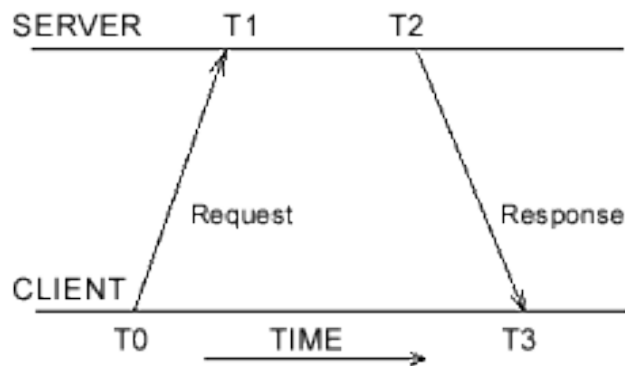
(2014 best paper award)

25 September
ISPCS 2014
Austin, Texas



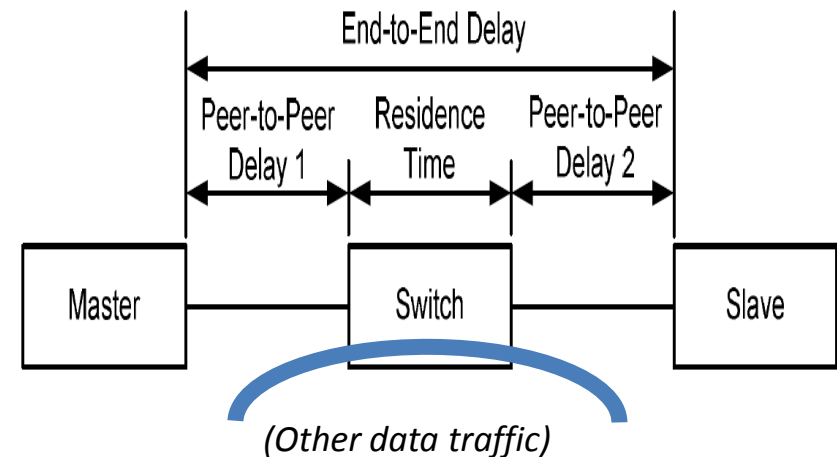
NTP

- Mature IETF standard
- Milliseconds accuracy
- Multiple time sources



PTPv2

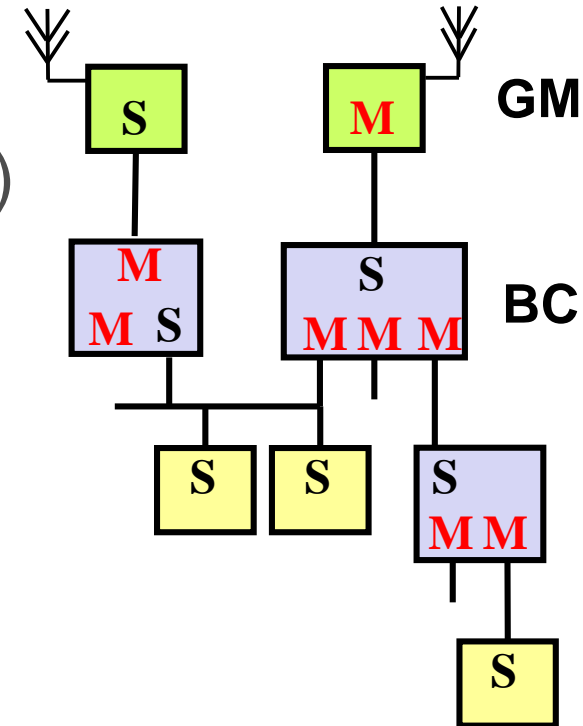
- Recent IEEE standard
- Micro-seconds accuracy
- Single time source



PTPv2 failures in Financial networks

- **Eurex, August 2013**

- Active GM sent *bad* time (leap seconds = 0)
- Backup GMs remain passive
- Slaves jumped by 35 seconds
- Trading halted => all customers affected



- **IMC, July 2011**

- Same problem as above: Single source

Byzantine robustness

- There are always corner cases with a single time source; clients need 3 sources to absorb 1 byzantine failure
- Mathematical proof -> Failure description -> Testlab Proof

1996 - Mathematical proof (Fetzer, Christian)

Integrating External and Internal Clock Synchronization

Christof Fetzer and Flaviu Cristian
Department of Computer Science & Engineering
University of California, San Diego
La Jolla, CA 92093-0114*
e-mail: {cfetzer, flaviu}@cs.ucsd.edu
<http://www-cse.ucsd.edu/users/{cfetzer,flaviu}>

June 4, 1996

Abstract

is a granular representation of real-time and is typically provided by a standard source of time such as NIST. Clocks can be externally or internally synchronized [1]. A clock is *externally* synchronized if at any point in real-time the distance between its value and reference time is bounded by an a priori given constant called *maximum external deviation*. A set of clocks is *internally* synchronized if at any point in real-time the distance between the values of two correct clocks in the set is bounded by an a priori given constant called the *maximum internal deviation* and each clock runs within a linear envelope of real time. Externally synchronized clocks are always internally synchronized by two times the maximum external deviation, but internally synchronized clocks are not always externally synchronized.

We address the problem of how to integrate fault-tolerant external and internal clock synchronization. In this paper we propose a new *external/internal* clock synchronization algorithm which provides both external and internal clock synchronization for as long as a majority of the reference time servers (servers with access to reference time) stay correct. When half or more of the reference time servers are faulty, the algorithm degrades to a fault-tolerant internal clock synchronization algorithm. We prove that at least $2F+1$ reference time servers are necessary for achieving external clock synchronization when up to F reference time servers can suffer arbitrary failures, thus the proposed algorithm provides maximum fault-tolerance. In this paper we also derive lower bounds for the best maxi-

The systems we consider in this paper consist of a

2012 - First failure description (Estrela, Bonebakker)

Challenges deploying PTPv2 in a Global Financial company

Pedro V. Estrela
IMC Financial Markets, Amsterdam, Netherlands
Email: pedro.estrela@imc.nl

Lodewijk Bonebakker
IMC Financial Markets, Amsterdam, Netherlands
Email: lodewijk.bonebakker@imc.nl

Abstract—This paper describes the challenges encountered when deploying PTPv2 on the worldwide network of a financial company, by upgrading nearly all servers in all data-centers over a period of two years, to achieve global microsecond level accuracy between any pair.

Acknowledging that PTP was initially designed as a LAN protocol and that all current time-keeping industry efforts are focused on PTP, the issues can be broadly divided into a) issues on the PTPv2 standard itself, b) issues that have to be addressed when PTP is expanded to work over WANs, and c) issues that caused the biggest operational impact on the (tested) implementations.

In all, this paper contributes concrete examples where PTP's byzantine robustness, scalability and efficiency characteristics range between absent to poor – and attempts to raise awareness on the steps needed to build PTP solutions with the characteristics that global users want.

I. INTRODUCTION

This paper describes the challenges encountered when deploying PTPv2 on the worldwide network of a financial company, in order to achieve microsecond level accuracy between any two servers (globally). For this, we will describe the issues discovered over the last two years, while deploying

Table I
A SUMMARY OF THE ACRONYMS USED IN THIS PAPER

ACL	Access Control Lists
BC	Boundary Clock
BMC	Best Master Clock
DC	Data-Center
FINRA	Financial Industry Regulatory Authority
GM	GrandMaster
IGMP	Internet Group Management Protocol
LAN	Local Area Network
MAN	Metropolitan Area Network
NE	Network Equipment
NIC	Network interface controller
NTP	Network Time Protocol
PDM-SM	Protocol Independent Multicast - Sparse Mode
RP	Rendezvous Point
TTL	Time To Live
UTC	Universal Coordinated Time

Taking these considerations into account, this paper divides the encountered issues into a) those that affect PTPv2 as it is defined today (i.e., for LANs), b) the issues that have to be addressed when PTP is expanded to work over WANs and c) the issues that caused the biggest operational impact on the (tested) implementations. In all, this paper attempts to raise

2014 - First proof + solution: (Estrela, Neusuess, Owczarek)

Using a multi-source NTP watchdog to increase the robustness of PTPv2 in Financial Industry networks

Pedro V. Estrela
IMC Financial Markets
Amsterdam, Netherlands
pedro.estrela@imc.nl

Sebastian Neusuess
Deutsche Börse AG
Frankfurt, Germany
Sebastian.Neusuess@deutsche-boerse.com

Wojciech Owczarek
NYSE Euronext
Belfast, UK
wowczarek@nyx.com

Abstract— This paper describes a fundamental single point of failure in the PTPv2 protocol that affects its robustness to failure in specific error scenarios. The architecture design of electing a single unique time source to a PTP domain – the PTP GrandMaster – makes this protocol vulnerable to byzantine failures.

Previous work has described this vulnerability from both a theoretical and practical point of view – and in particular how this affects the financial industry. This paper advances the discussion by contributing a description of the latest high-accuracy regulatory requirements on the financial industry, and by documenting new examples of failures in real-world customer-facing operations. It then describes an example of one of possible ways to increase PTP robustness while preserving its accuracy (using a multi-source NTP watchdog), and a laboratory test that shows how different protocol implementations are affected by this problem.

In all, the current paper attempts to raise awareness of the robustness requirements within the financial industry today. As only PTP is accurate enough for both current and upcoming regulatory requirements, we hope that these issues are addressed in the forthcoming PTPv3 standard based on a multi-source time source.

fundamental single point of failure that renders this protocol vulnerable to “byzantine failures” – the worst possible class of failures where failing GMs do not shutdown, but instead start to send misleading time information to their slaves.

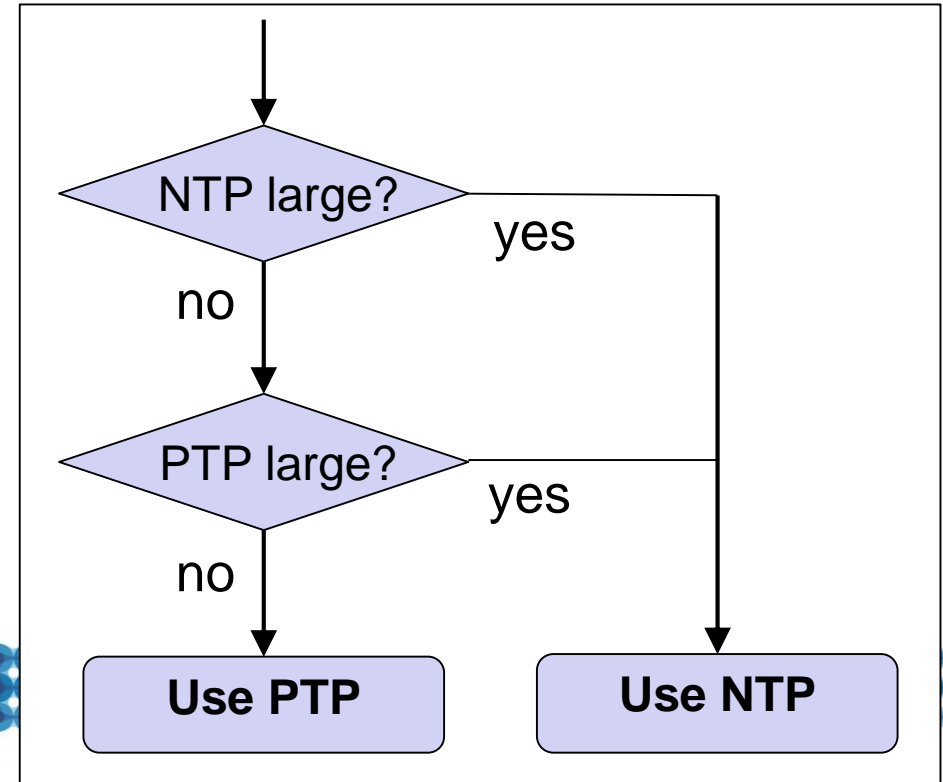
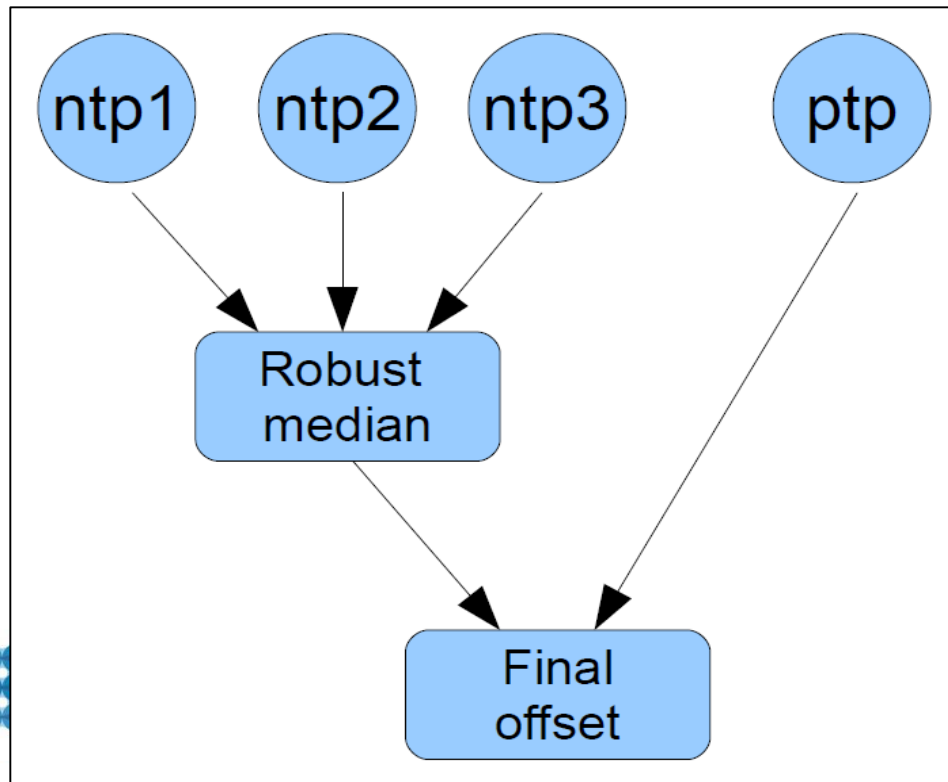
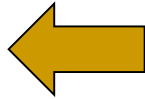
Previous work has described this exact vulnerability from both a theoretical [2] and practical point of view [3] – and in particular how this affects the financial industry [4].

To advance the discussion, this paper makes the following contributions:

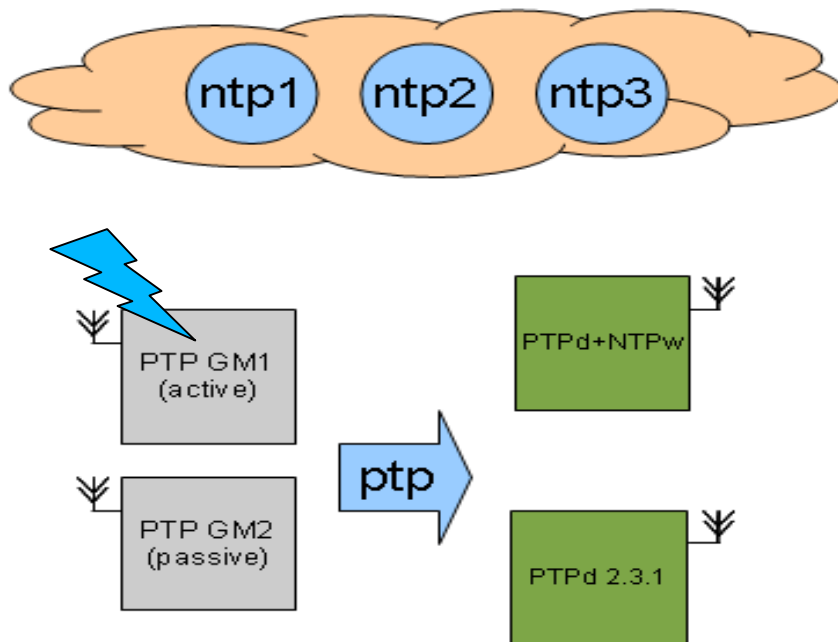
- a description of the latest regulatory requirements that are pushing higher accuracy obligations to the financial industry ([11]/[13]/[15])
- a description of new examples of failures in real-world customer-facing operations [10]
- an example of one of the possible ways to increase PTP robustness while preserving its accuracy (using a multi-source NTP watchdog to prevent failure scenarios)

Solution, using with NTP watchdog

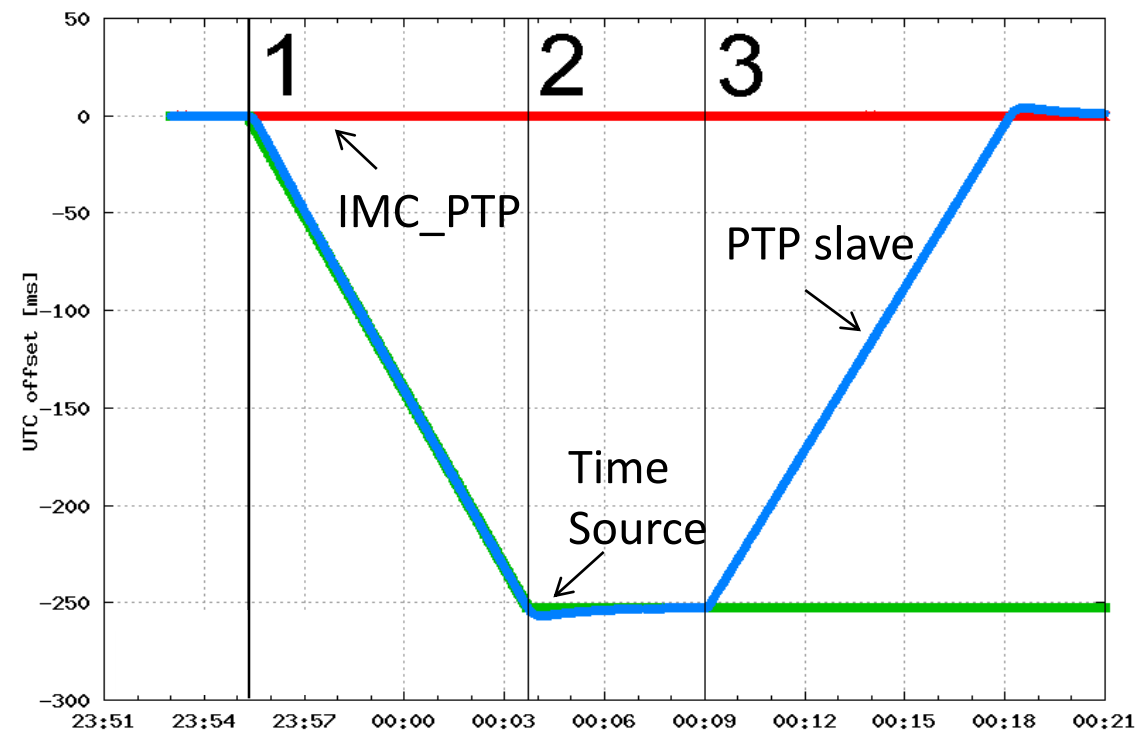
- NTP servers queried in parallel to PTP
- Robust median offset can override PTP offset:
 - -0.2 ms
 - +0.1 ms
 - +35000 ms
- PTP only touches the clock if allowed by the NTP watchdog



Testbed



Clock Error attack



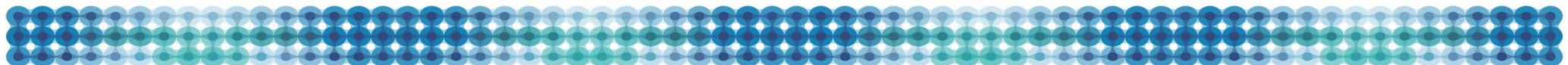
Conclusion

- IMC opportunities

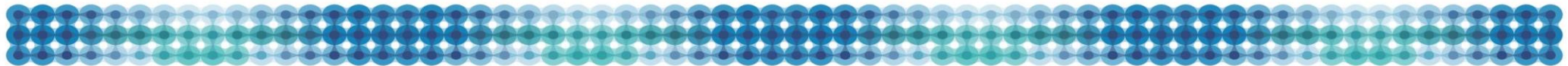
- Information Technology
 - Quantitative Trading
 - Both Internships and Full time opportunities
- } **Very heavy CS needs!**

- More questions?

- IMC: <http://www.imc.nl/yourcareer/opportunities>
- Scientific papers: <http://tagus.inesc-id.pt/~pestrela/ptp>

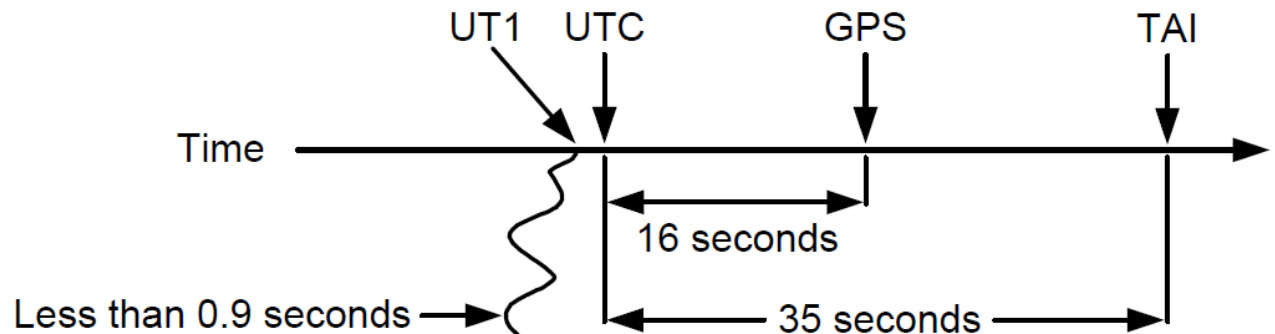


Extra Slides

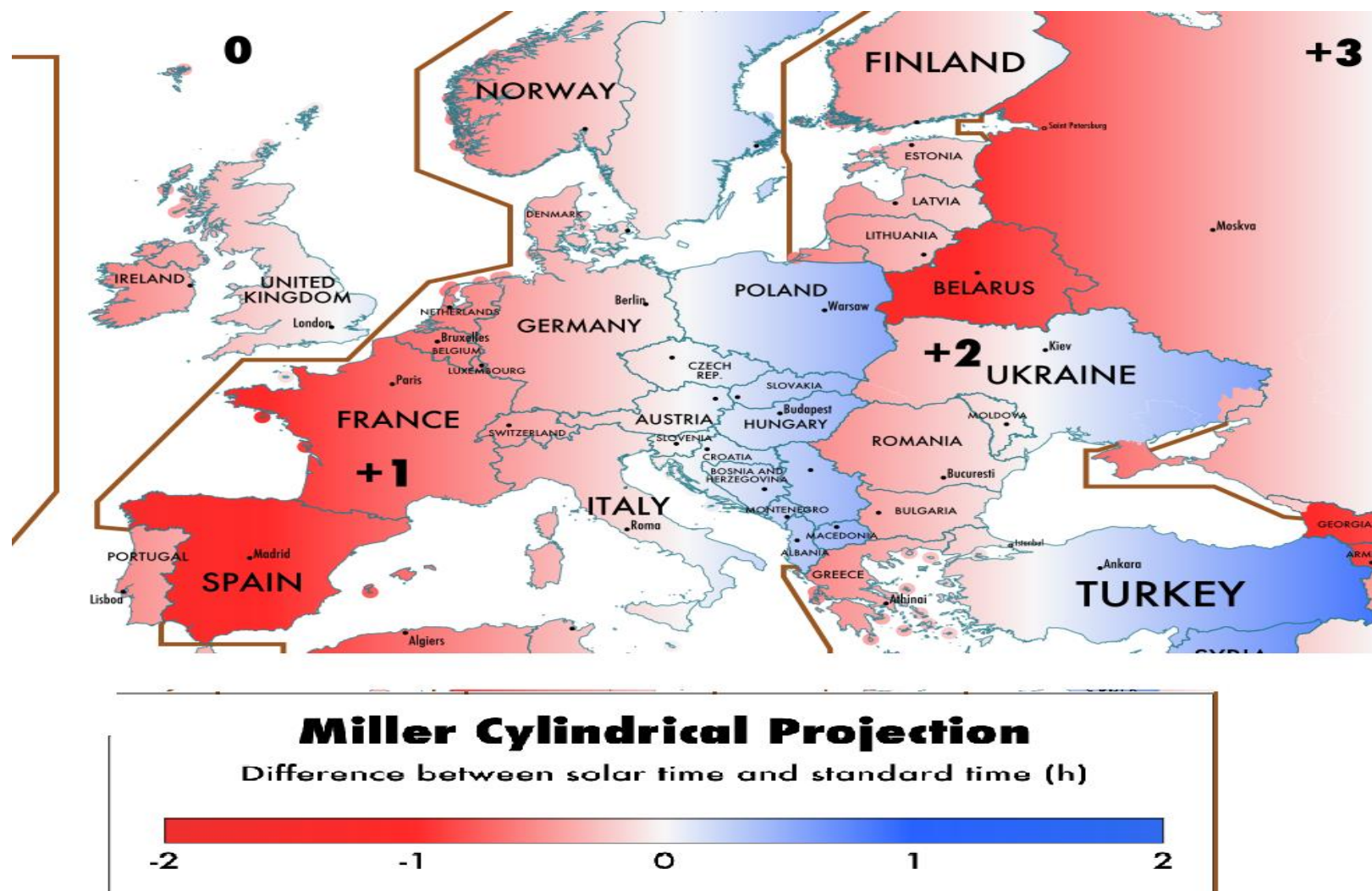


Leap seconds = Problems

- **Heraldsun:**
 - “Leap second crashes Qantas and leaves passengers stranded”
- **Cnet:**
 - “Leap second bug causes site software crashes”
- **Globalpost:**
 - “Weird Wide Web - Leap second causes flight delays and internet problems”
- **Buzzfeed:**
 - “How a second brought down half the Internet”
- **Wired:**
 - “Leap second glitch explained”



UTC - Solar Time @ noon



Source: <http://blog.poormansmath.net/the-time-it-takes-to-change-the-time>

Solution 1: no “0” leap seconds

- Leap seconds != “0”
 - Earth has been slowing-down, so we’ve been adding leap seconds
 - Banning (UTC valid = “1” / UTC offset = “0”) would avoid some of the problems for hundreds, if not thousands, of years

Solution 2: Fixed Leap seconds

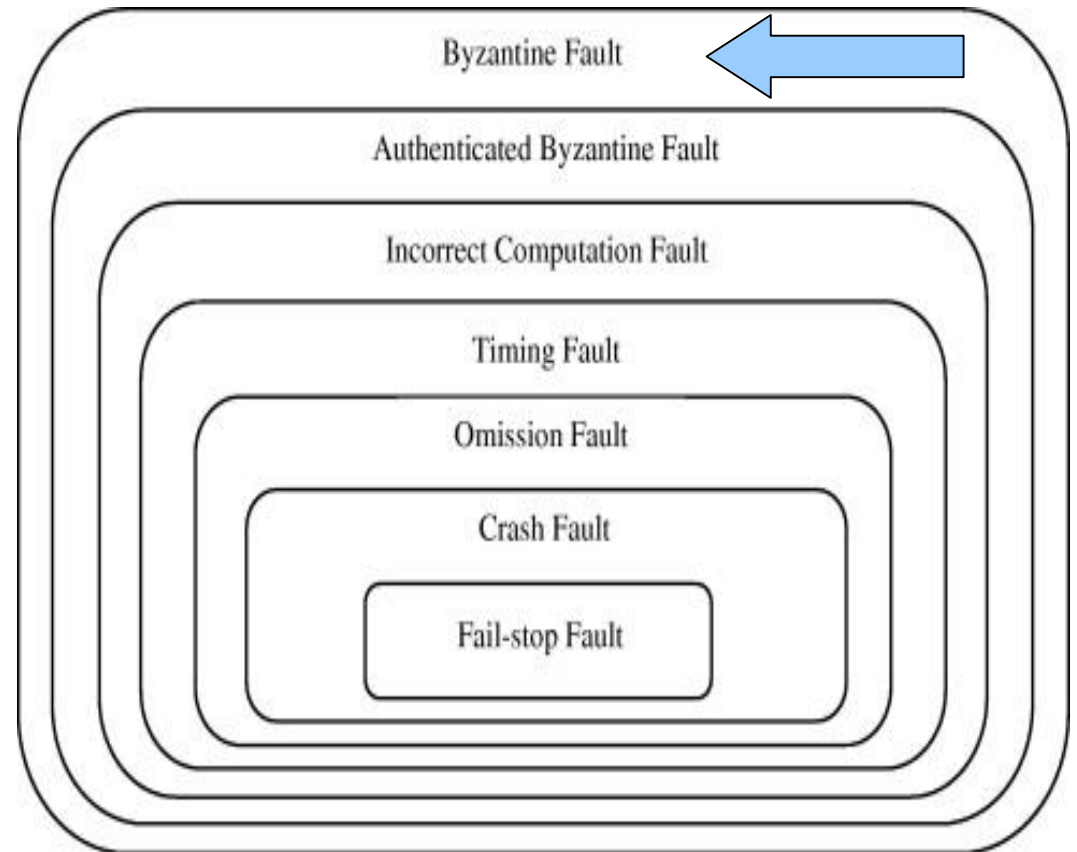
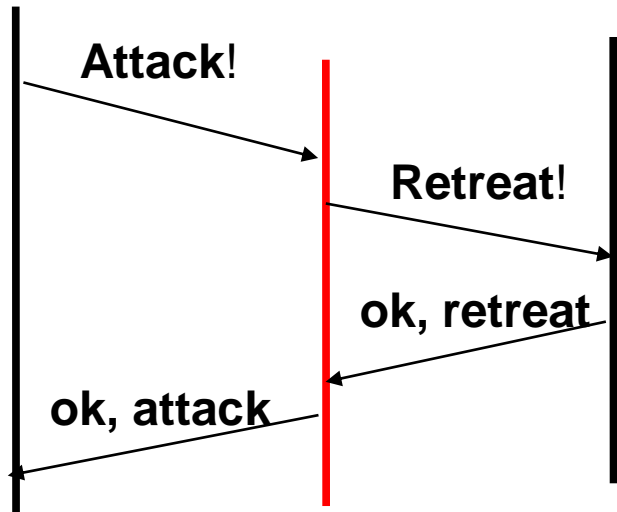
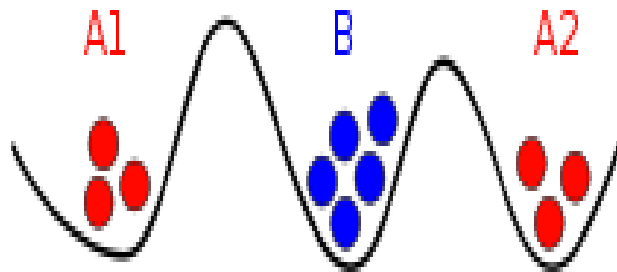
- Political decision
 - Recognize that “Daylight Saving Time” is a political decision
- Leap seconds = “35” forever
 - On the World RadioCommunications Conference 2015 (WRC-15), fix leap seconds to “35” forever
 - *(decision based on ITU studies happening now)*
- No Leap hour every 600 years
 - Idea: periodically, skip Daylight Savings Time for a year
 - <http://old.post-gazette.com/pg/05210/545823.stm>
 - <http://leapsecond.com/>

Solution 3: Rockets 😊

- Use rockets to control Earth's rotation
 - Either speed up, or speed down
 - We can even use a typical PTP “PI” servo
 - Could be an easier solution than the other two 😊

Byzantine Theory recap

Allied1 Enemy Allied2



Picture credits

- www.dilbert.com
- “Understanding and Applying Precision Time Protocol”, Steve T. Watt, Shankar Achanta, Hamza Abubakari, and Eric Sagen, Schweitzer Engineering Laboratories, Inc.
- “Mitigating GPS Vulnerabilities”, Shankar Achanta, Steve T. Watt, and Eric Sagen, Schweitzer Engineering Laboratories, Inc.
- ...

