



6th PhD School on Traffic Monitoring and Analysis TMA
2016 - Louvain La Neuve, Belgium, 5-6 April 2016

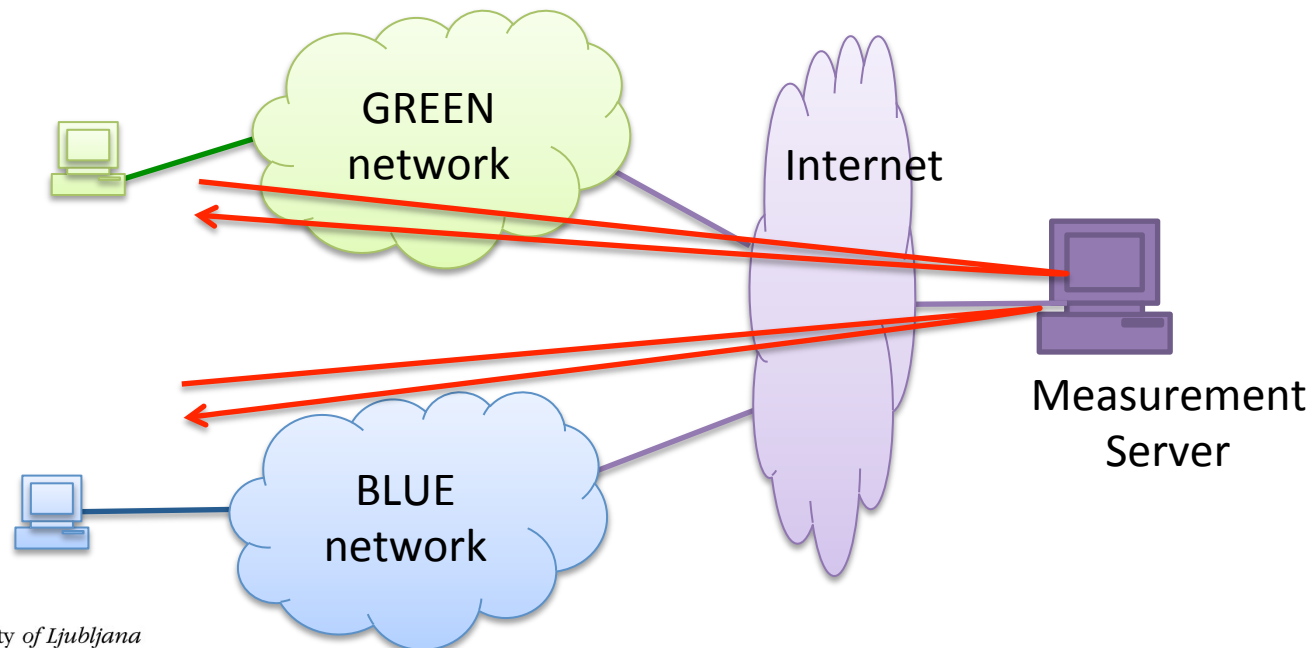
Toy Exercise - Assignment

Fabio Ricciato

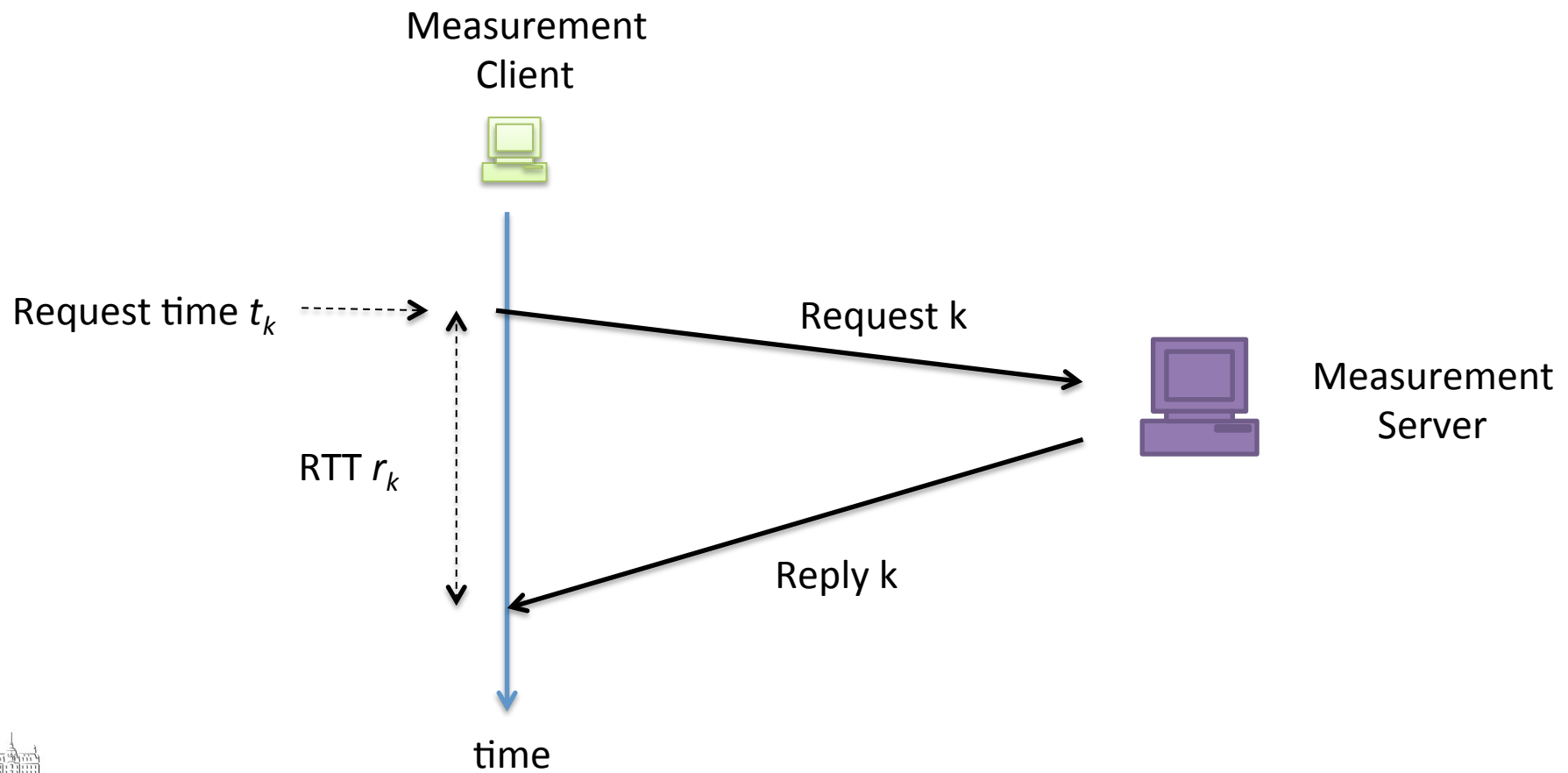


Scenario

- Scenario: the **Green** and **Blue** teams have measured Round-Trip Times (RTT) from their respective access networks via repeated `request/reply` to the same server



Methodology



Scenario

- They provide you
 - A succinct description of their measurement methodology
 - The raw measurement data
 - files dataB.txt, dataG.txt
 - two columns for Start time t_k and RTT r_k

StartTime	RTT
972.00000000	0.01172138
972.02004475	0.01476894
972.04011993	0.01303928
972.06021071	0.01749107
...	...



Measurement methodology described by Blue team

“We run a series of periodic request/reply measurements towards the common measurement server. For every request packet k we record the departure timestamp t_k and start a timer. Every request carries a unique ID in the payload that is replicated in the reply packet, in order to ensure correct association between the reply and the corresponding request. When the reply packet is received, we record the value of the timer r_k . If the reply is not received within a maximum predefined timeout, we mark the request as “lost” and write “-1” in the output file. Consecutive measurements are spaced by 20 ms. The experiment started at 9:16 AM of 11.3.2016.”



Measurement methodology described by Green Team

“We have followed the same measurement methodology described by the Blue Team, with the same spacing interval of 20 msec between consecutive measurements. But we started a bit later, around 9:21 AM of the same day.”



Your task

- Process the data and ...
- ... determine which network has the **best performances**
- motivate your answer, provide quantitative performance metric values in support of your answer.



- Have fun 😊
- for questions & comments email to:

fabio.ricciato@fri.uni-lj.si





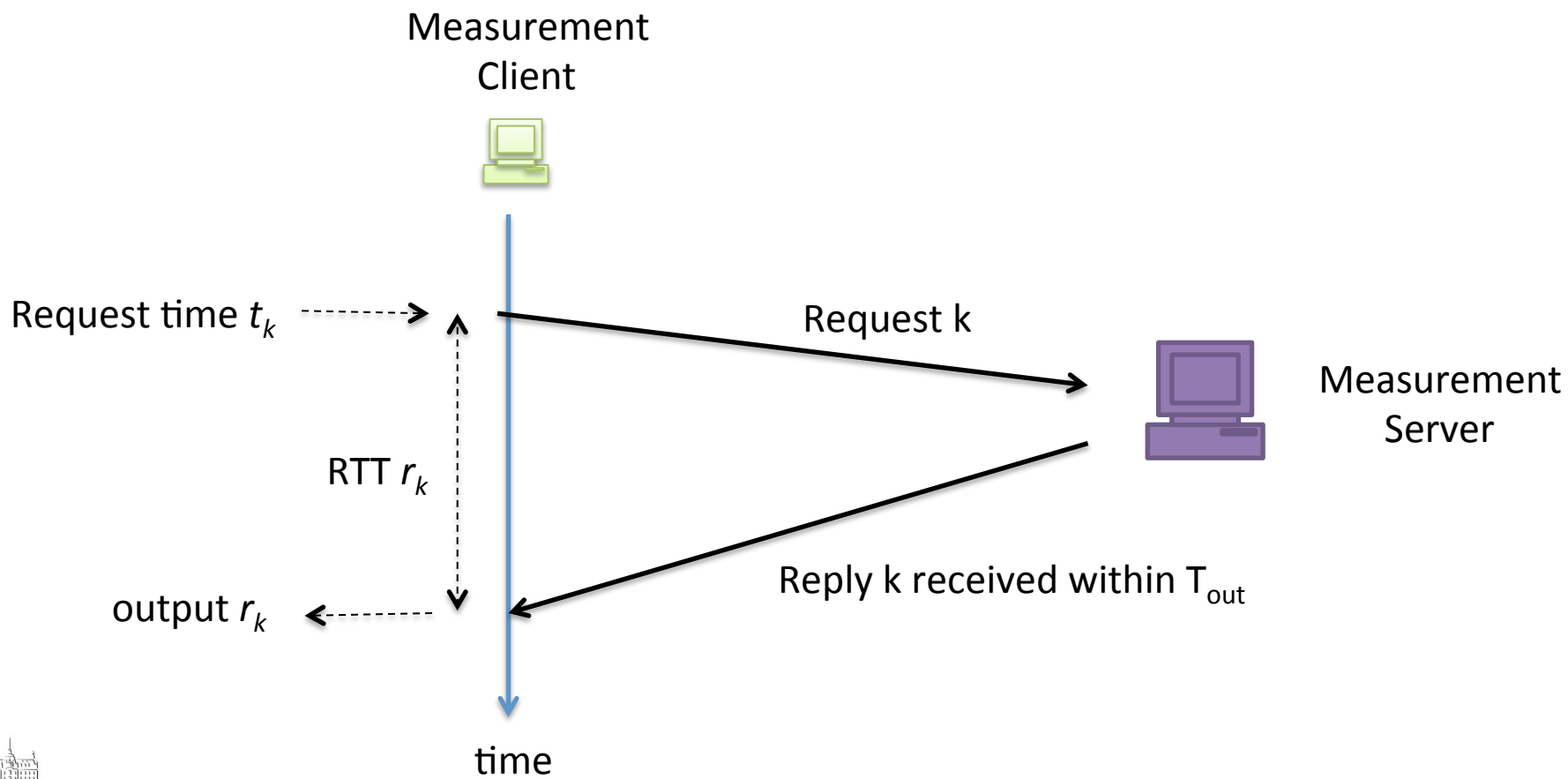
6th PhD School on Traffic Monitoring and Analysis TMA
2016 - Louvain La Neuve, Belgium, 5-6 April 2016

Toy Exercise

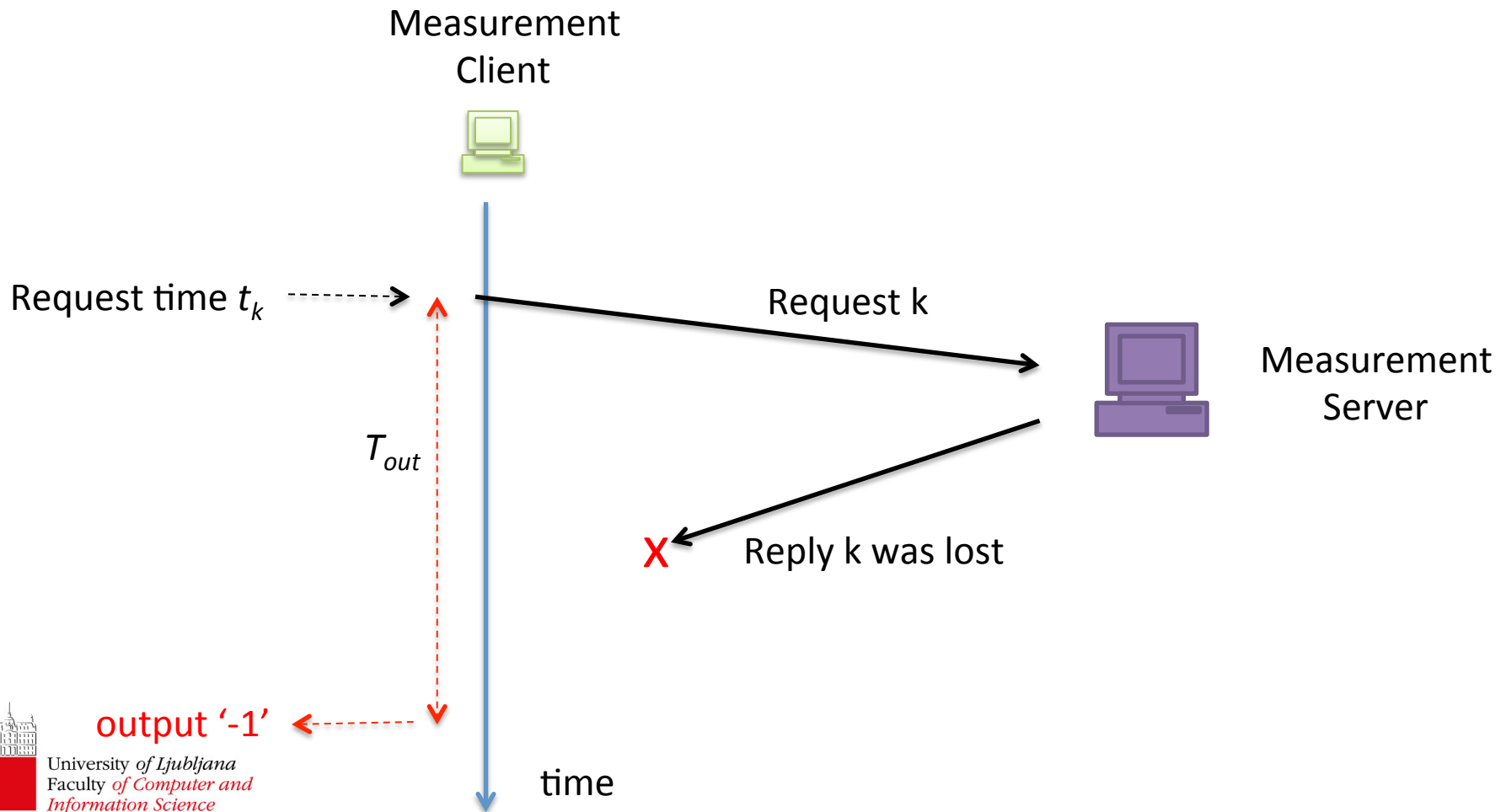
Sketch of Solution



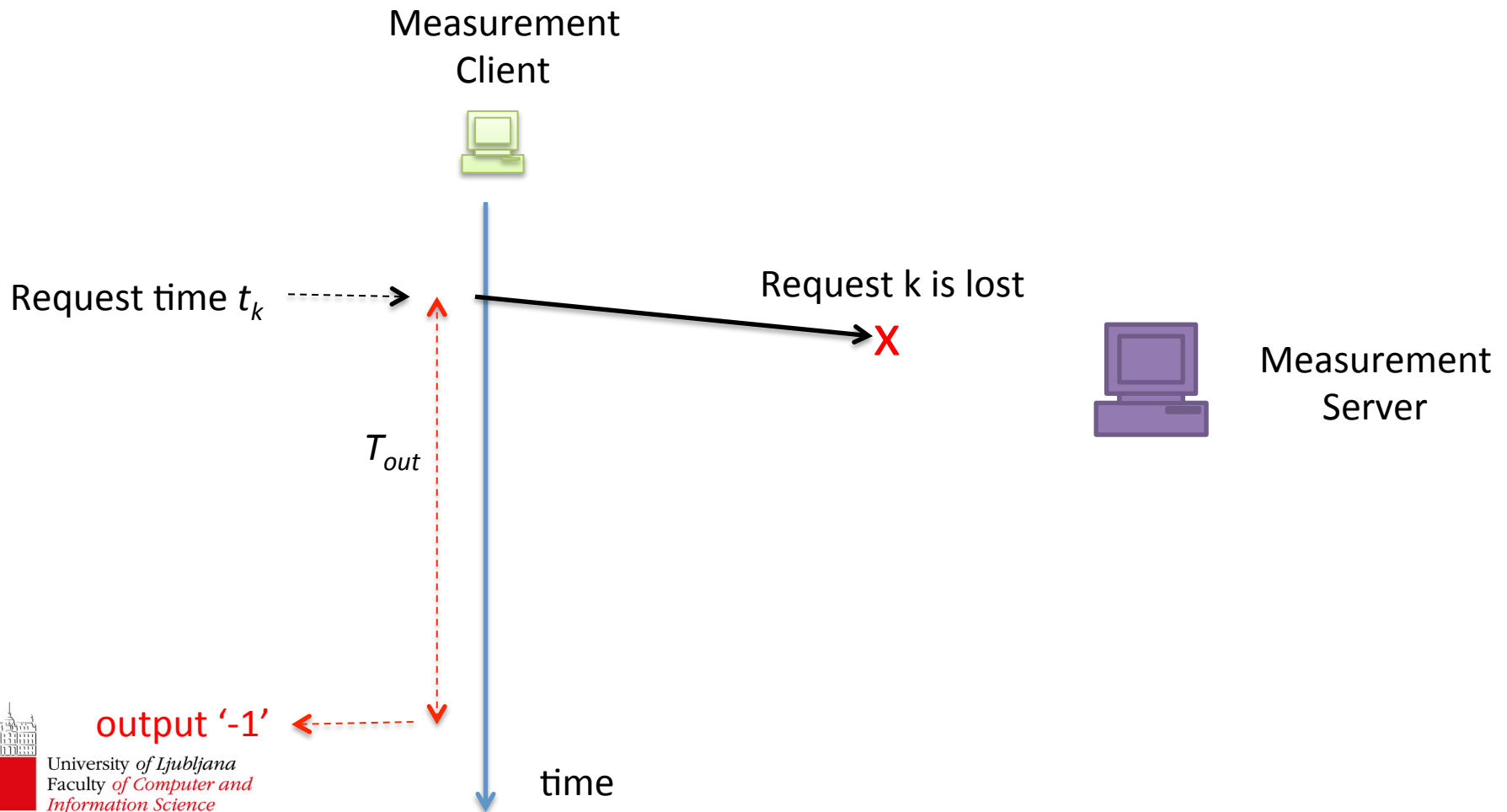
RTT sample



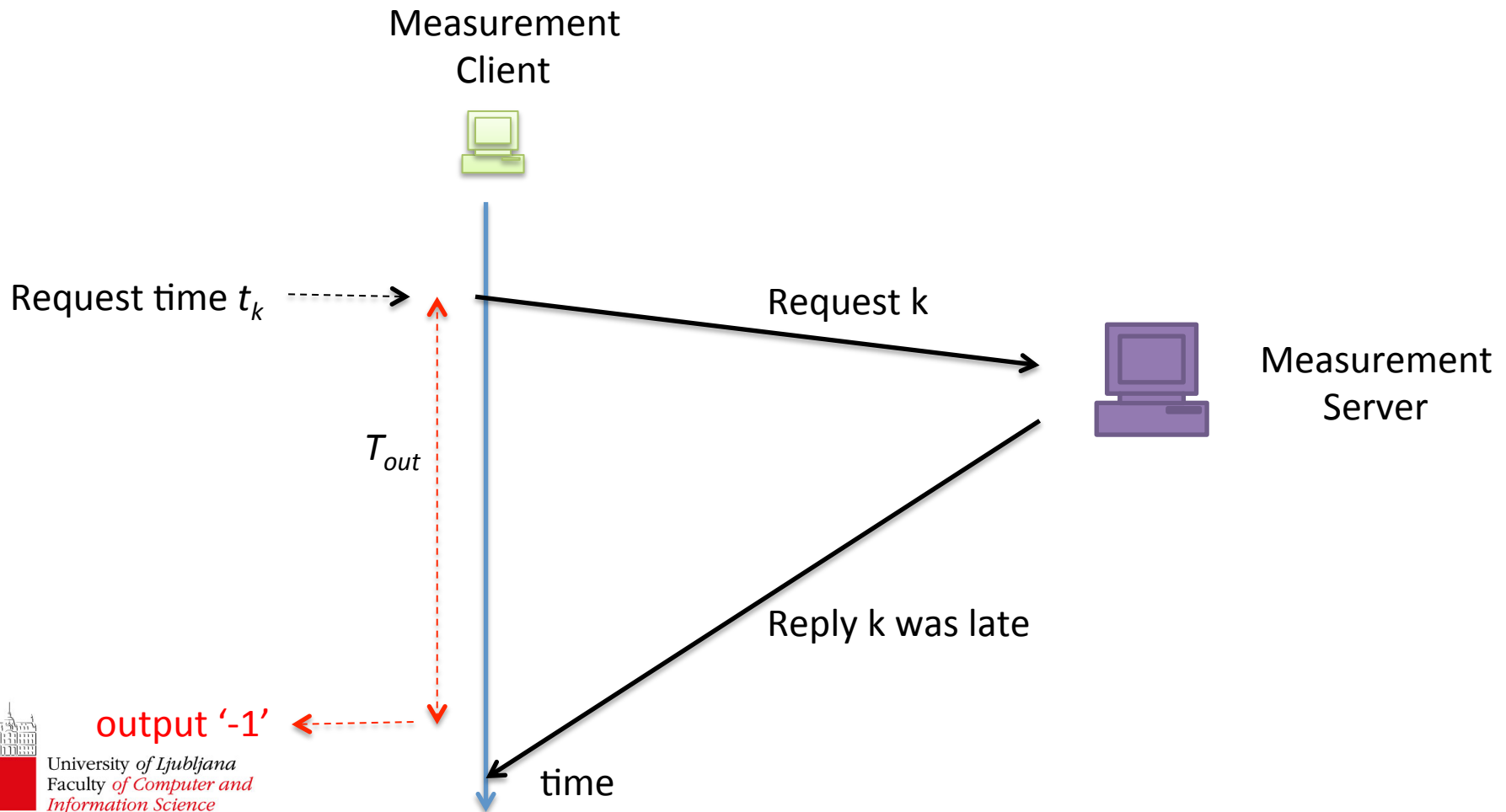
Timeout expired



Timeout expired



Timeout expired



A First look at the two datasets

Blue dataset

- $N=30,000$ total records
- $n_0= 949$ records w/o valid RTT ('-1')

Green dataset

- $N=30,000$ total records
- $n_0= 3466$ records w/o valid RTT ('-1')



A naive answer

Blue dataset

- N=30,000 total records
- $n_0 = 949$ records w/o valid RTT ('-1')

Green dataset

- N=30,000 total records
- $n_0 = 3466$ records w/o valid RTT ('-1')

$$u_0 = \frac{\# \text{ unanswered requests}}{\# \text{ all requests}} = \frac{n_0}{N}$$

- $u_0 = 949/30,000 = 0.0316$

- $u_0 = 3466/30,000 = 0.1155$

Blue network has lower $u_0 \rightarrow$ Blue network is better



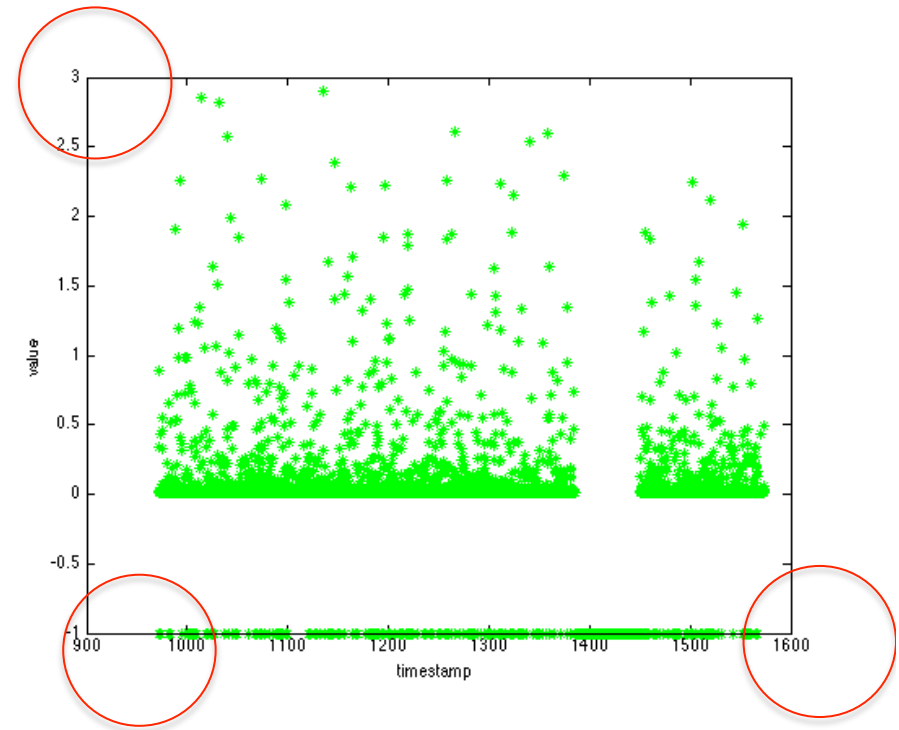
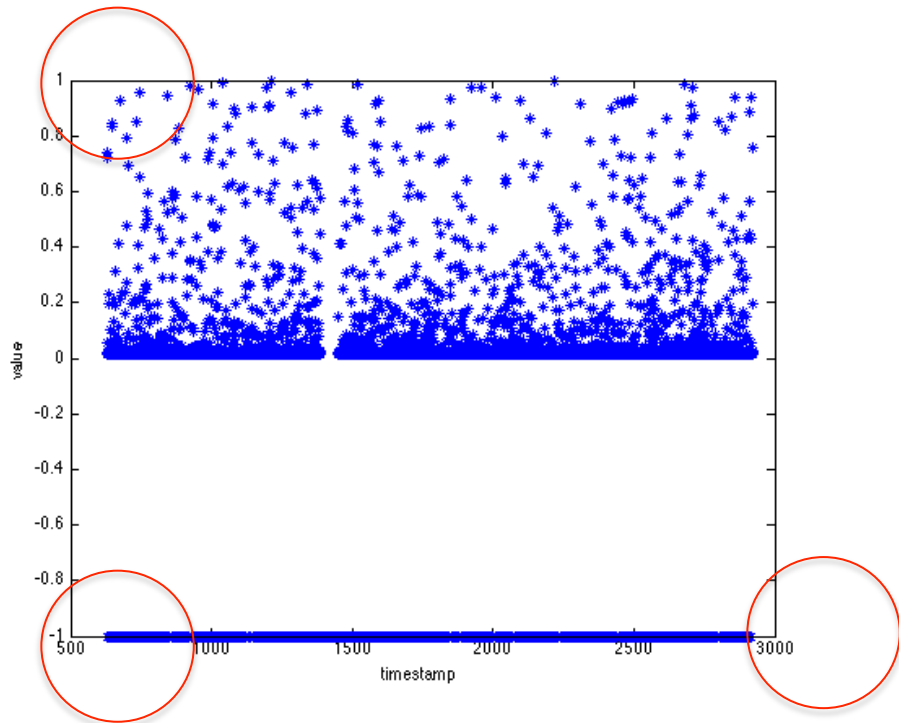
Preliminary look

628.0000000	0.0130000
628.0330000	0.0160000
628.0690000	0.0180000
628.1080000	0.0140000
628.1430000	0.0140000
628.1770000	0.0150000
628.2120000	0.0150000
628.2480000	-1.0000000
629.2680000	0.0140000
629.3020000	0.0180000
629.3410000	0.0150000
629.3770000	0.0160000
629.4130000	0.0150000
629.4490000	0.0180000
629.4870000	0.0170000
629.5240000	0.0160000
629.5610000	0.0180000
629.5990000	0.7420000
630.3610000	0.0570000
630.4380000	-1.0000000
631.4580000	0.7250000
632.2040000	0.0150000
632.2400000	0.0180000
632.2780000	0.0140000

972.02004475	0.01476894
972.04011993	0.01303928
972.06021071	0.01749107
972.08028098	0.01502918
972.10033678	0.01586265
972.12040874	0.01742269
972.14046197	0.01306065
972.16047813	0.01737577
972.18056606	0.01701450
972.20063363	0.01365929
972.22068388	0.01770575
972.24068468	0.01799749
972.26073387	0.33428613
972.28078814	0.01463612
972.30086112	0.01851273
972.32094276	0.01797620
972.34100671	0.01418419
972.36109110	0.01613566
972.38115101	0.01559531
972.40123778	0.01541753
972.42124092	0.01317477
972.44133156	0.01357658
972.46139773	0.01760751
972.48145926	0.01772676



Preliminary look



timestamps are local, not absolute



A First look at the two datasets

Blue dataset

- $N=30,000$ total records
- $n_0= 949$ records w/o valid RTT ('-1')
- $\max(r_k) = 0.998$
 - T_{out} likely set to 1 sec

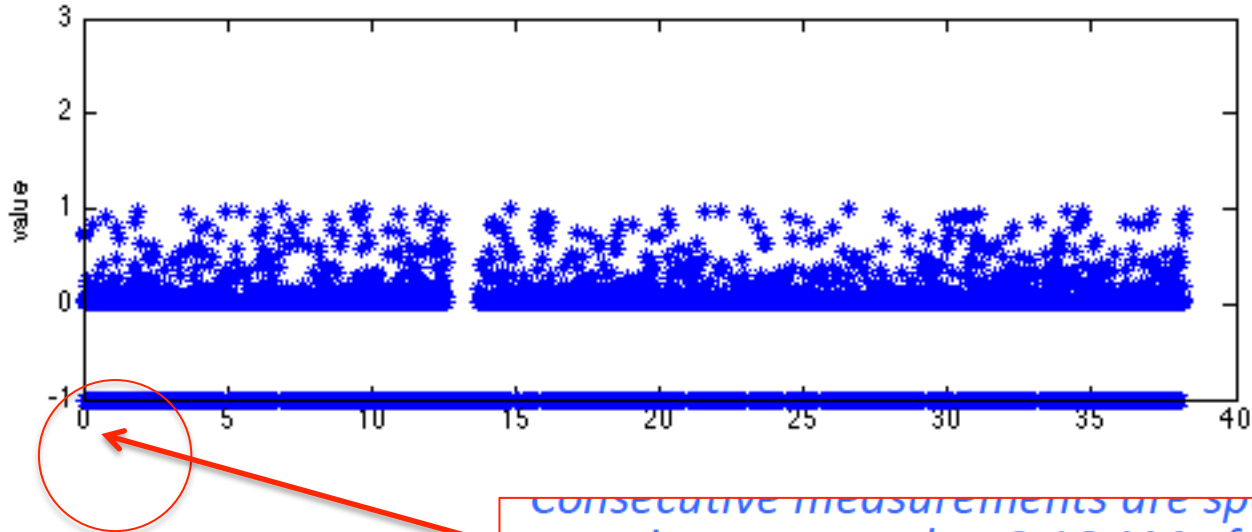
Green dataset

- $N=30,000$ total records
- $n_0= 3466$ records w/o valid RTT ('-1')
- $\max(r_k)=2.896$
 - T_{out} likely set to 3 sec



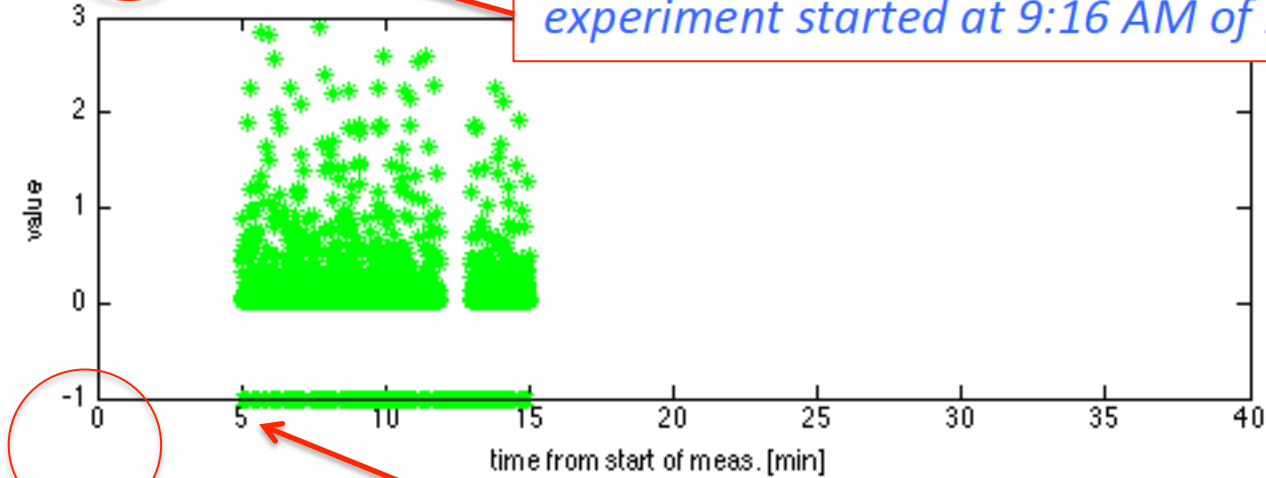
Adjusting timestamps (approximately)

removing
initial offset



*Consecutive measurements are spaced by 20 ms
experiment started at 9:16 AM of 11.3.2016."*

rescaling to
minutes



But we started a bit later, around 9:21 AM

A First look at the two dataset

Blue dataset

- N=30,000 total records
- $n_0 = 949$ records w/o valid RTT ('-1')
- $\max(r_k) = 0.998$
 - T_{out} likely set to 1 sec
- **total duration 10 min**

Green dataset

- N=30,000 total records
- $n_0 = 3466$ records w/o valid RTT ('-1')
- $\max(r_k) = 2.896$
 - T_{out} likely set to 3 sec
- **total duration 38 min**



A First look at the two dataset

Blue dataset

- **N=30,000** total records
- $n_0 = 949$ records w/o valid RTT ('-1')
- $\max(r_k) = 0.998$
 - T_{out} likely set to 1 sec
- **total duration 10 min**

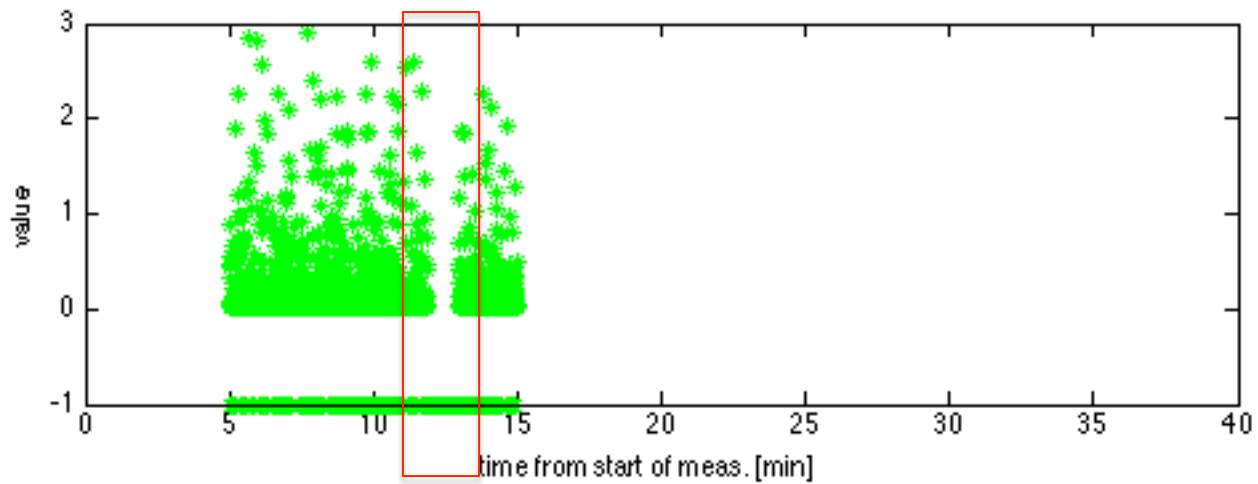
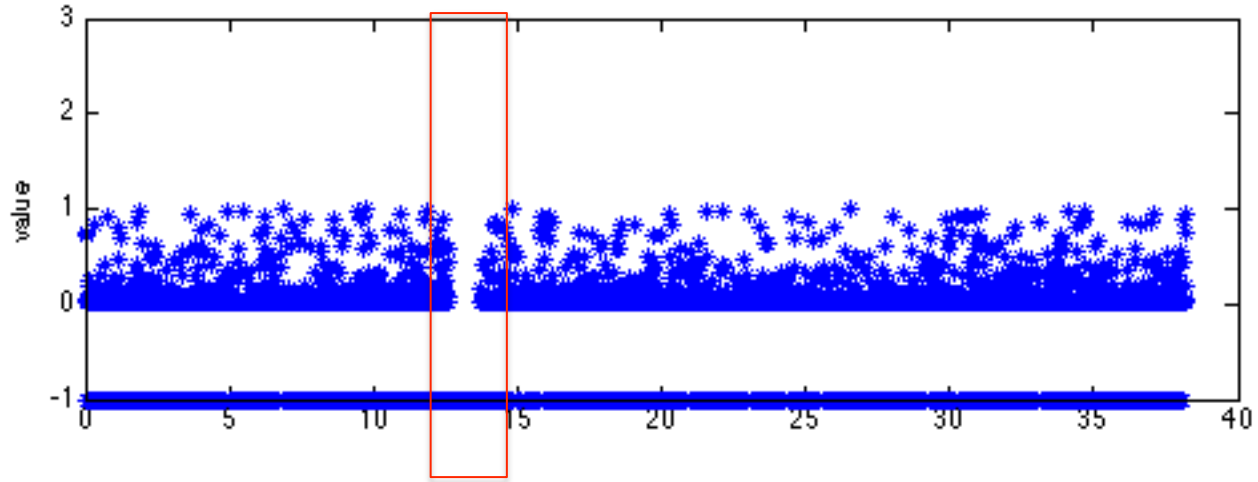
Green dataset

- **N=30,000** total records
- $n_0 = 3466$ records w/o valid RTT ('-1')
- $\max(r_k) = 2.896$
 - T_{out} likely set to 3 sec
- **total duration 38 min**

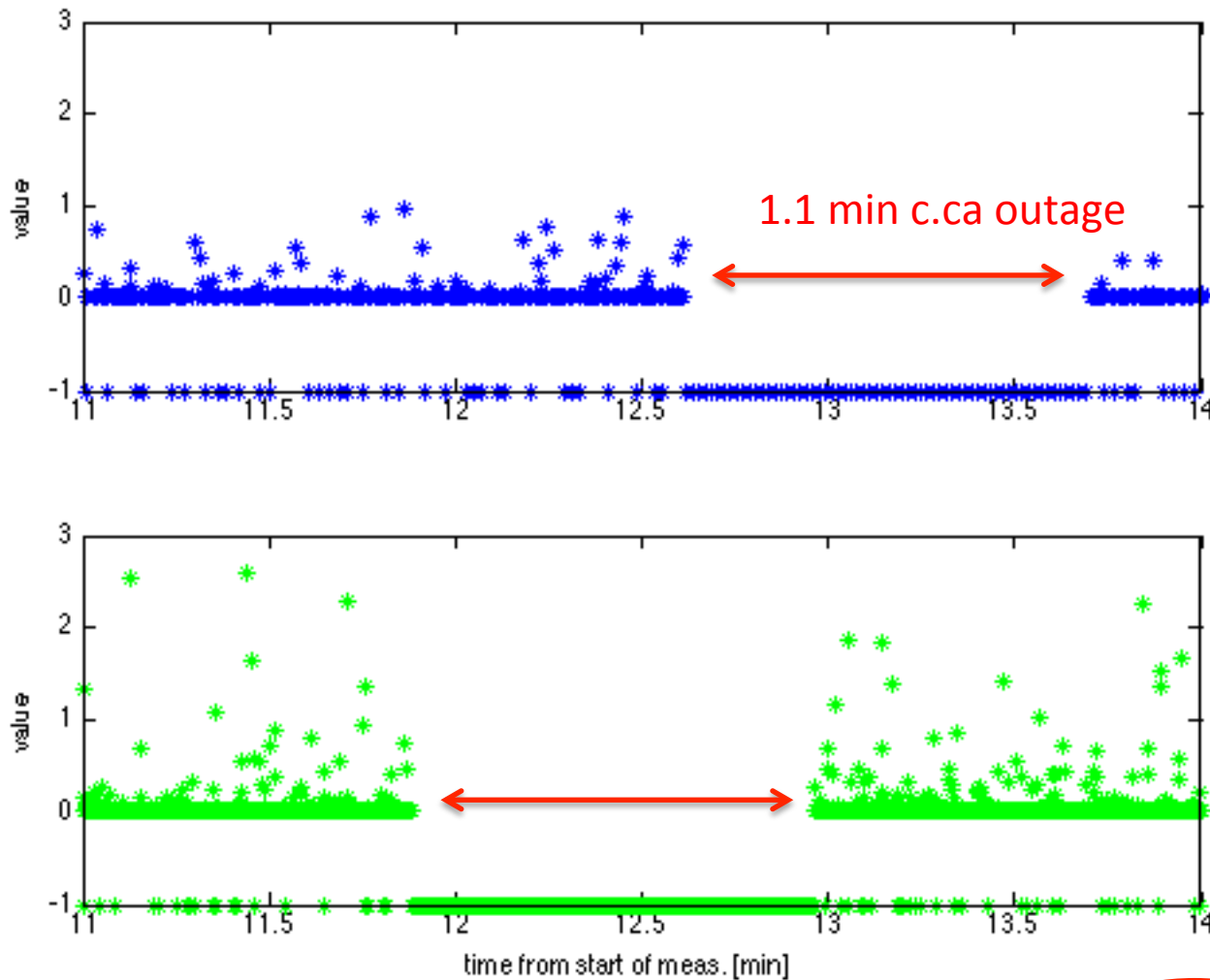
Same number of records (30,000), same spacing (20 ms),
but different duration ??????



outage ?



zooming into outage



probably the same outage ?

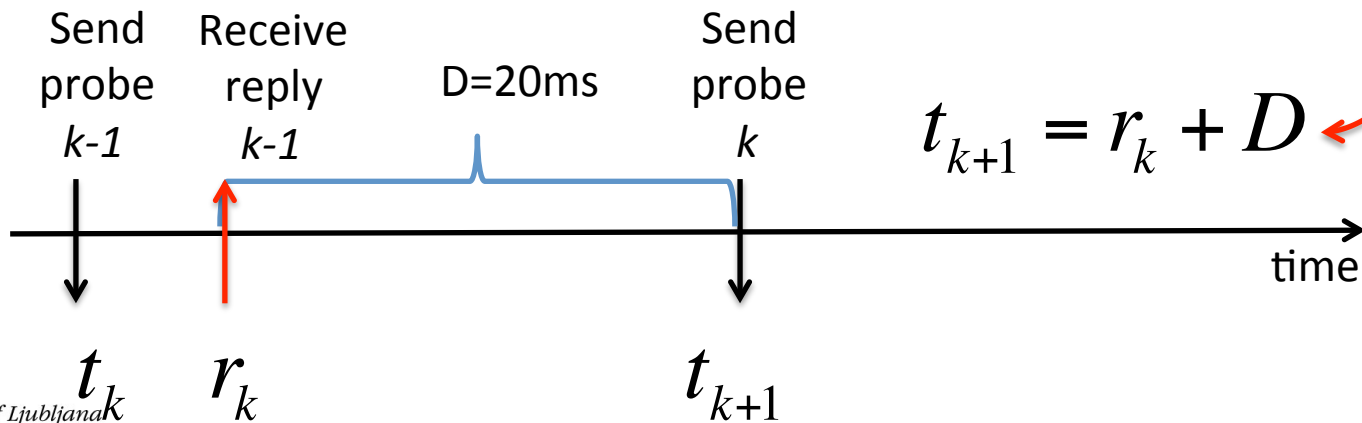
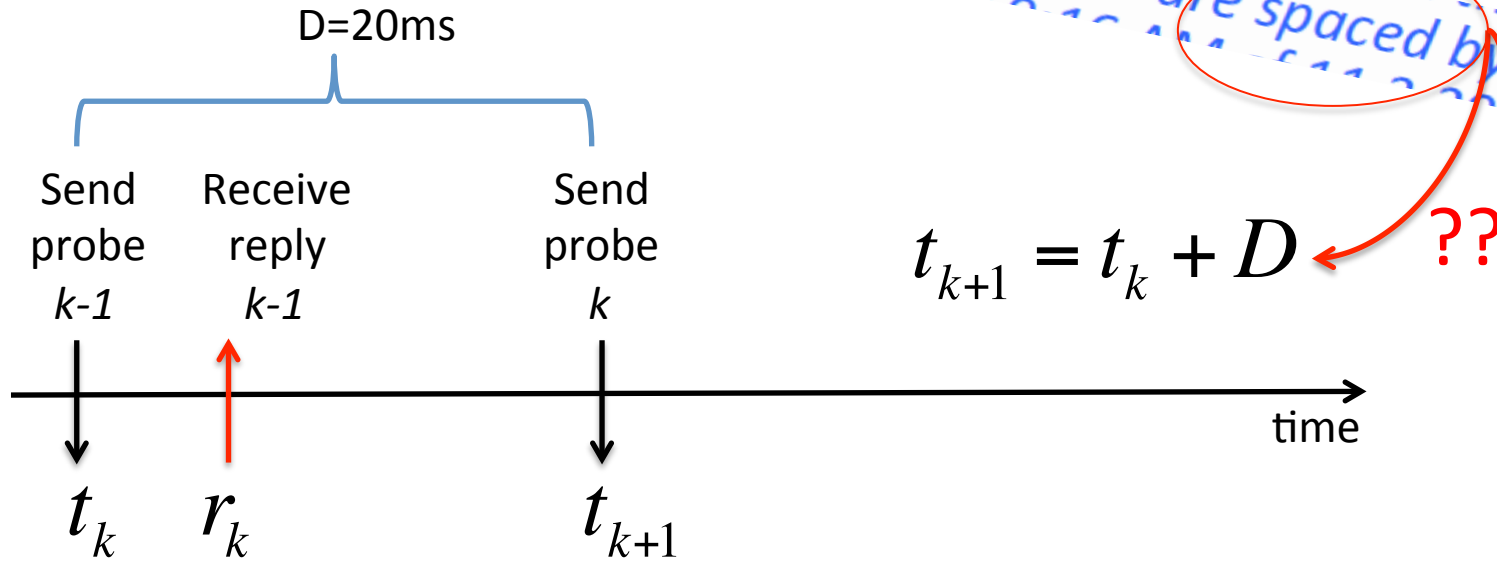
probably not due to causes external to the network-under-test, hence not relevant for the comparison

But we started a bit later, around 9:21 AM

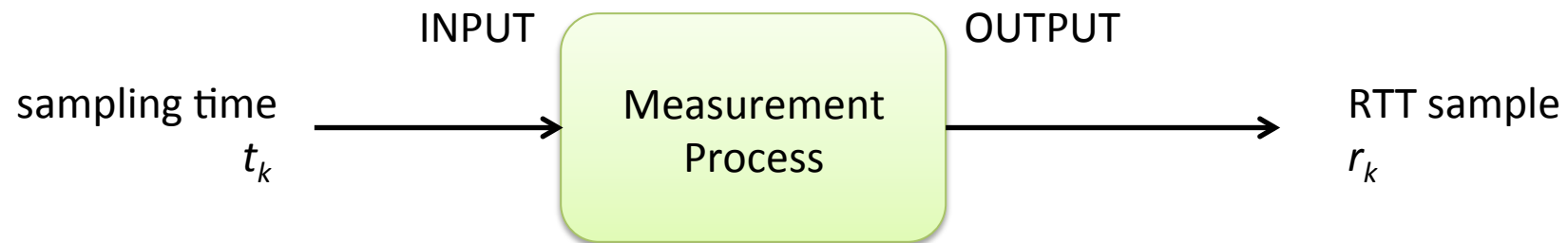


Ambiguity

Consecutive measurements are spaced by 20 ms.



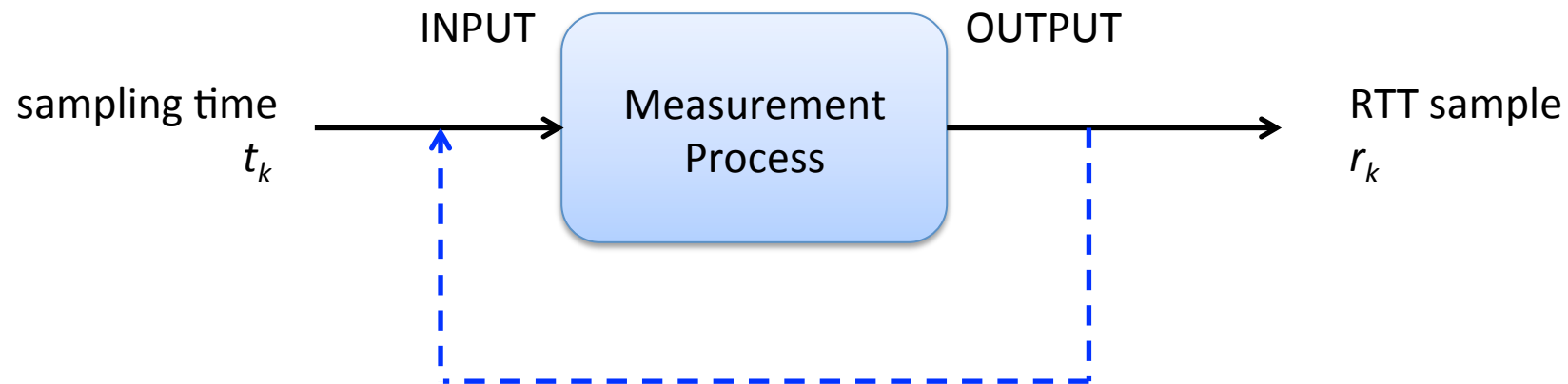
Measurement process Input/Output



$$t_{k+1} = t_k + D$$



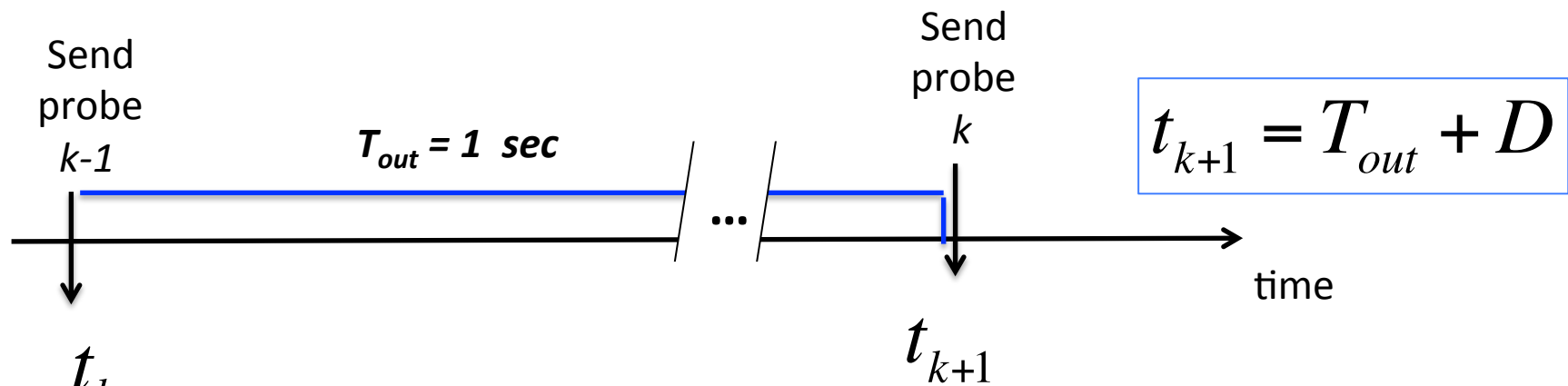
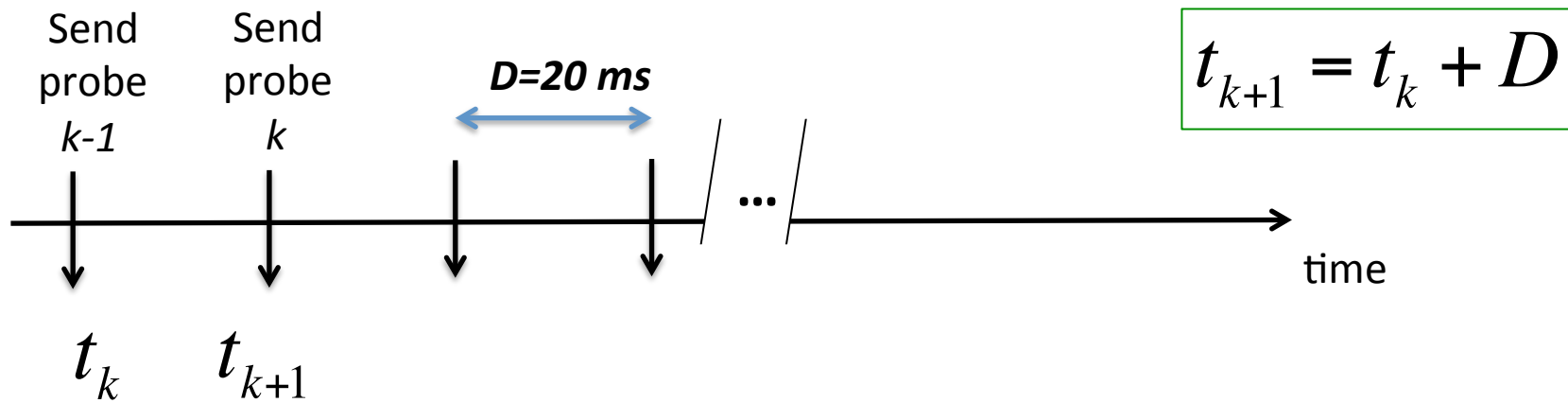
Measurement process Input/Output



$$t_{k+1} = t_k + r_k + D$$

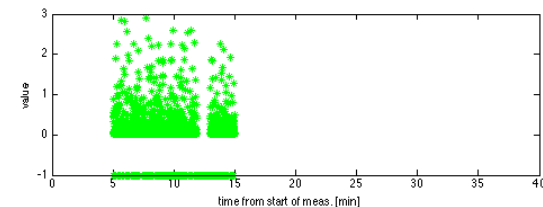
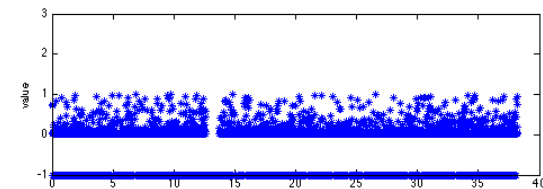
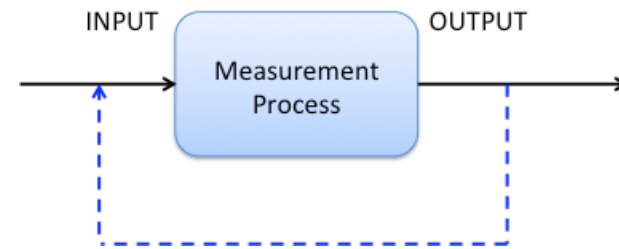


during the outage ...



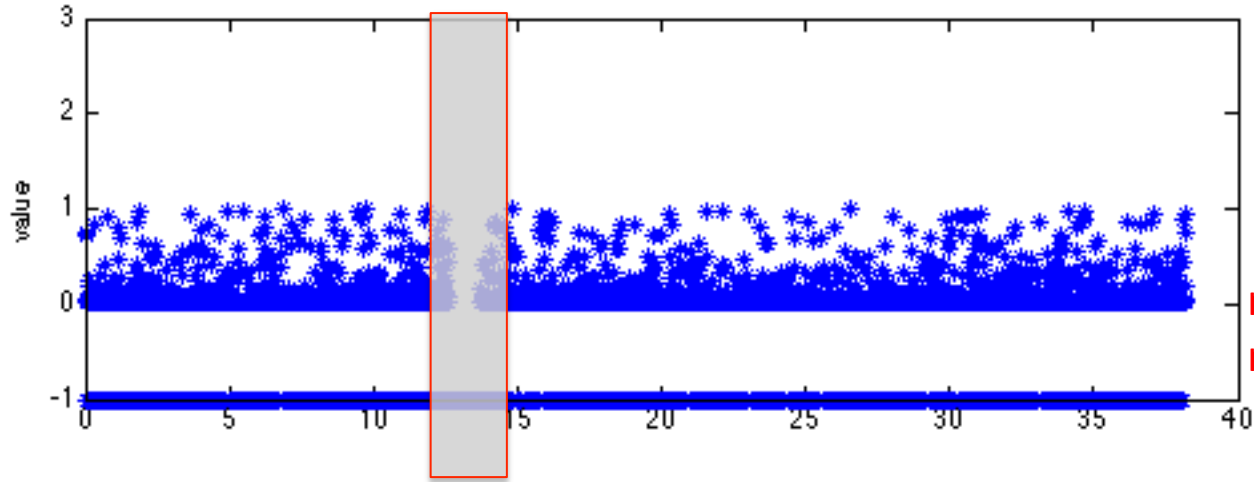


- (uncontrolled) correlations between output and input of measurement system
- Risk of distortion (bias)
 - under- or over-representation of certain phenomena
- other side effects

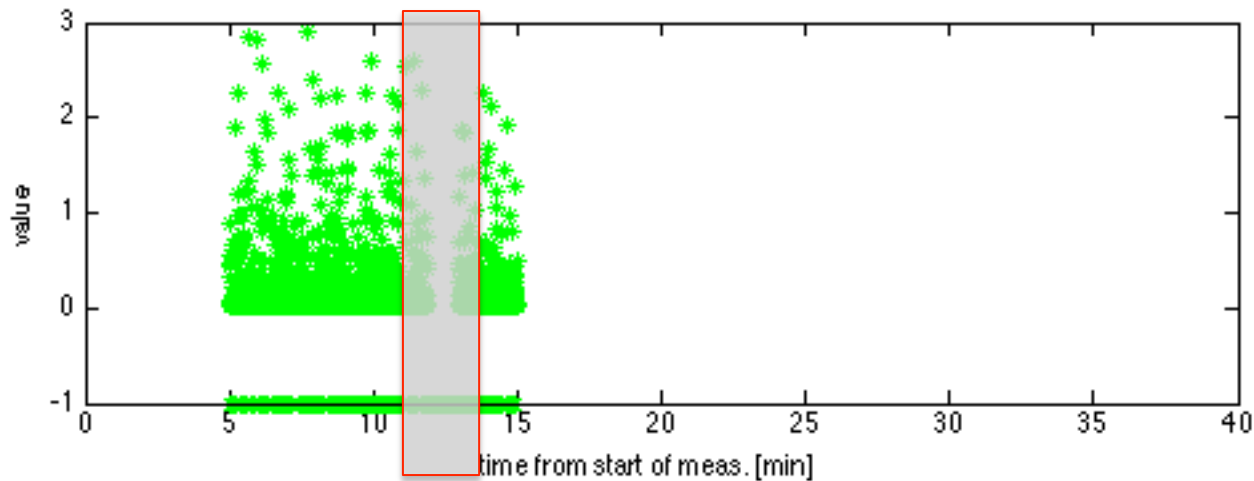


discarding data during outage

Discarding
data
during
outage



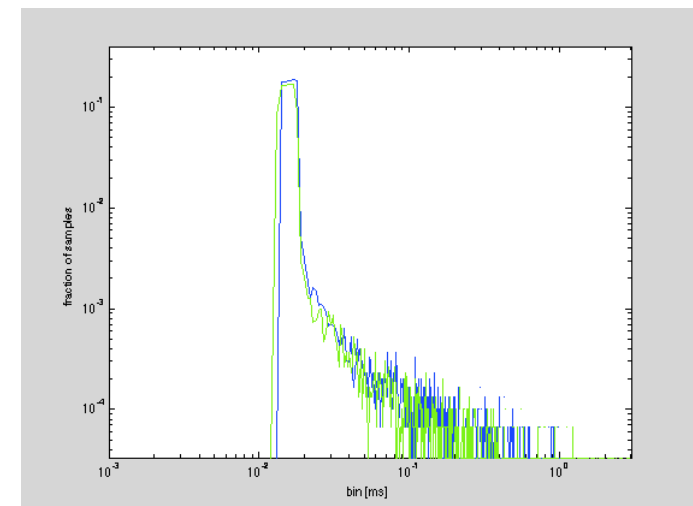
rescaling to
minutes



Let's look at the RTT distribution

- Histograms?
 - How to bin ? Linearly, logarithmically ...
 - What bin size ?
 - In any case you loose resolution (bin aggregation)
 - ...

- Avoid all that by looking at Cumulative Distributions!



Definition of ECDF, ECCDF

- CDF: Cumulative Distribution Function
- CCDF: Complementary CDF
- ECDF: Empirical CDF

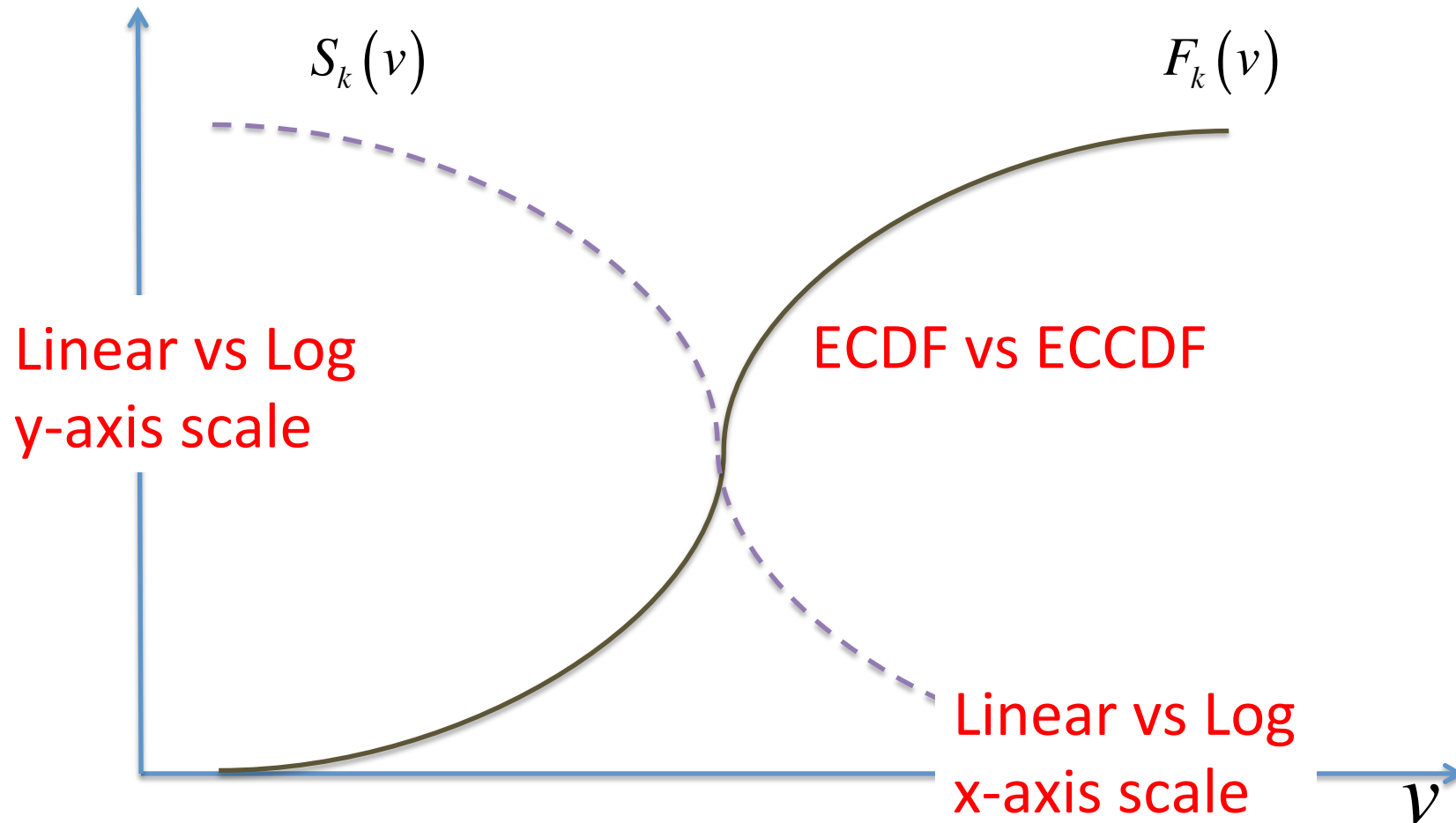
$$F_k(v) = \frac{\text{number of elements } \leq v}{\text{total number of elements}} = \frac{1}{k} \sum_{i=1}^k \mathbf{1}\{x_i \leq v\}$$

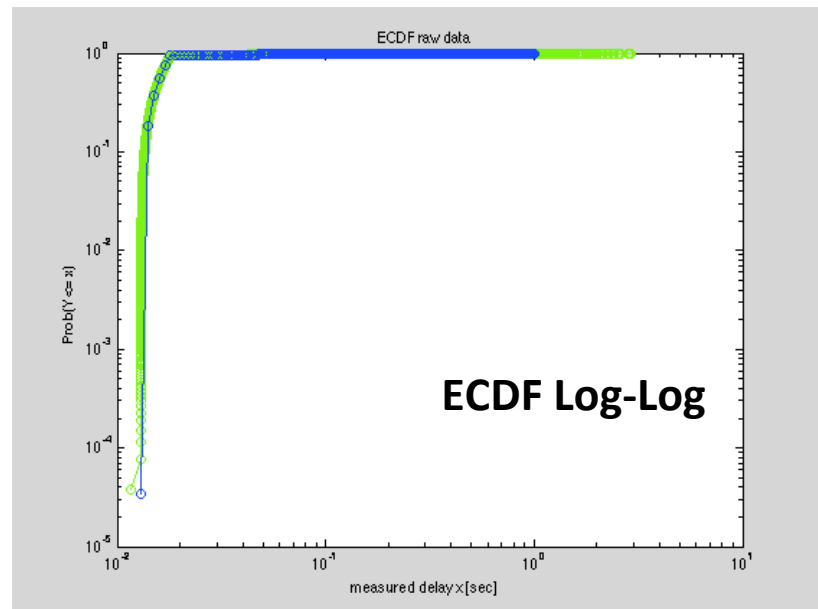
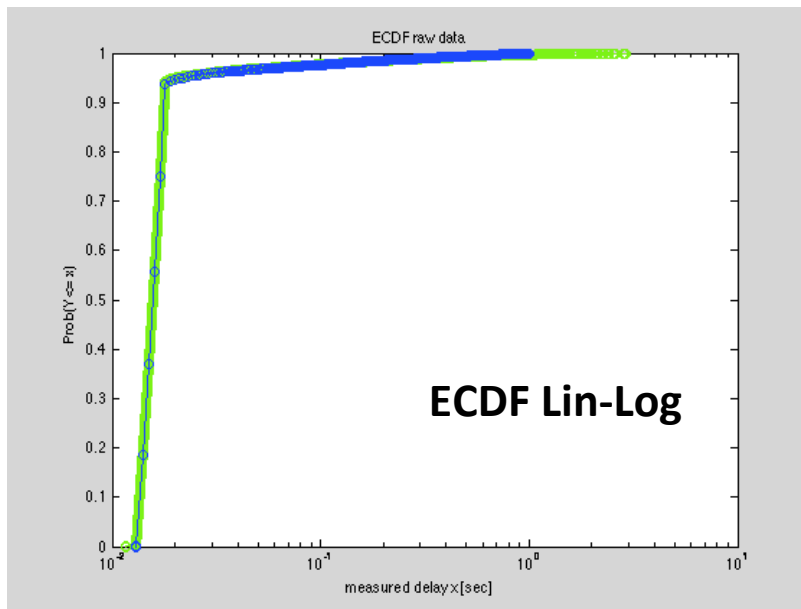
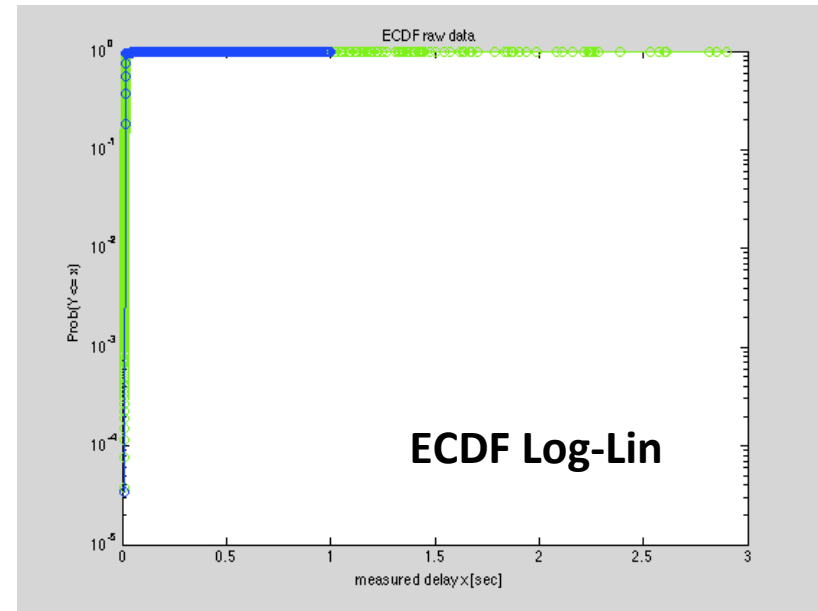
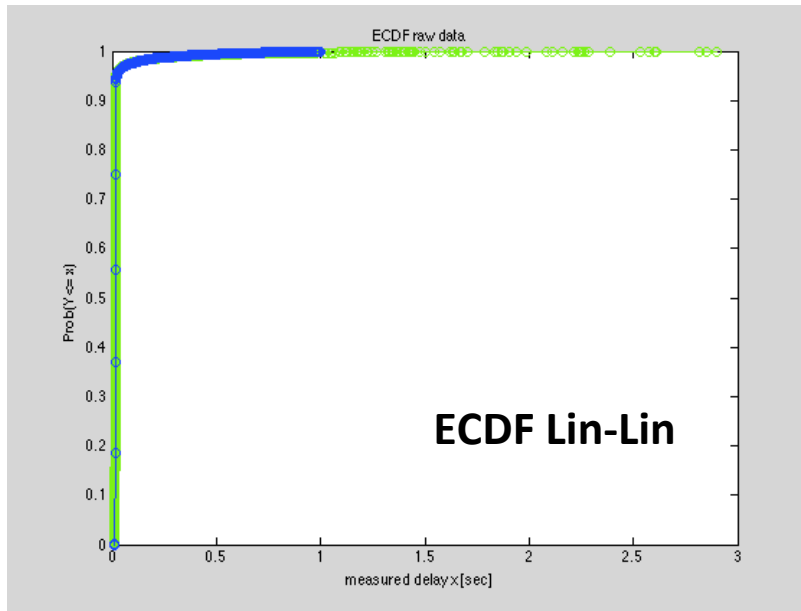
- ECCDF: Empirical Complementary CDF

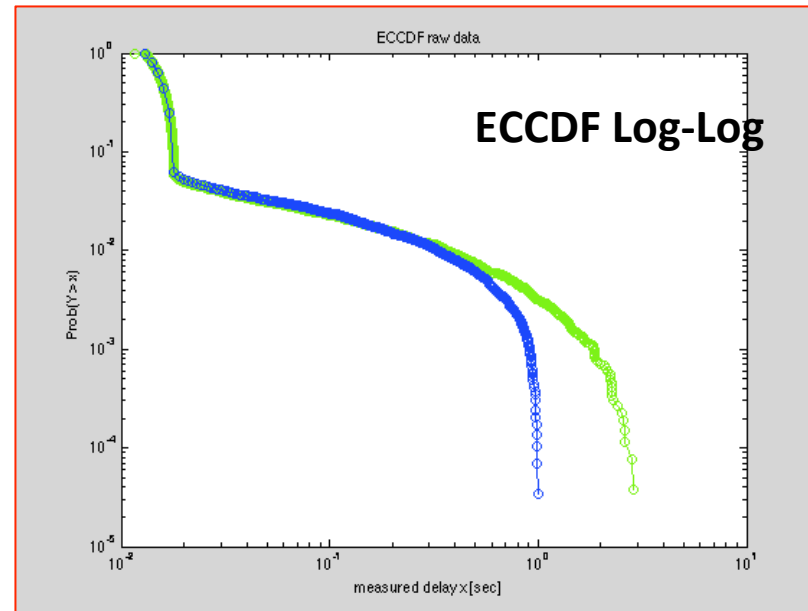
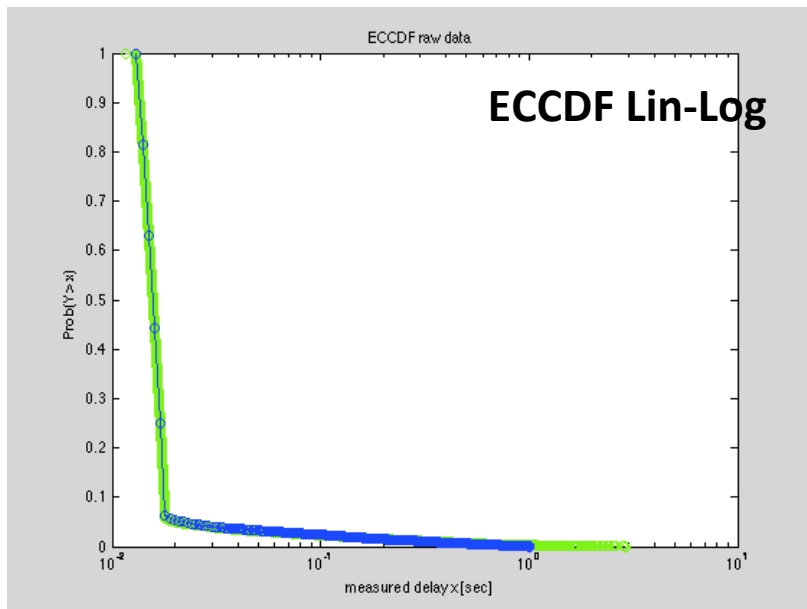
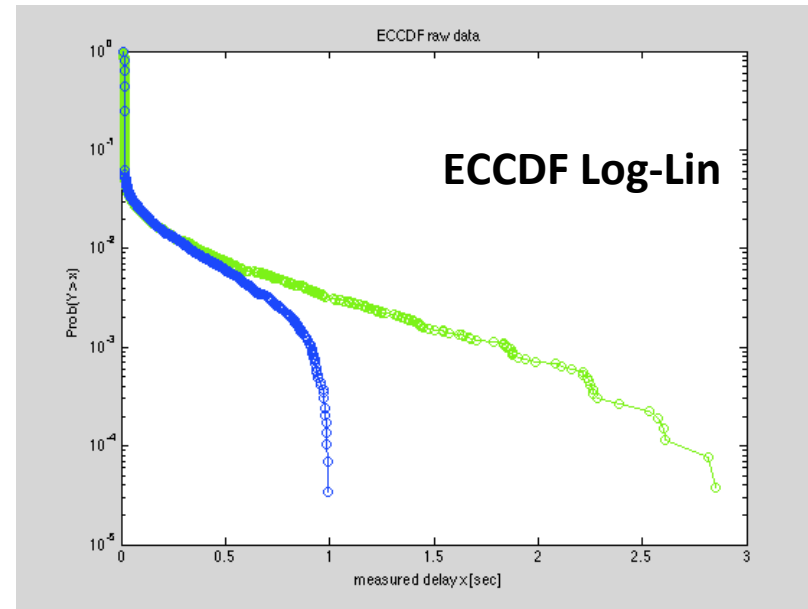
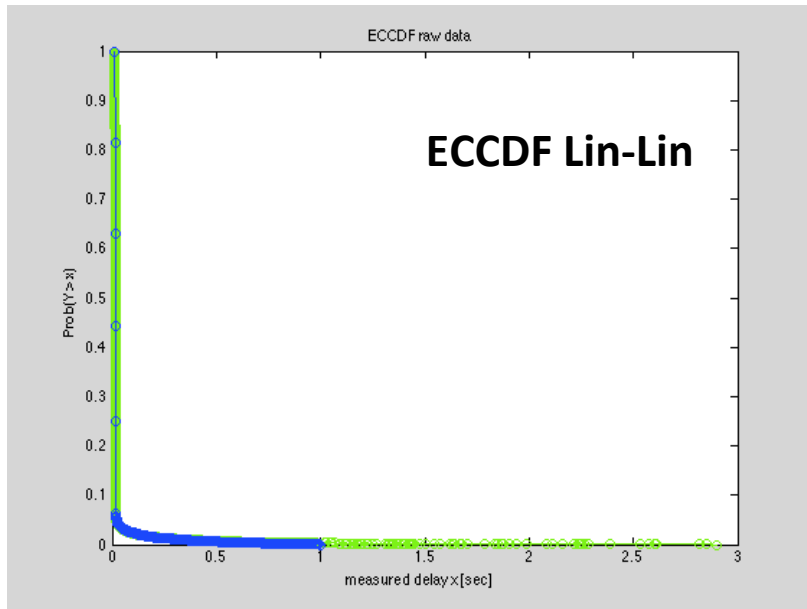
$$S_k(v) = \frac{\text{number of elements } > v}{\text{total number of elements}} = \frac{1}{k} \sum_{i=1}^k \mathbf{1}\{x_i > v\}$$

$$S_k(v) = 1 - F_k(v)$$

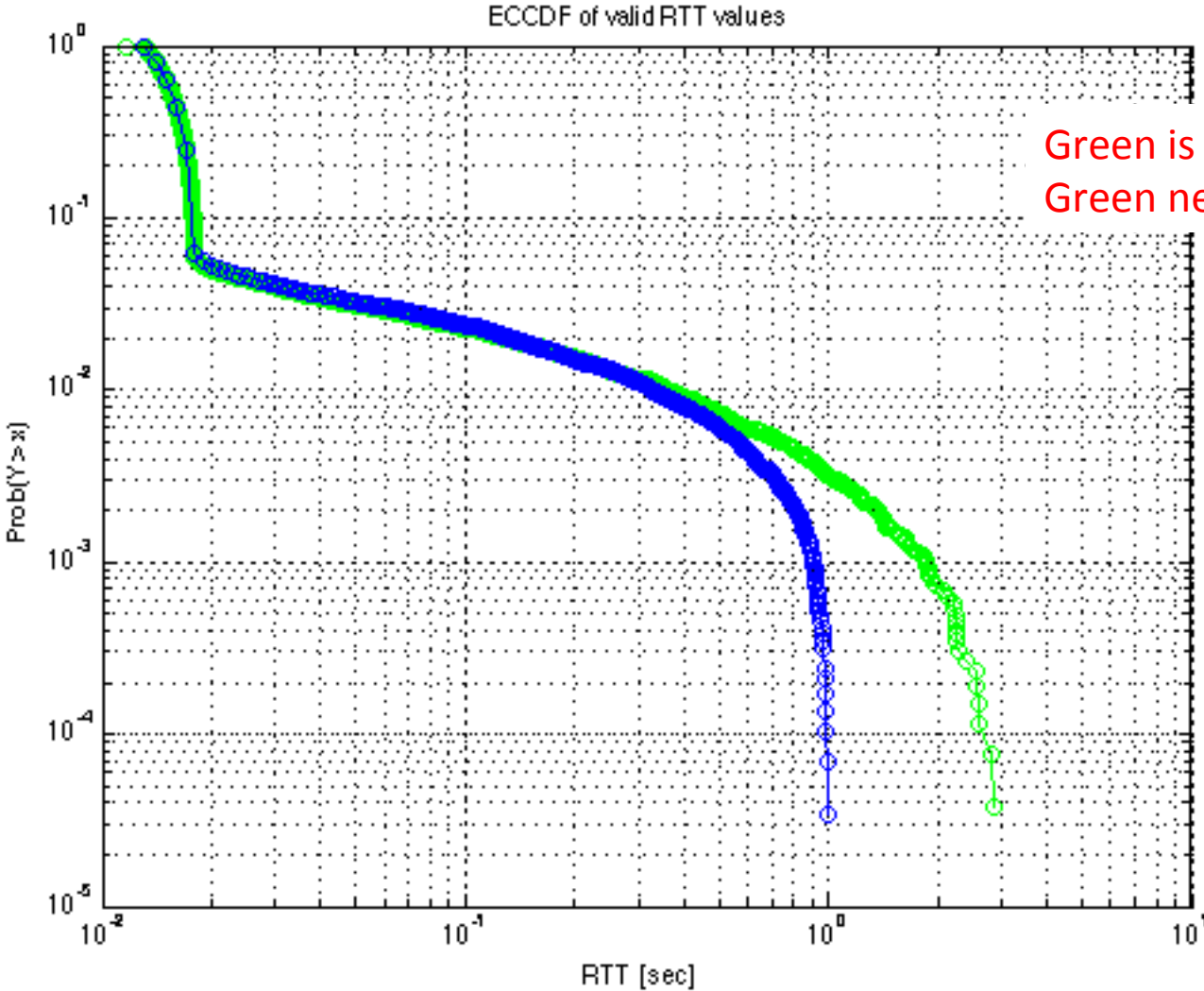
8 different ways of plotting a cumulative distribution





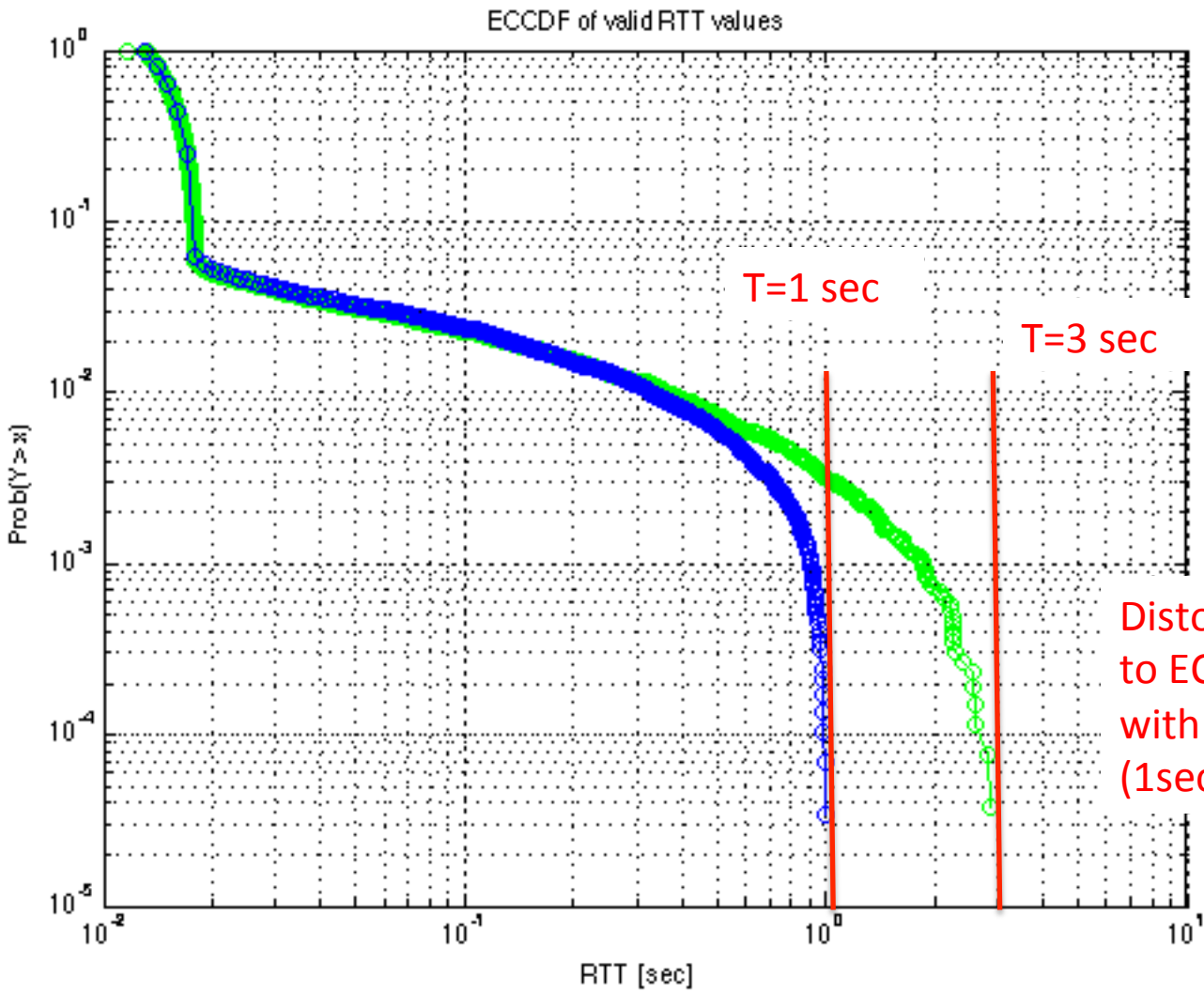


ECCDF of valid RTT samples (loglog)



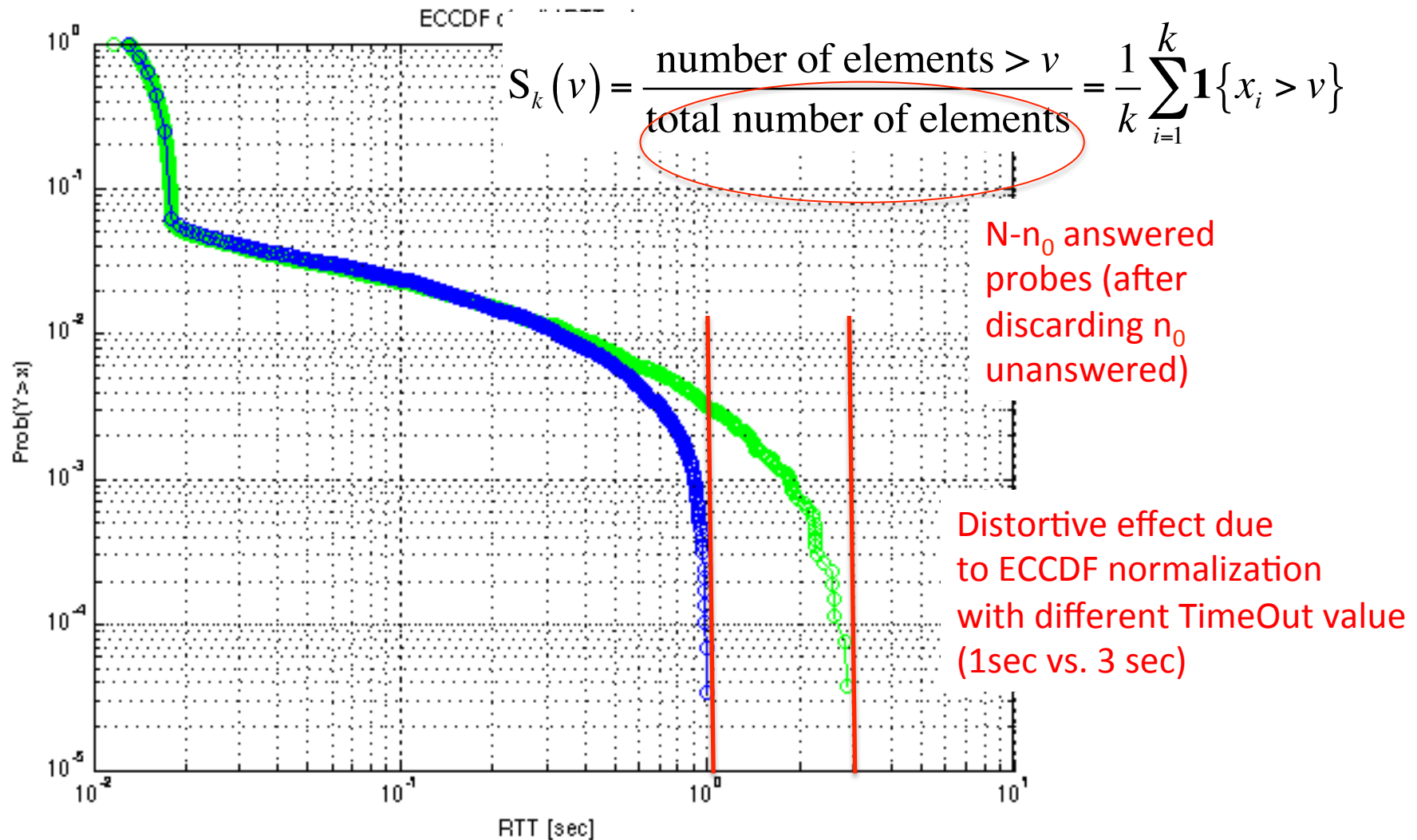
Green is above/right of Blue
Green network has larger RTT ???

ECCDF of valid RTT samples (loglog)



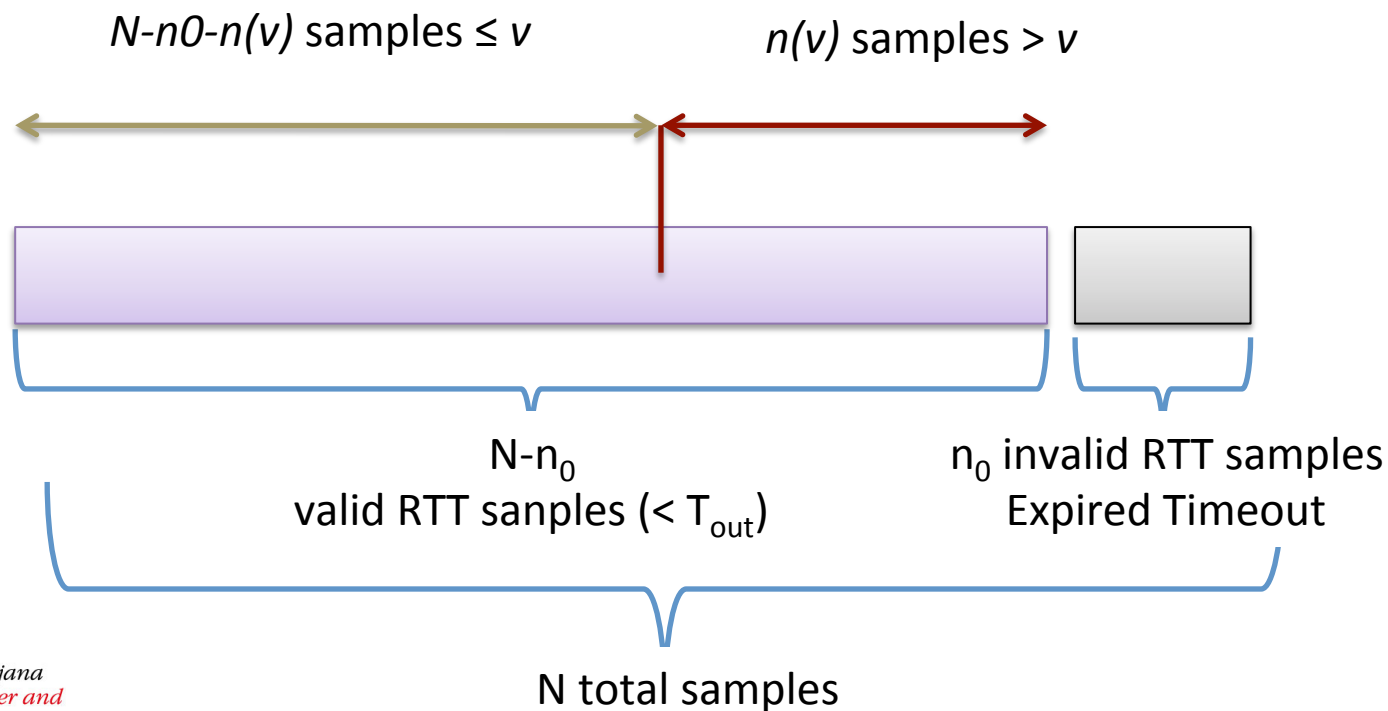
Distortive effect due to ECCDF normalization with different TimeOut value (1sec vs. 3 sec)

ECCDF of valid RTT samples (loglog)



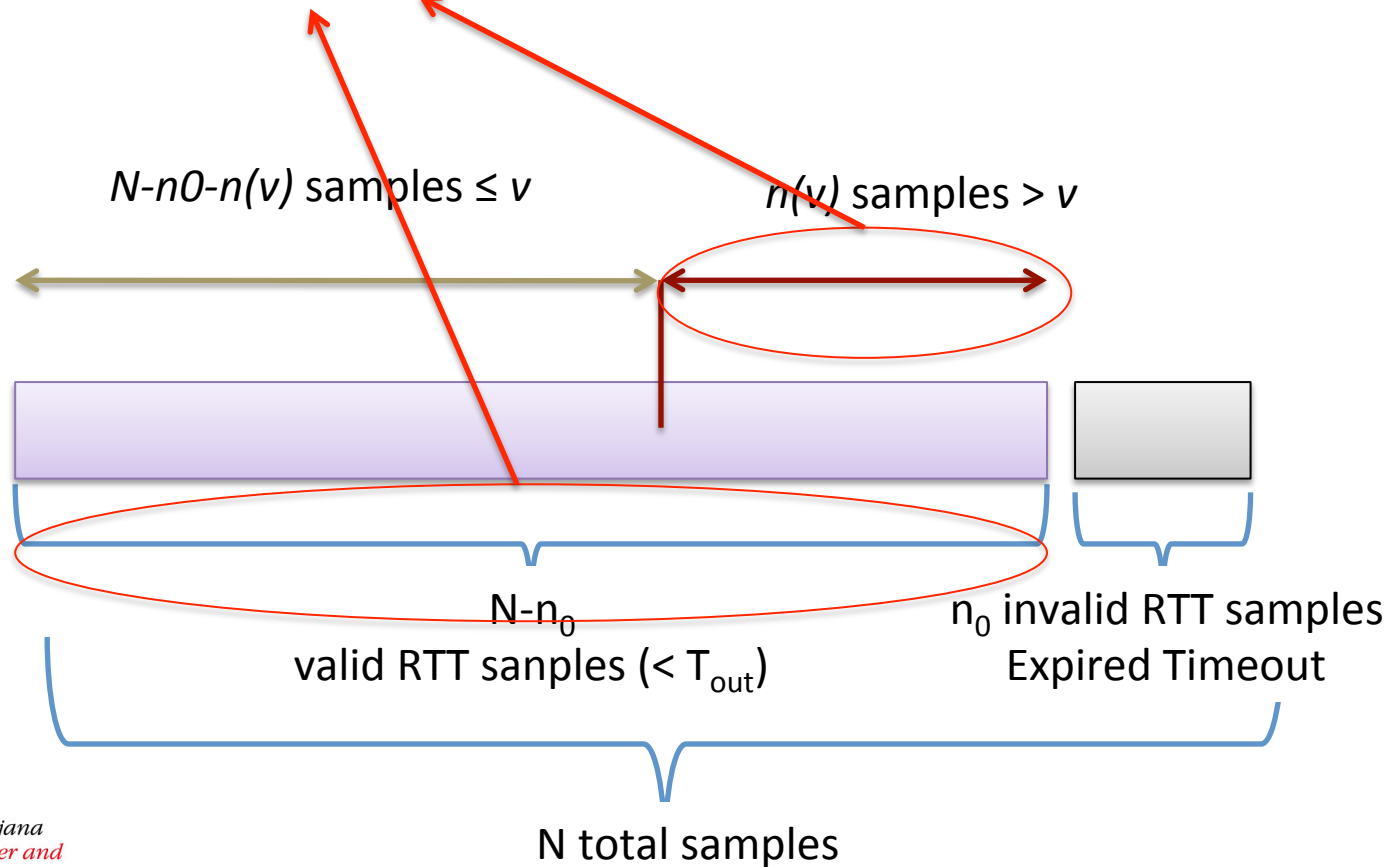
ECCDF of valid RTT samples (loglog)

$$S_k(v) = \frac{\text{number of elements } > v}{\text{total number of elements}} = \frac{1}{k} \sum_{i=1}^k \mathbf{1}\{x_i > v\}$$



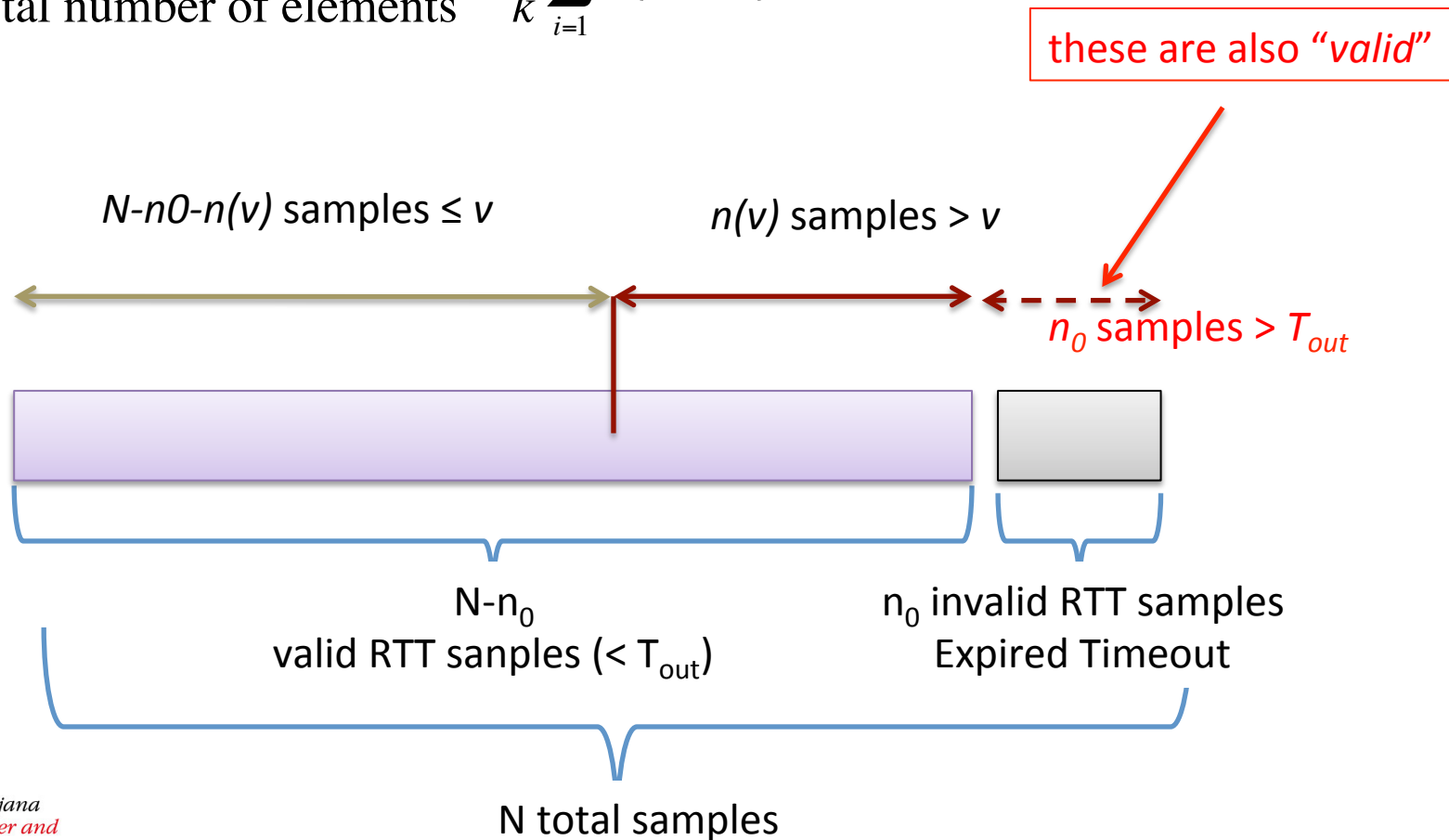
ECCDF of valid RTT samples (loglog)

$$S_k(v) = \frac{\text{number of elements } > v}{\text{total number of elements}} = \frac{1}{k} \sum_{i=1}^k \mathbf{1}\{x_i > v\}$$



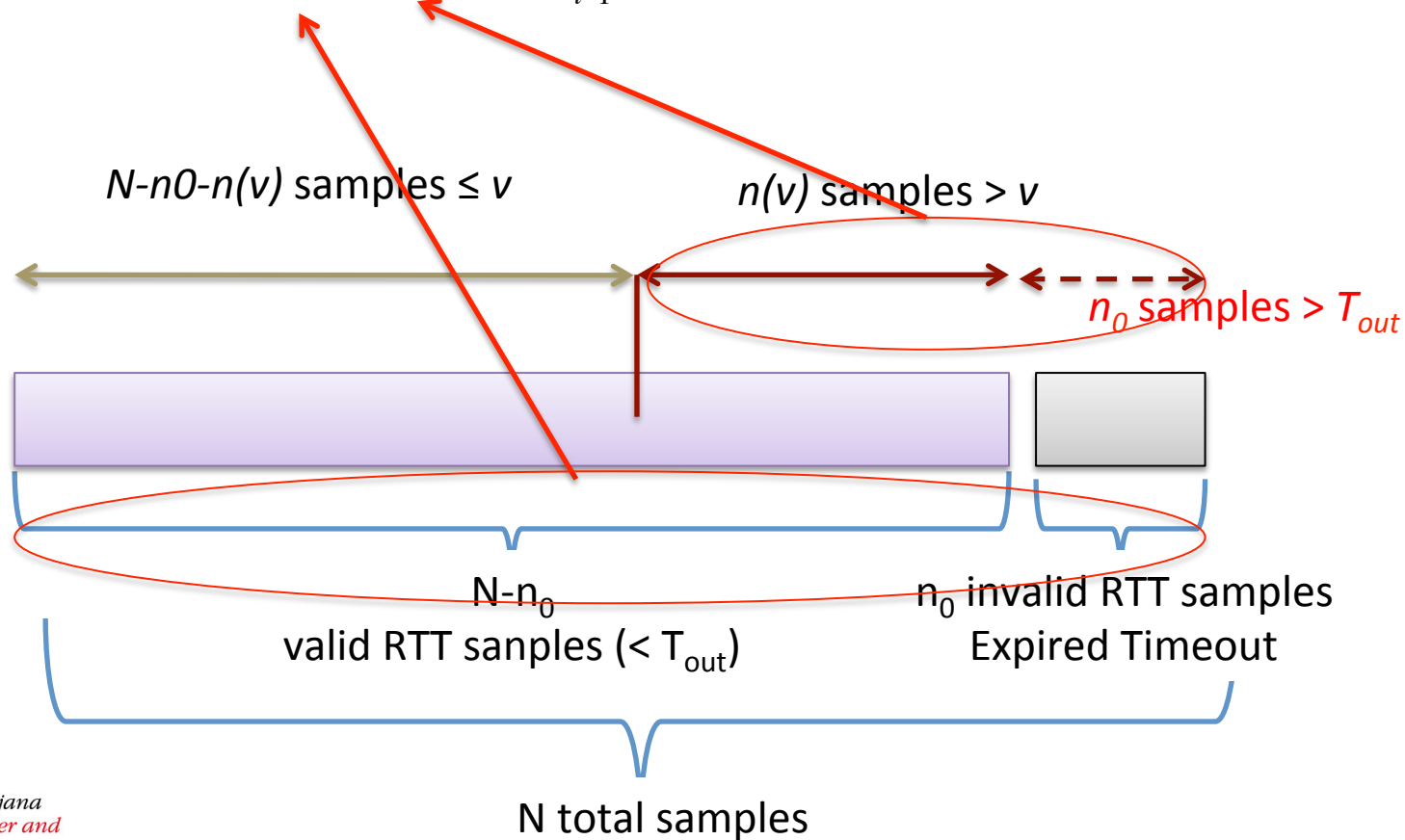
ECCDF of valid RTT samples (loglog)

$$S_k(v) = \frac{\text{number of elements } > v}{\text{total number of elements}} = \frac{1}{k} \sum_{i=1}^k \mathbf{1}\{x_i > v\}$$

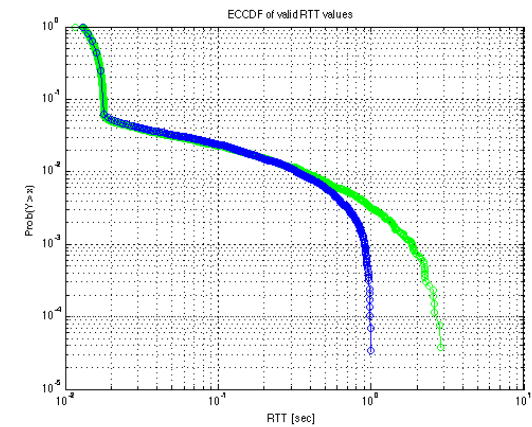
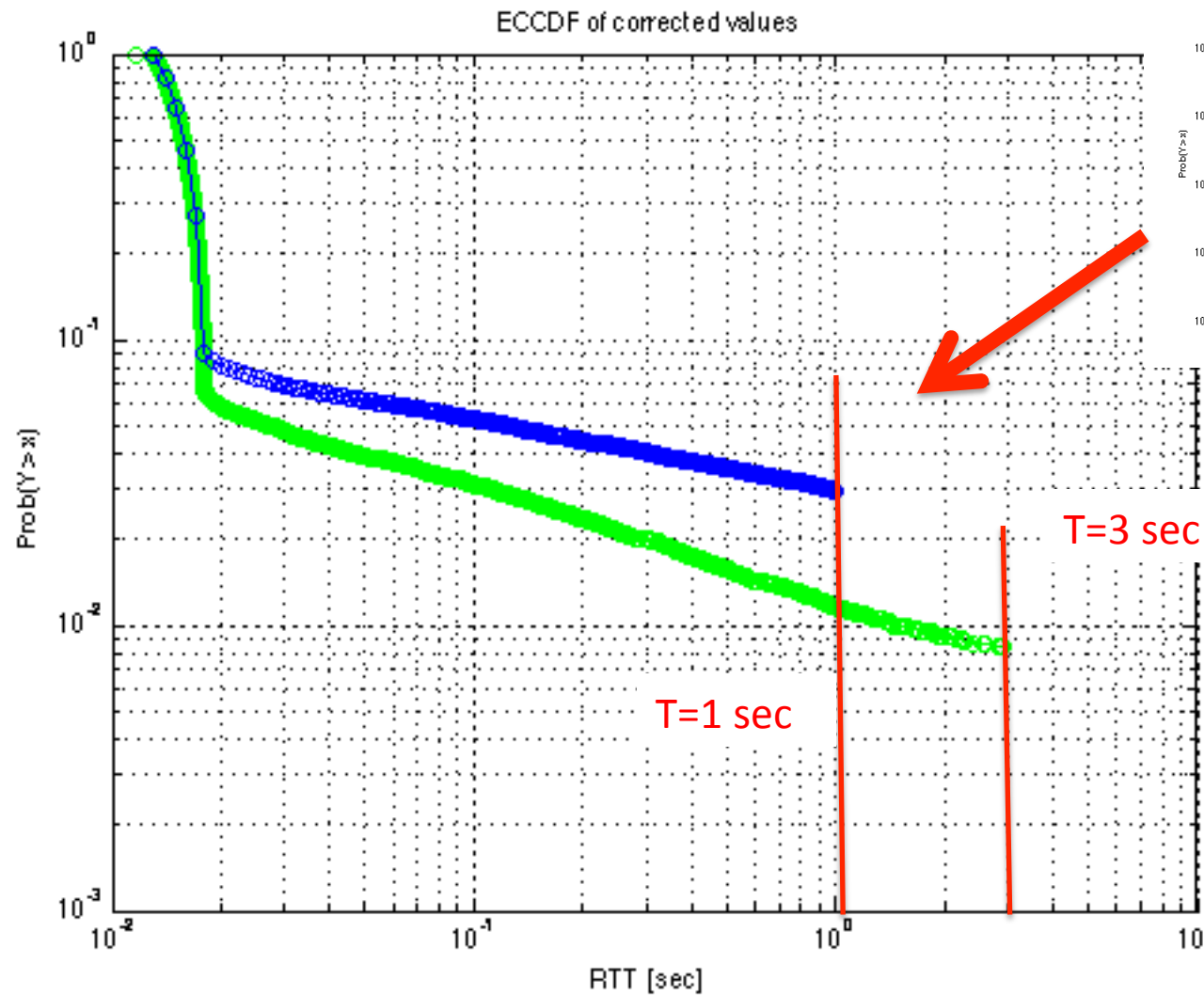


all ECDF of ~~valid~~ RTT samples (loglog)

$$S_k(v) = \frac{\text{number of elements } > v}{\text{total number of elements}} = \frac{1}{k} \sum_{i=1}^k \mathbf{1}\{x_i > v\}$$

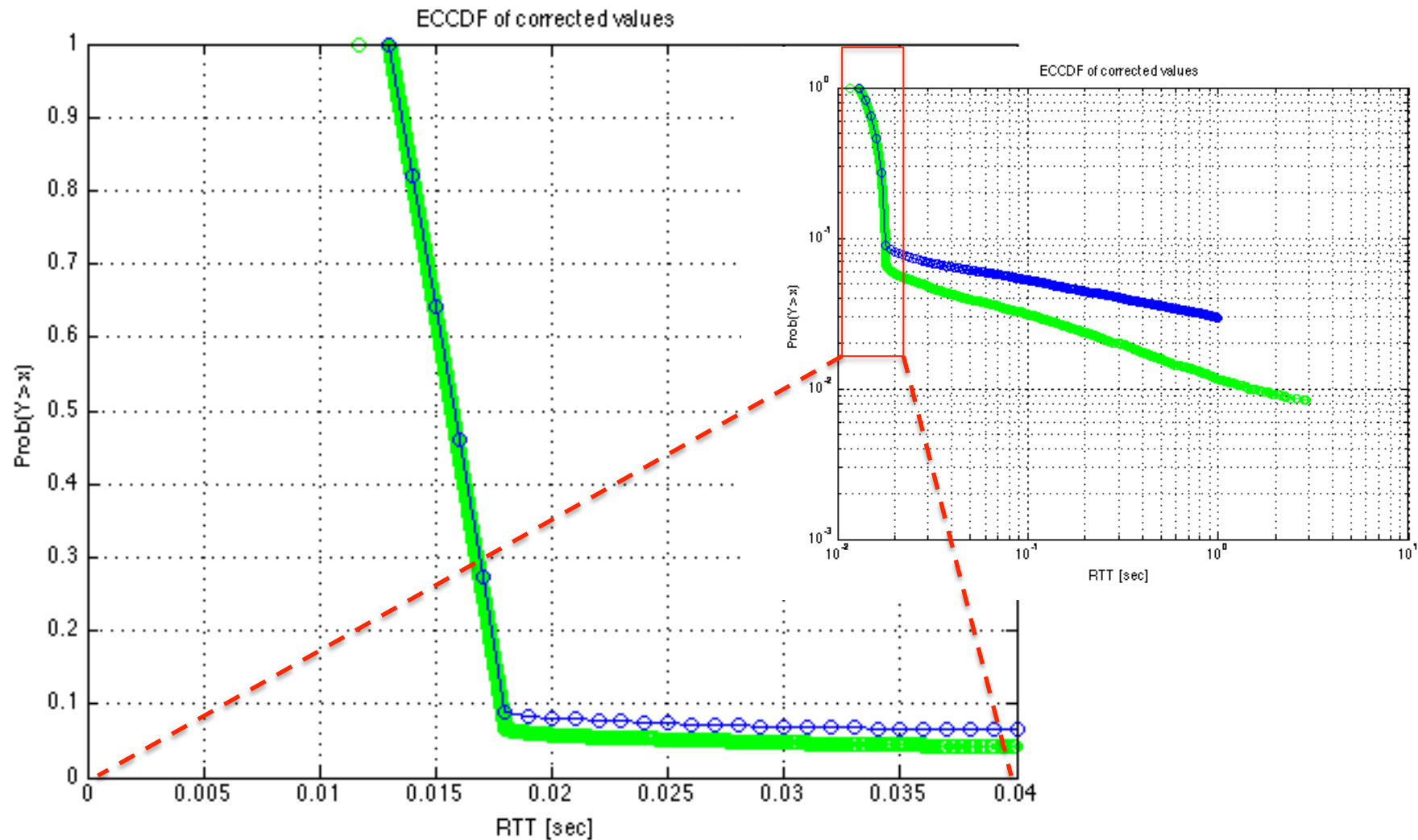


ECCDF of all RTT samples (loglog)

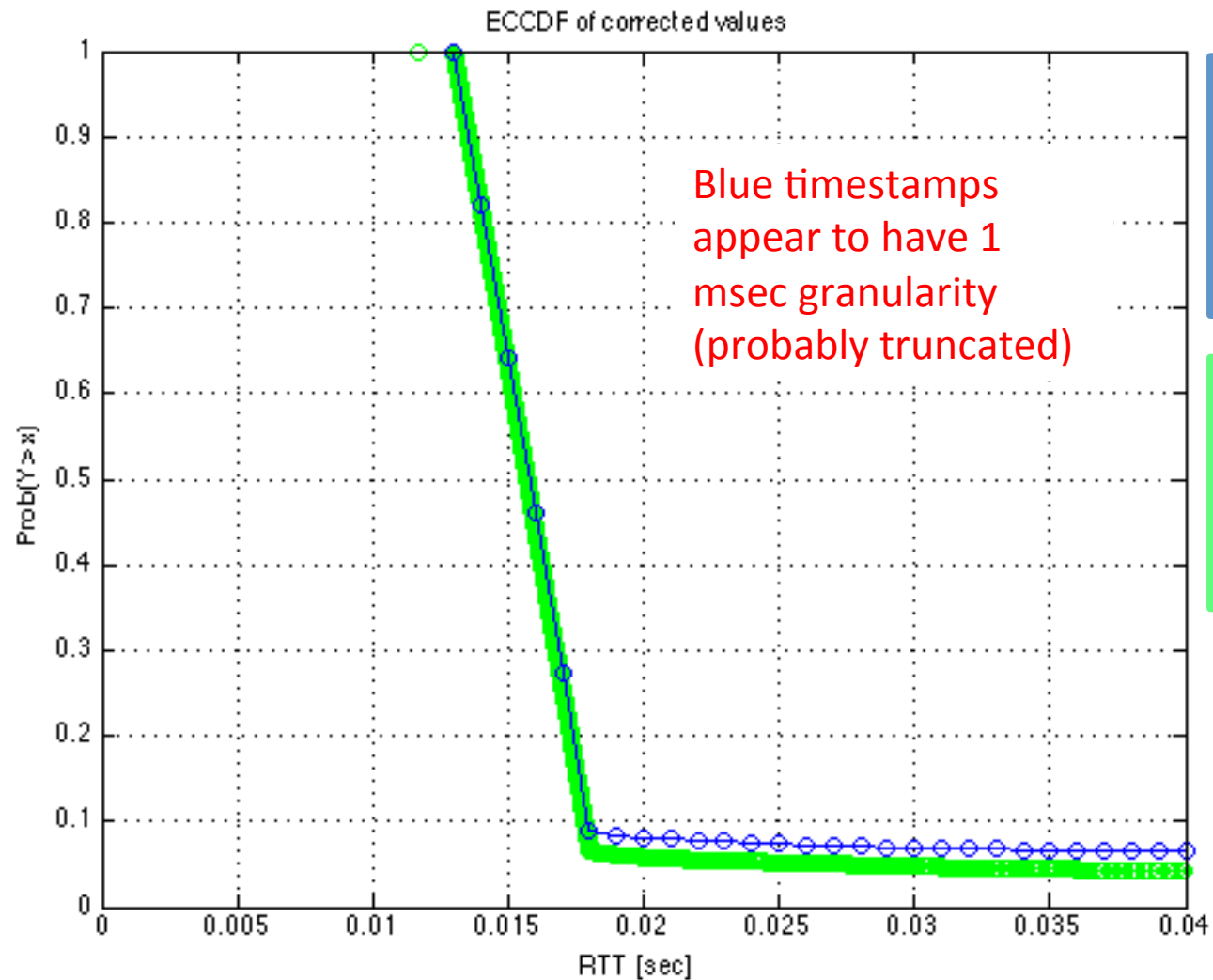


Contrary to above,
now Green is
below/left of Blue
→ Green network
has **lower** RTT !

zooming into the mass



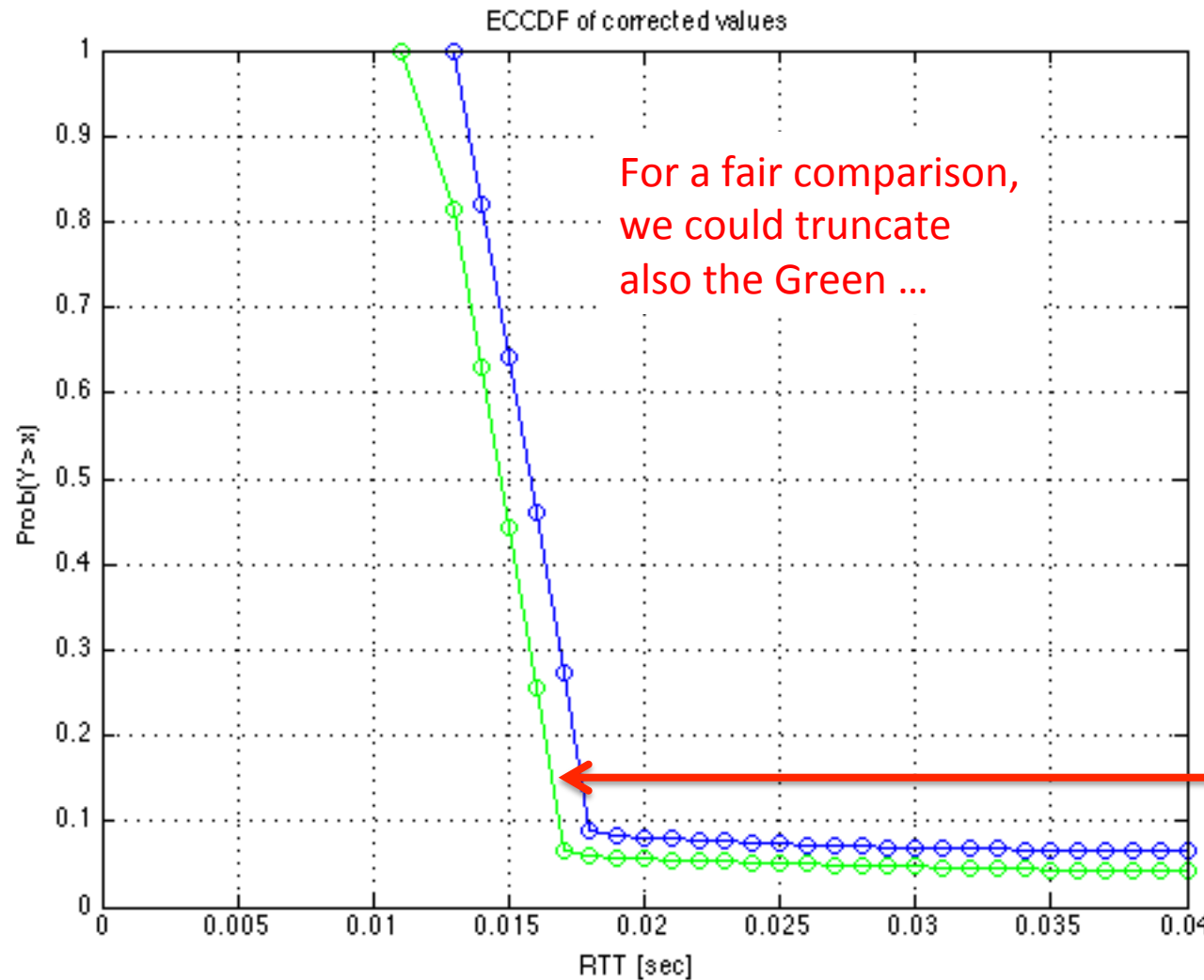
zooming into the mass



StartTime	RTT
628.0000000	0.01300000
628.0330000	0.01600000
628.0690000	0.01800000
628.1080000	0.01400000

StartTime	RTT
972.0000000	0.01172138
972.02004475	0.01476894
972.04011993	0.01303928
972.06021071	0.01749107

zooming into the mass

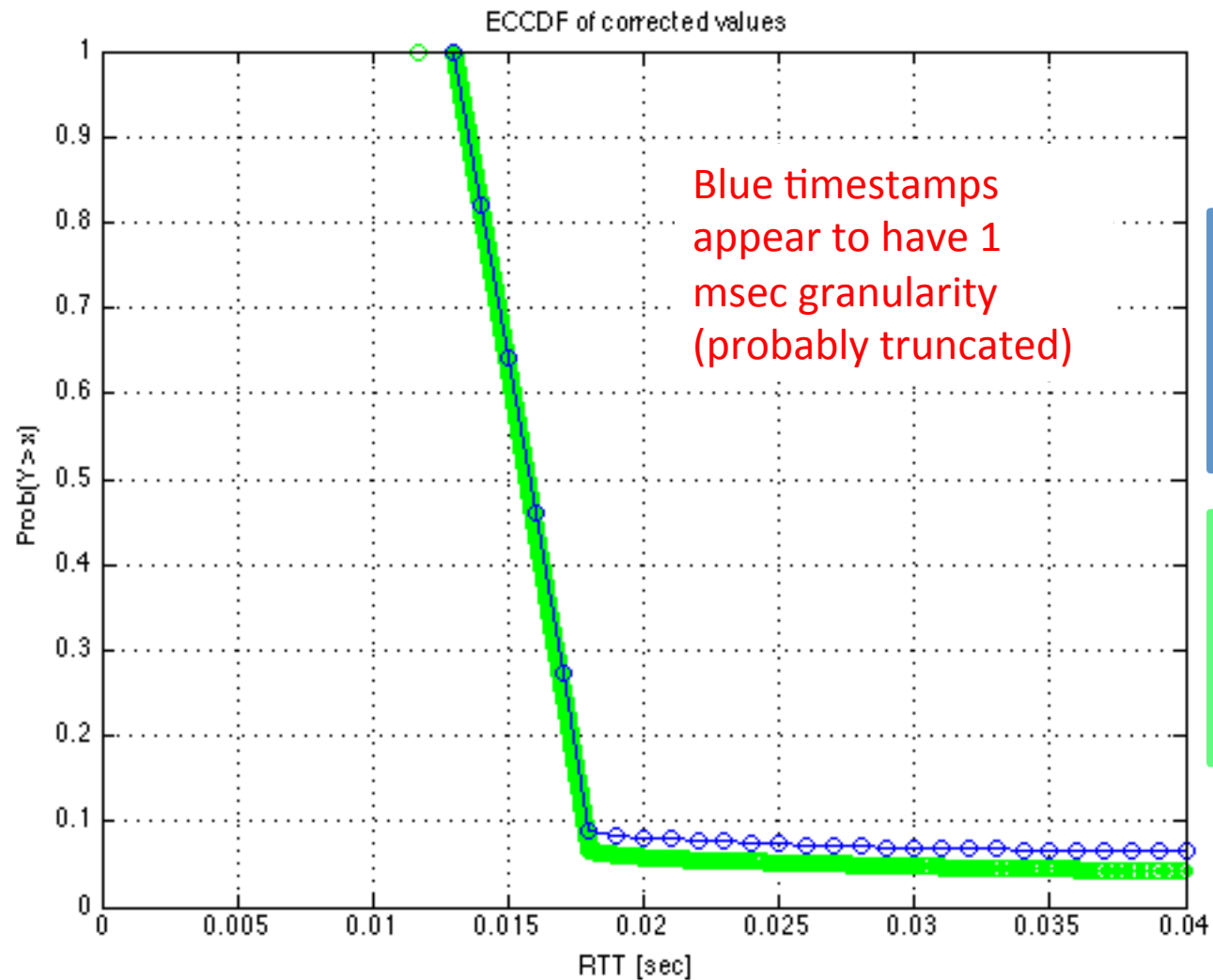


StartTime	RTT
628.0000000	0.01300000
628.0330000	0.01600000
628.0690000	0.01800000
628.1080000	0.01400000

StartTime	RTT
972.0000000	0.01172138
972.02004475	0.01476894
972.04011993	0.01303928
972.06021071	0.01749107

StartTime	RTT
972.0000000	0.01100000
972.02004475	0.01400000
972.04011993	0.01300000
972.06021071	0.01700000

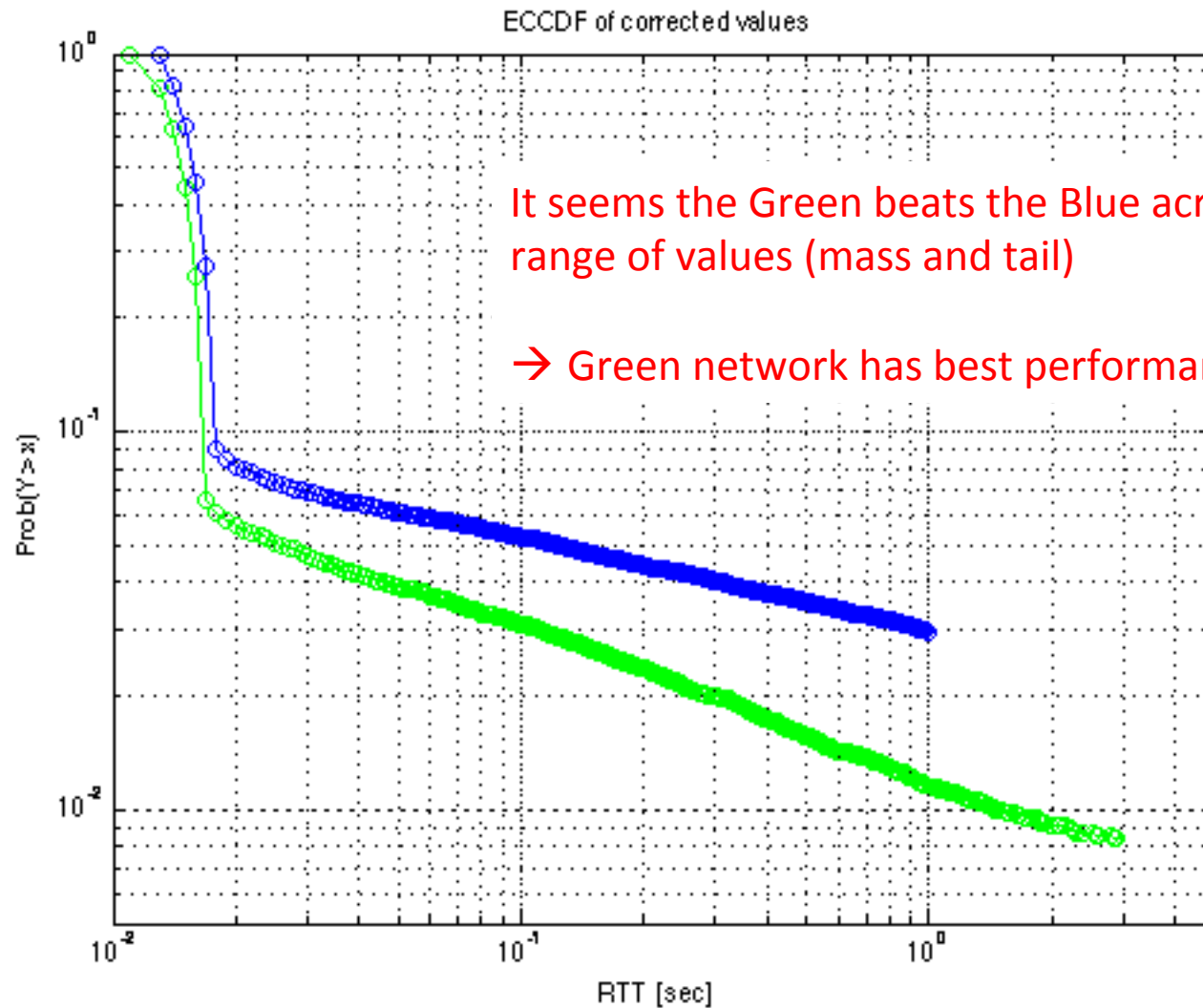
zooming into the mass



StartTime	RTT
628.0000000	0.01300000
628.0330000	0.01600000
628.0690000	0.01800000
628.1080000	0.01400000

StartTime	RTT
972.0000000	0.01172138
972.02004475	0.01476894
972.04011993	0.01303928
972.06021071	0.01749107

ECCDF of all samples (loglog)



Take-away message 1

- Meta-information = meta-data, parameter values, context data (outage), *detailed description of measurement collection method ...*
- Meta-information is important
 - meta-data not less important than “data”
- Meta-information is not always available
 - just missing, erroneous or ambiguous

... are spaced by 20 ms.



within a maximum predefined timeout,

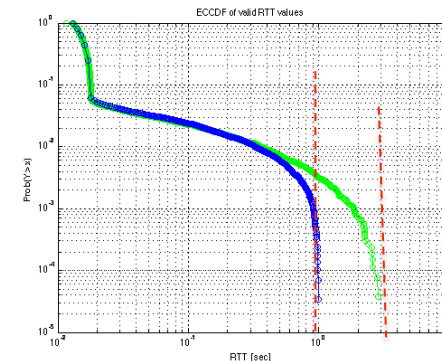
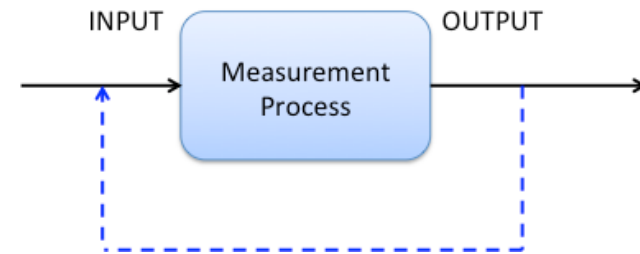
Take-away message 2

- Try to collect and record as much meta-data and context-information as possible
- Sometimes, missing information can be inferred (guessed?) from the data
 - in our example: outage, timeout values
 - reverse engineering the collection method: “20 ms spacing”, 1ms truncation



Take away message 3

- Watch against spurious correlations, sources of bias, distortion, artifacts ...
- in the *collection* phase ...
- ... and during the *analysis*



Take away message 4

- Different problems / artifacts might interact in subtle ways
 - e.g. outage + I/O feedback → over/under representation



Final words

- Torture your data. But first caress them!
- Have fun with your future dialogues with real data 😊
- For questions & feedback email to:
fabio.ricciato@fri.uni-lj.si

