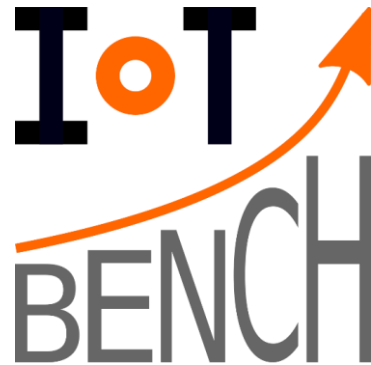


Towards a Methodology for Experimental Evaluation in Low-Power Wireless Networking



Romain Jacob

Usman Raza

Lothar Thiele

Carlo Alberto Boano

Marco Zimmerling

CPS-IoTBench Workshop

April 15, 2019

Includes material from Hanspeter Schmid and Alex Huber

“ We need a benchmark for IoT networking. ”

“ We need a benchmark for IoT networking. ”

⇔ Comparing performance

“ We need a benchmark for IoT networking. ”

⇔ Comparing performance

⇒ Repeatable experiments

“ We need a benchmark for IoT networking. ”

- ⇔ Comparing performance
- ⇒ Repeatable experiments
- ⇒ Formalize the experimental methodology

Fact 1	The RF environment affects performance of low-power protocol in unpredictable ways.*
+	
Fact 2	Real RF environment cannot be controlled.
=	
Consequence	Performance variability is expected.

* Either unfeasible or unpractical to model

6



How do you handle variability?

“Repeat your experiment.”

Easy, right?

How do you handle variability?

“Repeat your experiment.”

Easy, right?

Killer
questions

How long should be your experiment?

How many times should you repeat it?

How do you handle variability?

“Repeat your experiment.”

Easy, right?

How long should be your experiment?

How many times should you repeat it?



“Run many long tests.”

Not so easy...

Let us assume

You ran “many long tests”

1 M samples of XYZ
spread over a
large period of time



How do you synthesize your results?

“Use statistics.”

Let us assume

You ran “ many long tests ”

1 M samples of XYZ
spread over a
large period of time

How do you synthesize your results?

“ Use statistics. ”

Literally

A piece of data obtained
from a large quantity of data

Mean, median,
standard deviation, etc.

Beware!

Descriptive
statistics

≠

Predictive
statistics

What is the collected data like

What the collected data
allows to **infer** about
future/other data (unknown)

Beware!

Descriptive
statistics

≠

Predictive
statistics

What is the collected data like

What the collected data
allows to **infer** about
future/other data (unknown)

The “interesting case”

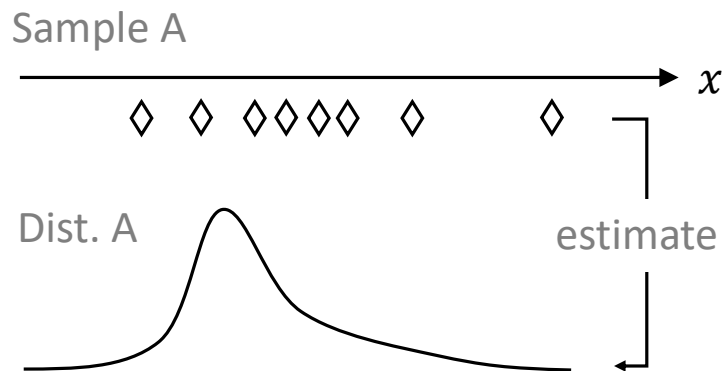
13

Descriptive statistics compare the samples



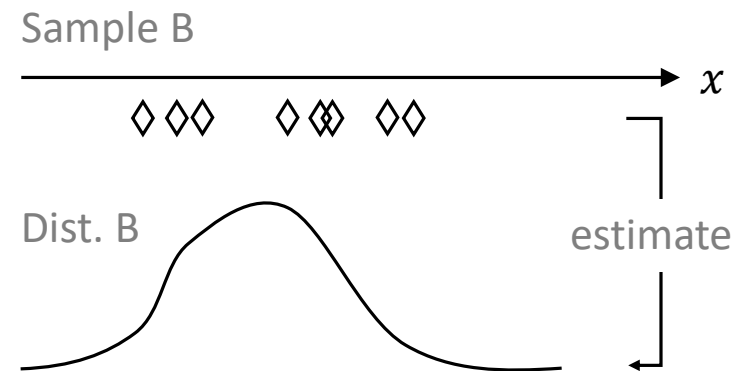
Conclusion Sample A is “better” than sample B.

Predictive statistics
(aims to) compare the **underlying distributions**



Conclusion

A is “better” than B.



If one sample A and B, then likely
sample A is “better” than sample B.

Descriptive
statistics

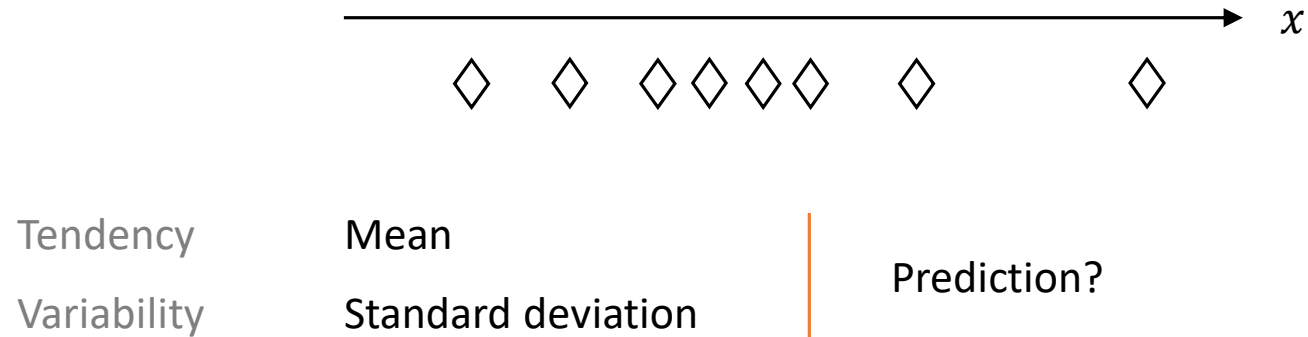
Sample A is “better”
than sample B.

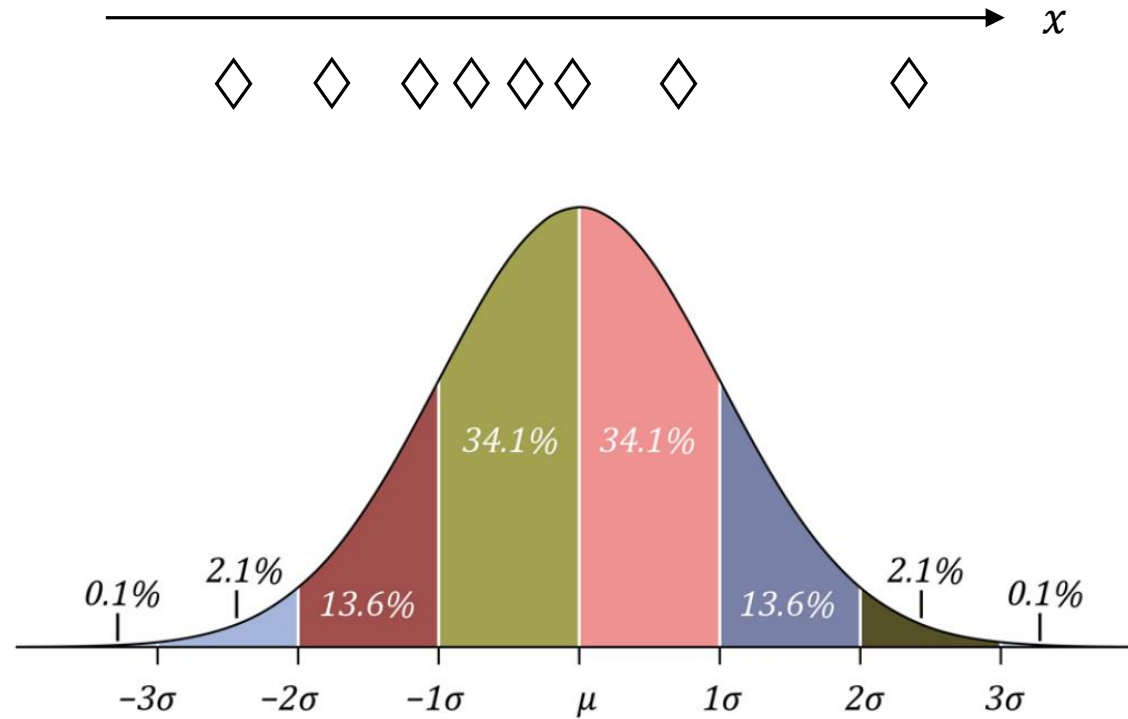
\neq

Predictive
statistics

If one sample A and B,
then likely sample A
is “better” than sample B.

Much stronger statement
but also harder to make





If we use mean and standard deviation in this context,
we make **two mistakes**

First error

The standard deviation of the **sample** is not
the standard deviation of the **underlying distribution**.

Use confidence intervals

Confidence interval (CI) ?

Informally A numerical **interval**
in which **lies the true value** (which you don't know)
of some parameter **with some probability** (or confidence level)

Example $[a, b]$ is a 95% CI for the mean of x
 \Leftrightarrow The probability that the true mean value of x
is included in $[a, b]$ is larger or equal to 95%

If we use mean and standard deviation in this context,
we make two mistakes

First error

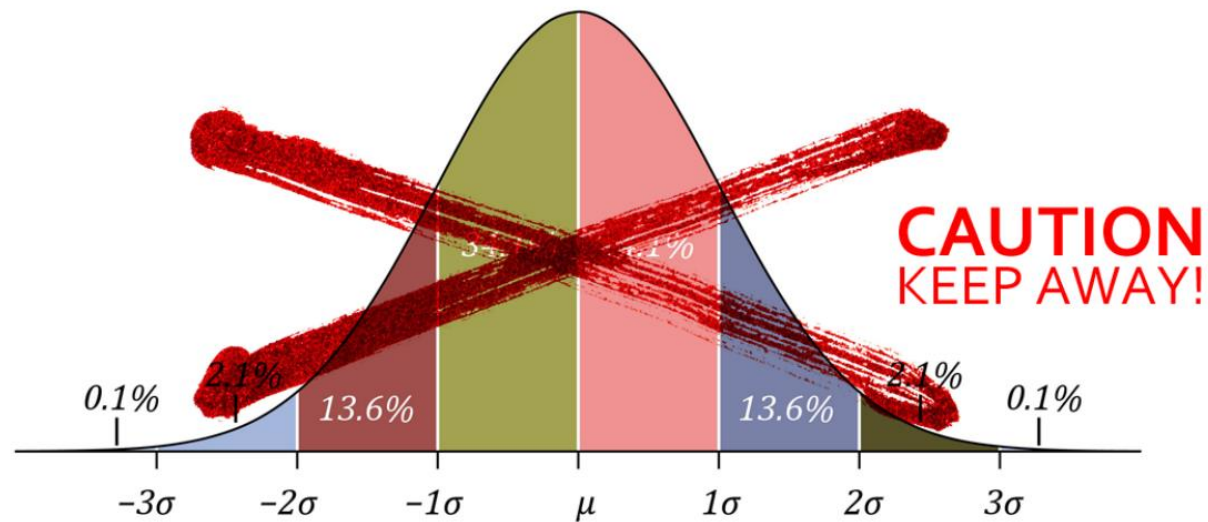
The standard deviation of the **sample** is not
the standard deviation of the **underlying distribution**.

Use confidence intervals

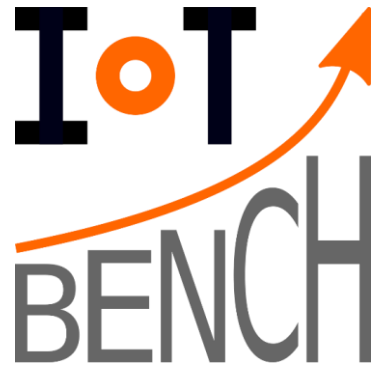
Second error

The underlying **distribution is not normal!**

The standard deviation does not help making predictions
Use **non-parametric statistics**



Towards a Methodology for Experimental Evaluation in Low-Power Wireless Networking

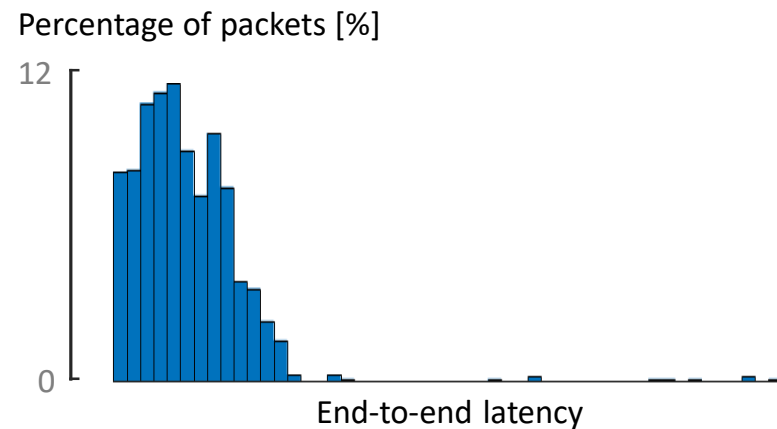
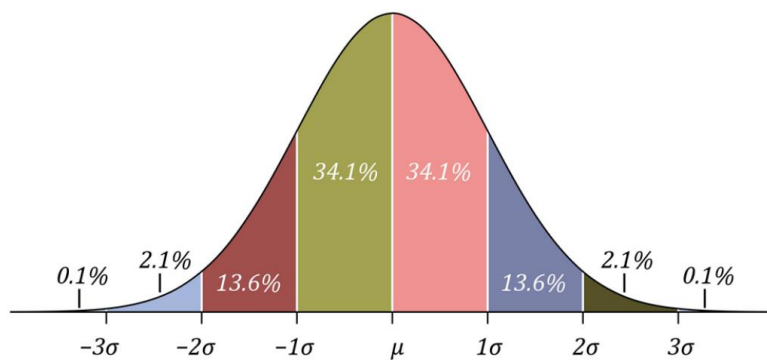


Know your data

Use non-parametric statistics

Formalizing low-power wireless
experimental evaluation

Performance measurements in computer science
are typically **not normally distributed**



Performance measurements in computer science
are typically **not normally distributed**

Use non-parametric statistics
based on distribution percentiles

Great for
predictive statistics

p -th percentile
or P_p

$p\%$

of the distribution is below

$(1 - p)\%$

of the distribution is above

Percentiles are powerful predictive statistics

Simple to use

Can compute CI for any percentile with any confidence

Distribution independent

Estimates are valid regardless of the underlying distribution

Robust

Estimates are not skewed by outliers

Confidence intervals

$$\mathbf{P} \left\{ x_m \leq M \leq x_{N-m+1} \right\} = 1 - 2 \sum_{k=0}^{m-1} \binom{N}{k} \frac{1}{2^N}$$

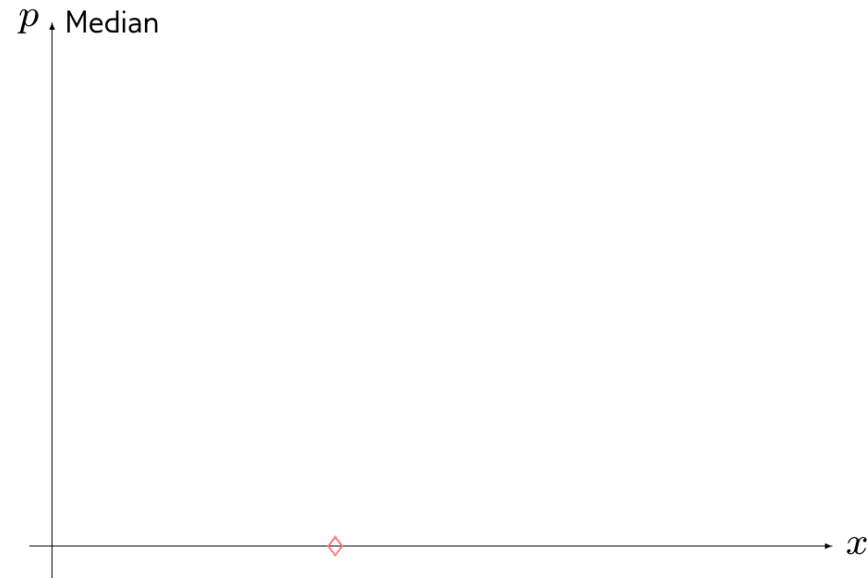
$$\mathbf{P} \left\{ x_m \leq P_p \right\} =$$

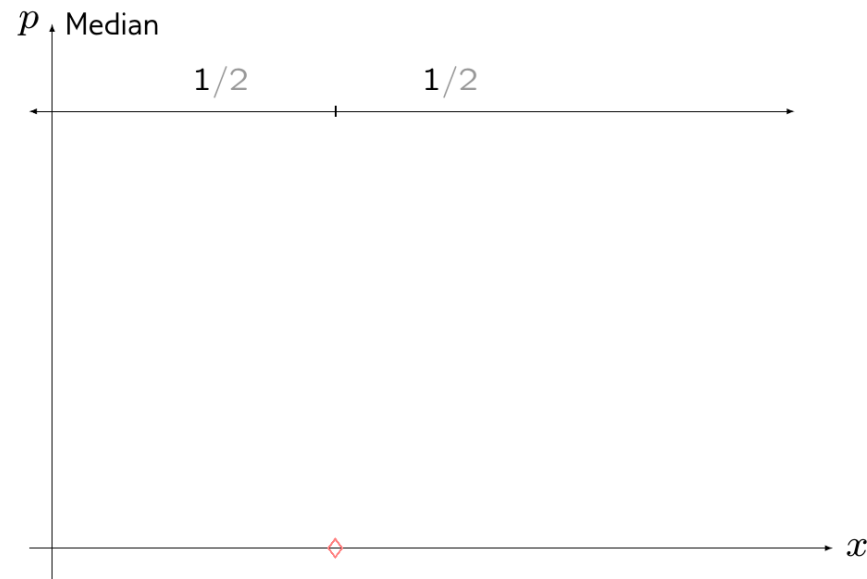
$$P \left\{ x_{N-m+1} \geq P_{1-p} \right\} = 1 - \sum_{k=0}^{m-1} \binom{N}{k} p^k (1-p)^{N-k}$$

Confidence intervals

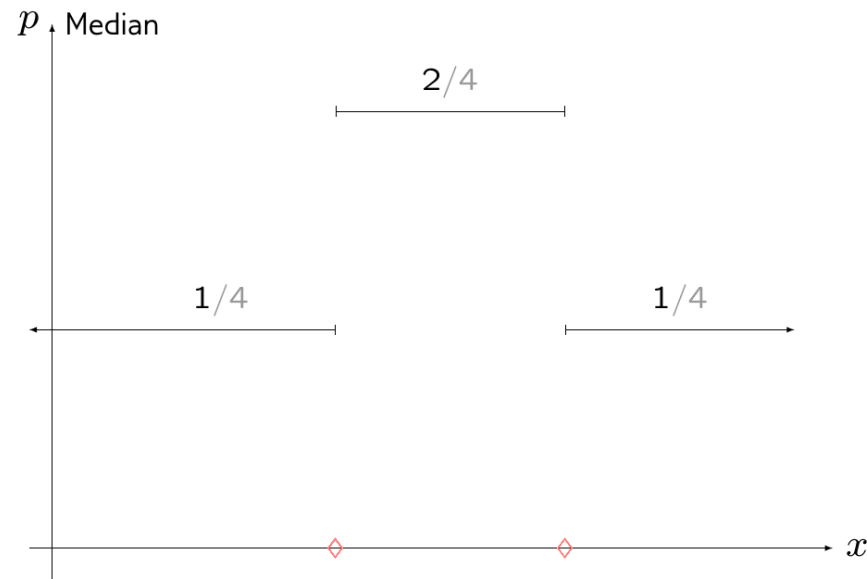
$$\mathbb{P} \left\{ x_m \leq M \leq x_{N-m+1} \right\} = 1 - 2 \sum_{k=0}^{m-1} \binom{N}{k} \frac{1}{2^N}$$

$$\begin{aligned} \mathbb{P} \left\{ x_m \leq P_p \right\} &= \\ \mathbb{P} \left\{ x_{N-m+1} \geq P_{1-p} \right\} &= 1 - \sum_{k=0}^{m-1} \binom{N}{k} p^k (1-p)^{N-k} \end{aligned}$$

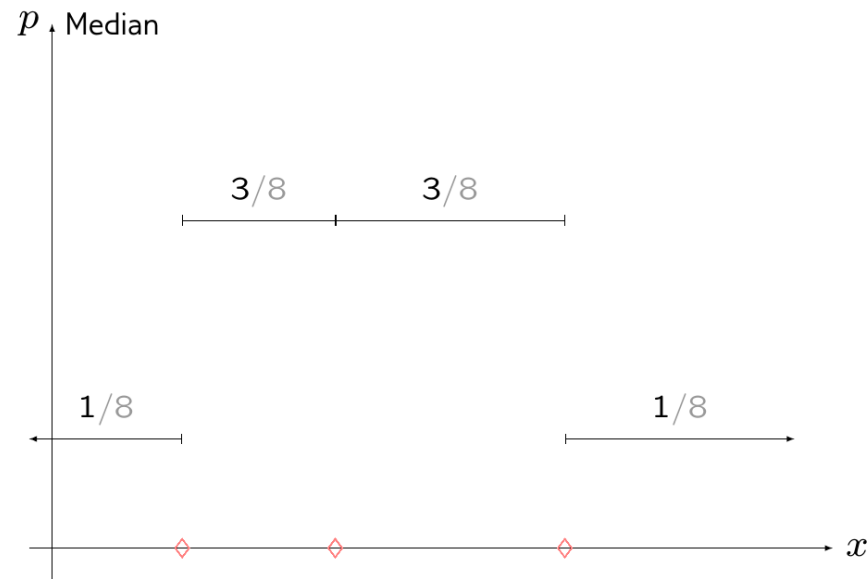




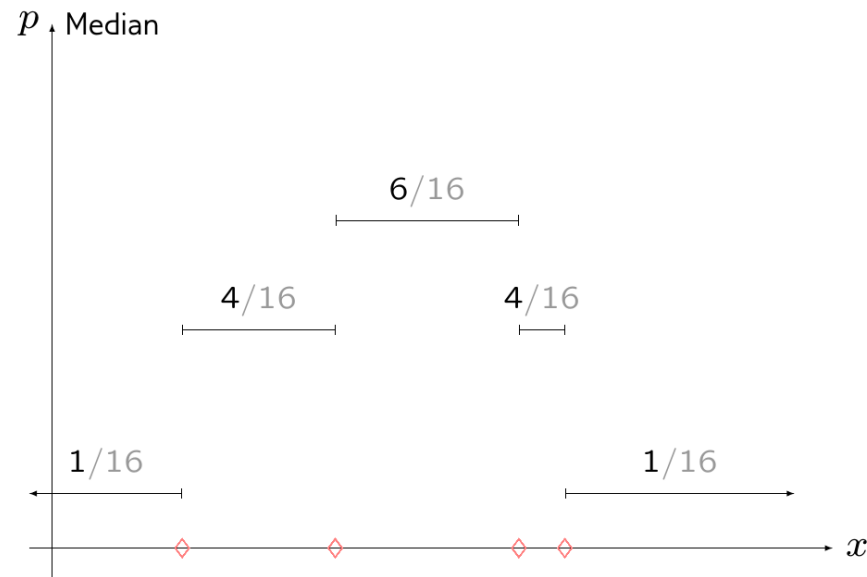
Hypothesis Samples are **i.i.d.**



Hypothesis Samples are **i.i.d.**



Hypothesis Samples are **i.i.d.**



Binomial distribution

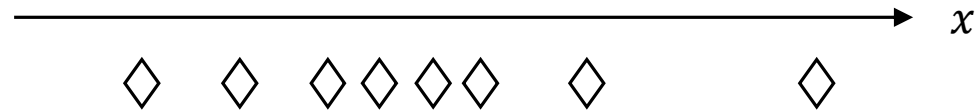
$$\mathbb{P} \{x_k \leq P_p \leq x_{k+1}\} = \binom{N}{k} p^k (1-p)^{N-k}$$

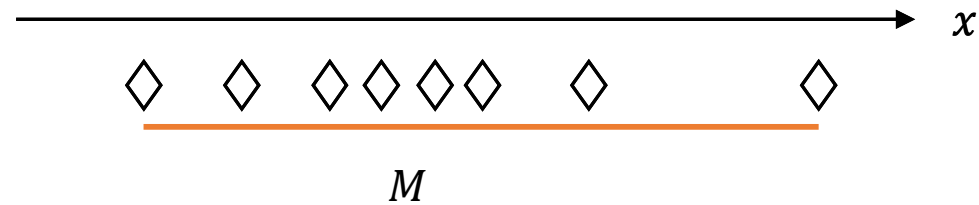
Confidence intervals

$$\mathbb{P} \left\{ x_m \leq M \leq x_{N-m+1} \right\} = 1 - 2 \sum_{k=0}^{m-1} \binom{N}{k} \frac{1}{2^N}$$

$$\mathbb{P} \left\{ x_m \leq P_p \right\} =$$

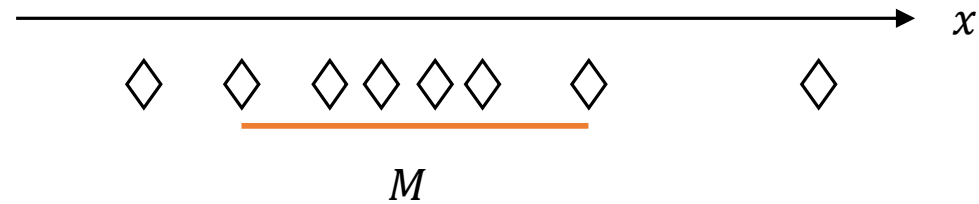
$$\mathbb{P} \left\{ x_{N-m+1} \geq P_{1-p} \right\} = 1 - \sum_{k=0}^{m-1} \binom{N}{k} p^k (1-p)^{N-k}$$





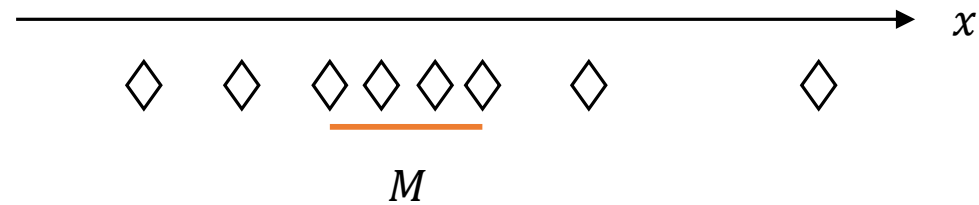
$$P = 99,29 \%$$

when 2 x 0 points are excluded



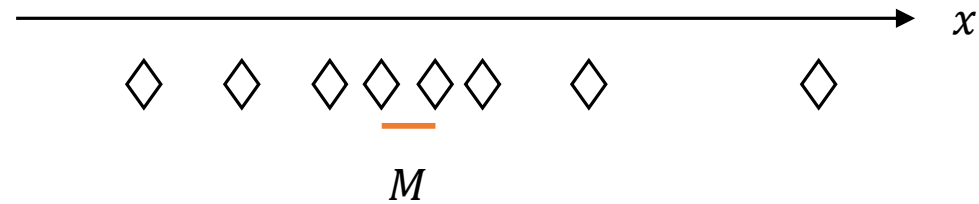
$$P = 92,97 \%$$

when 2 x 1 points are excluded



$$P = 71,09 \%$$

when 2 x 2 points are excluded



$$P = 27,34 \%$$

when 2 x 3 points are excluded

Percentiles are powerful predictive statistics

Simple to use

Can compute CI for any percentile with any confidence

Distribution independent

Estimates are valid regardless of the underlying distribution

Robust

Estimates are not skewed by outliers

Percentiles are powerful predictive statistics

Simple to use

Can compute CI for
any percentile with **any confidence**

Distribution independent

Estimates are valid regardless
of the underlying distribution

Robust

Estimates are not
skewed by outliers

For any confidence c , any percentile p ,

$$N \geq \frac{\log(1 - c)}{\log(1 - p)}$$

Example

$$c = 0,95$$

$$p = 0,01 \text{ or 1-th percentile}$$

$$\Rightarrow N \geq 299$$

Thus

We can derive the **minimal number of samples** required for estimating any percentile with any confidence

Thus

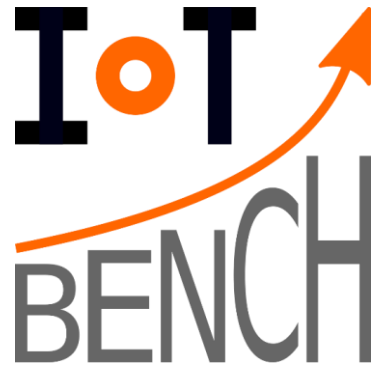
We can derive the **minimal number of samples** required for estimating any percentile with any confidence

So now

How long should be your experiment?
How many times should you repeat it?

Can be answered **rationally**.

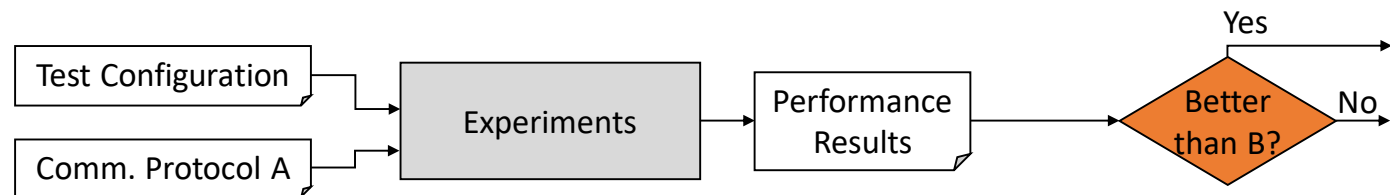
Towards a Methodology for Experimental Evaluation in Low-Power Wireless Networking

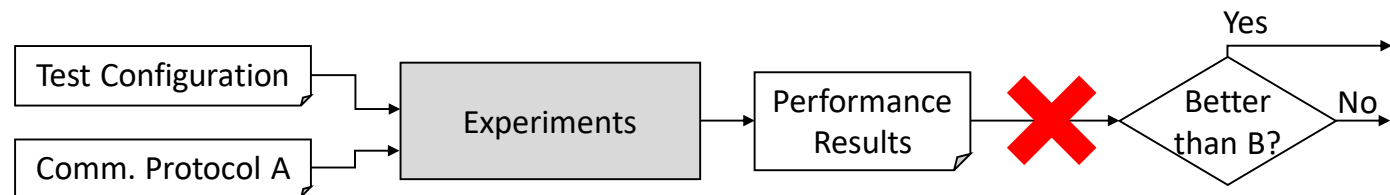


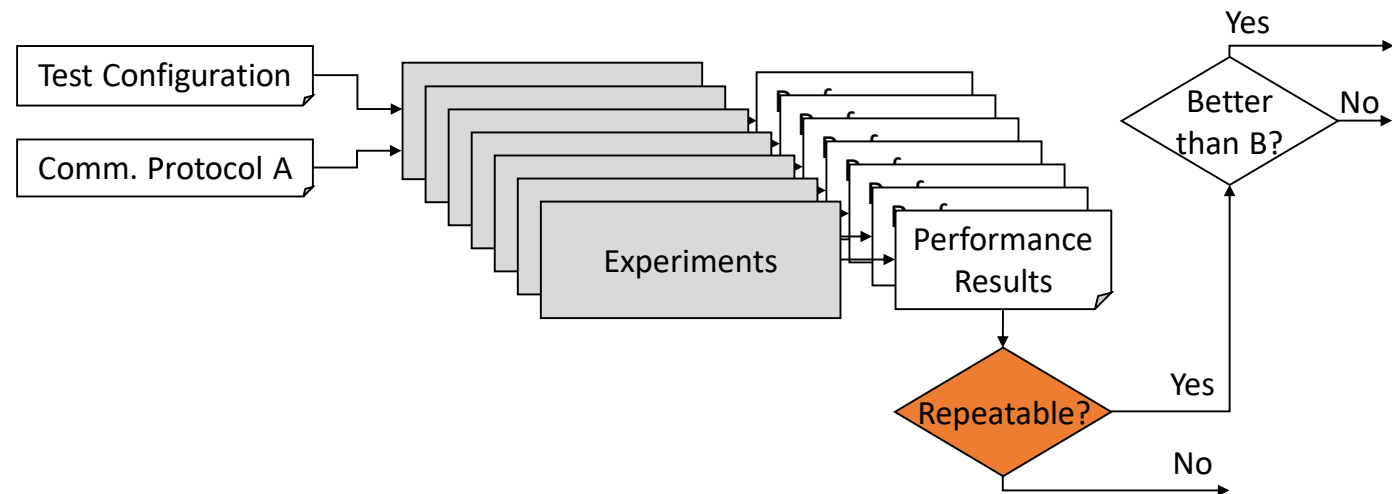
Know your data

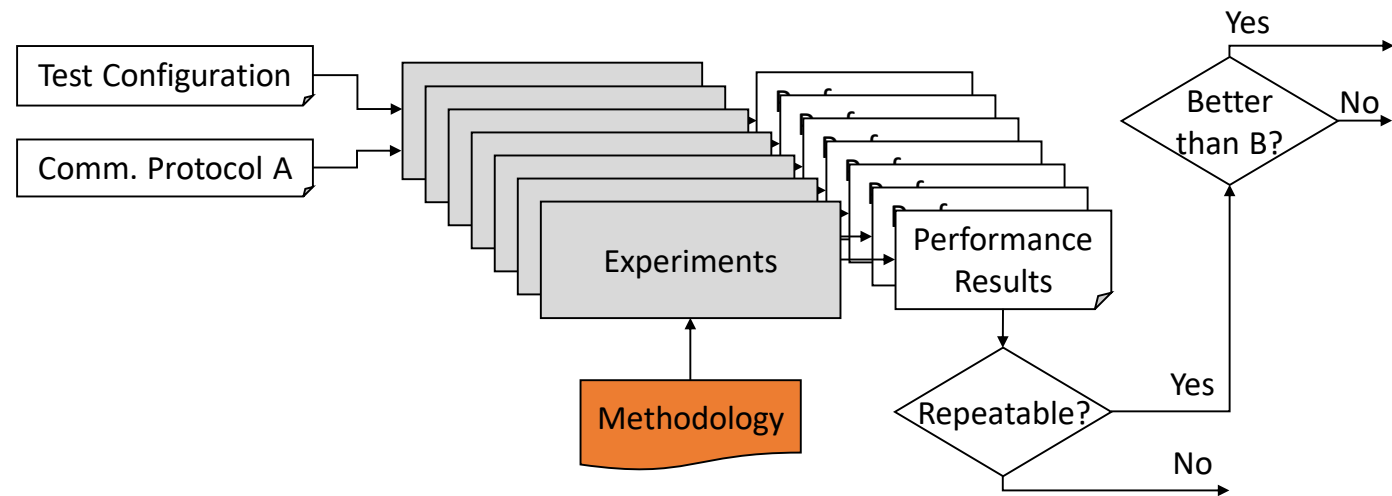
Use non-parametric statistics

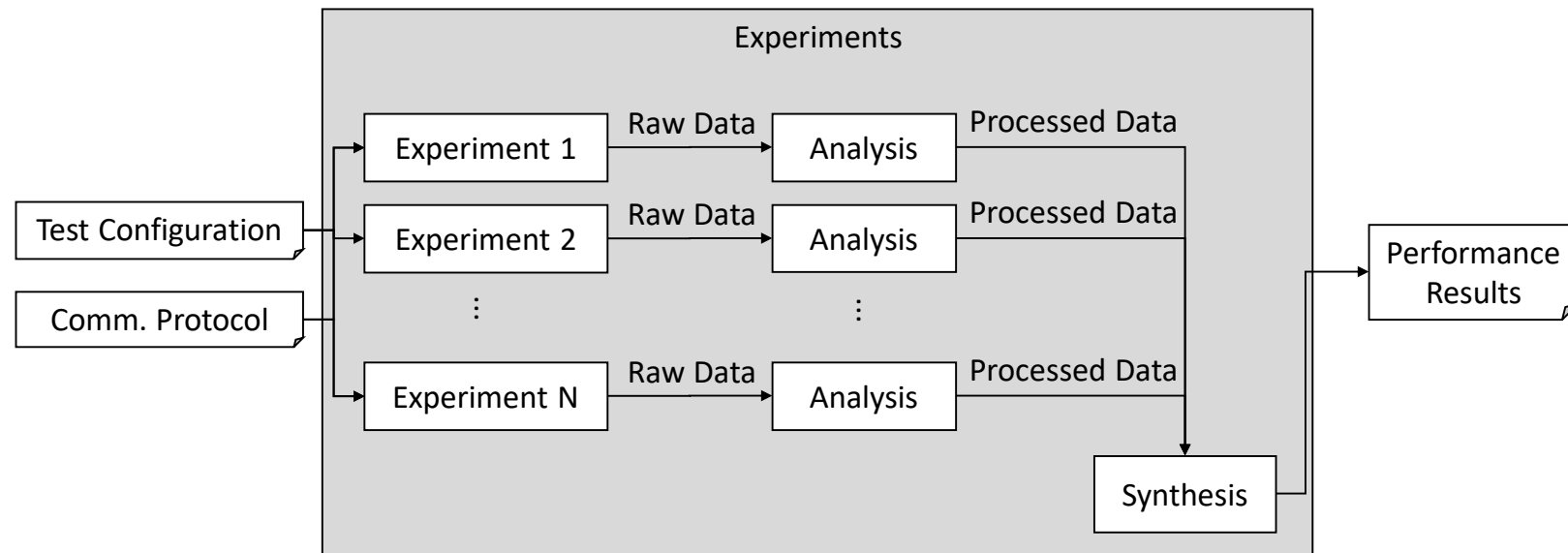
Formalizing low-power wireless
experimental evaluation

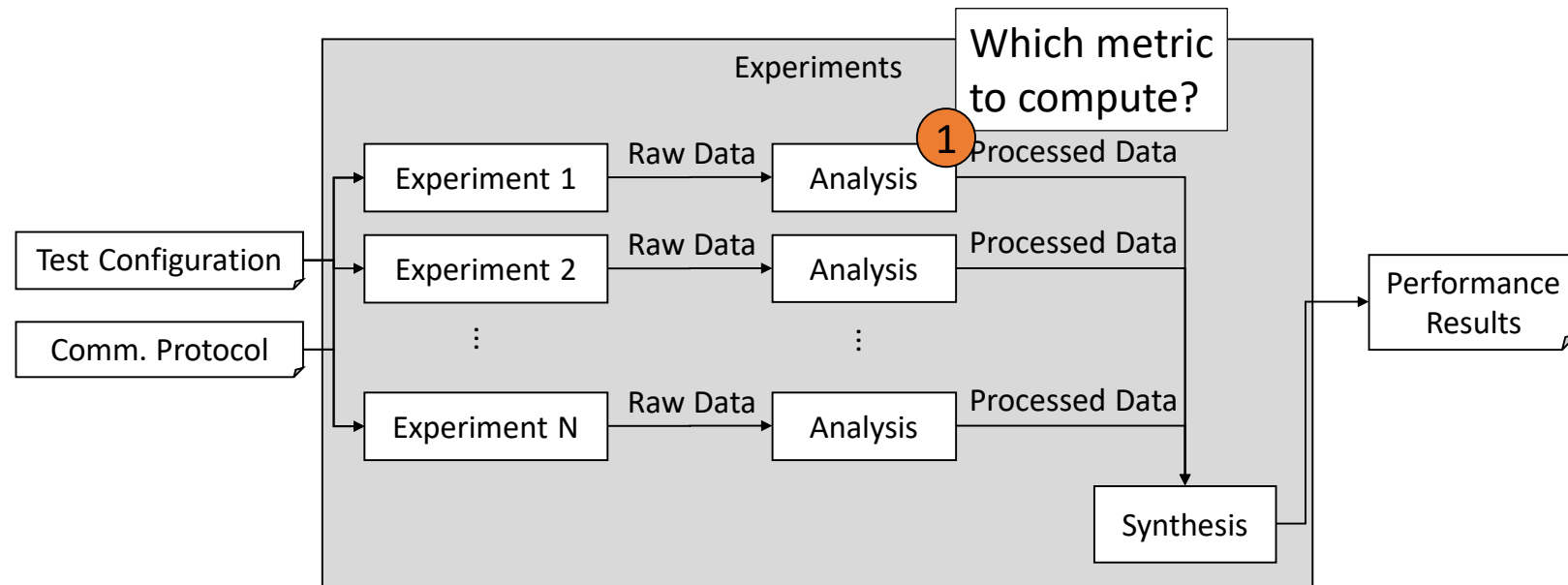


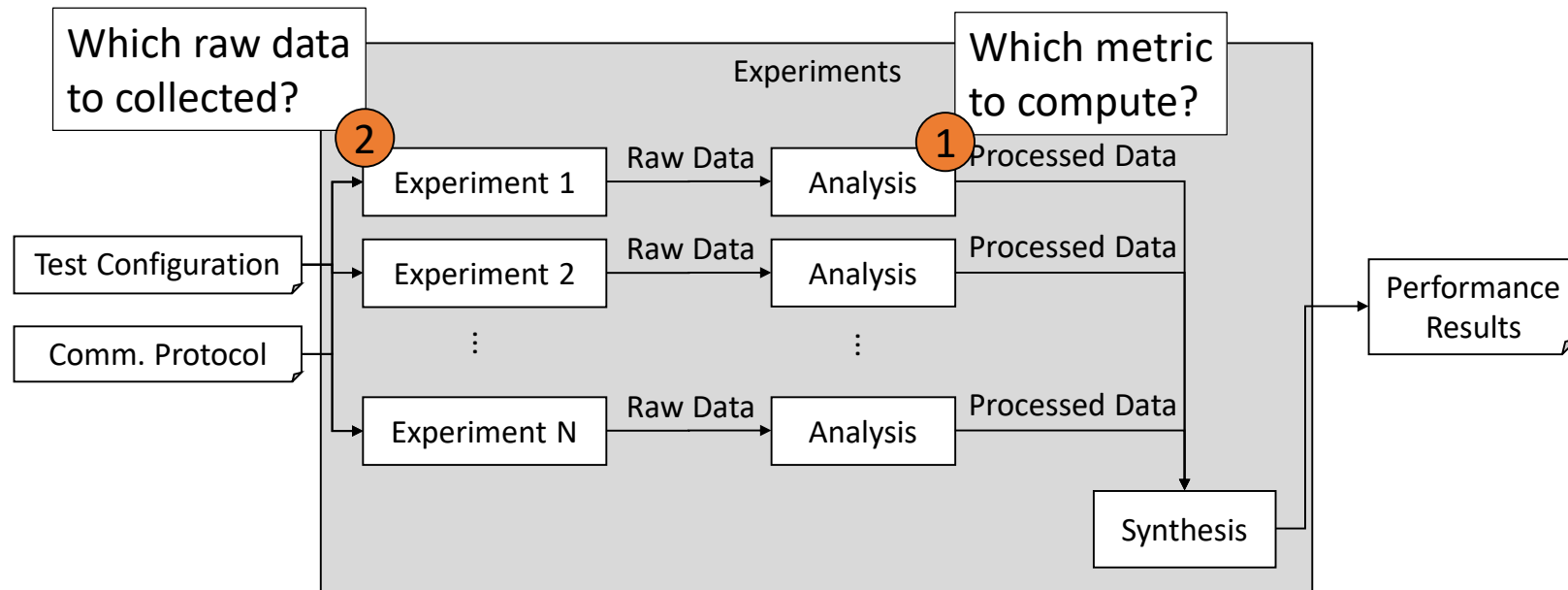


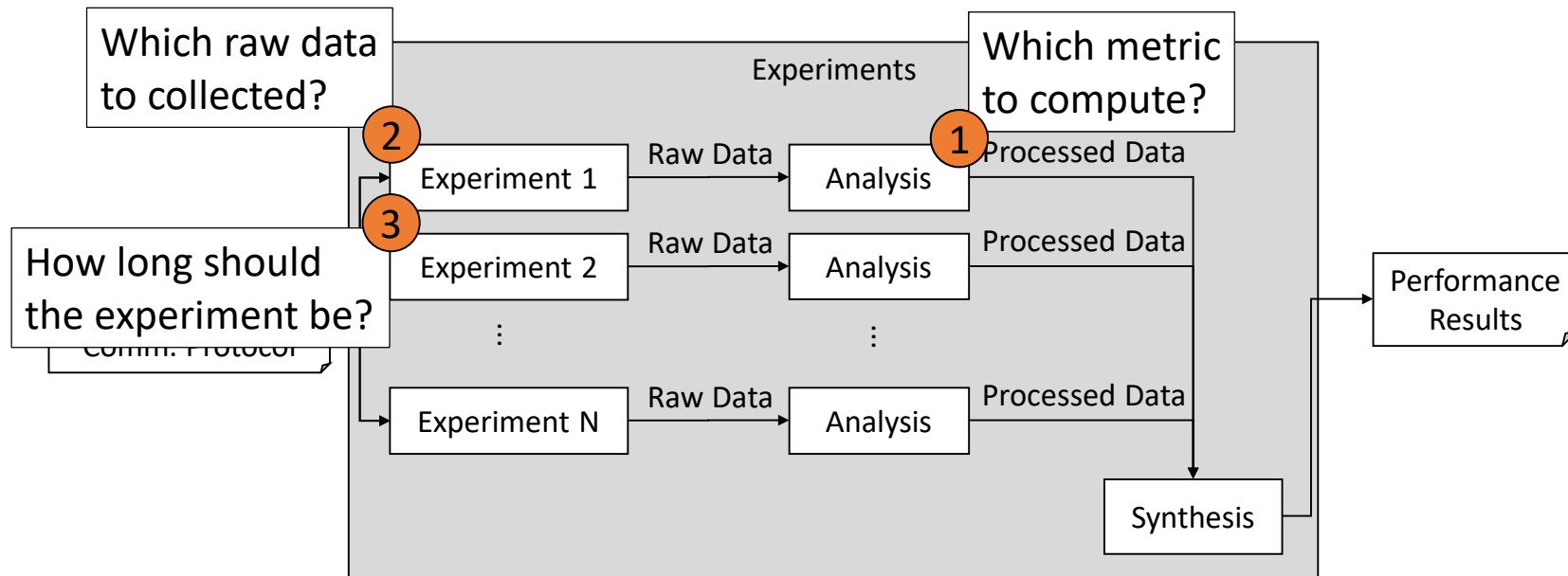


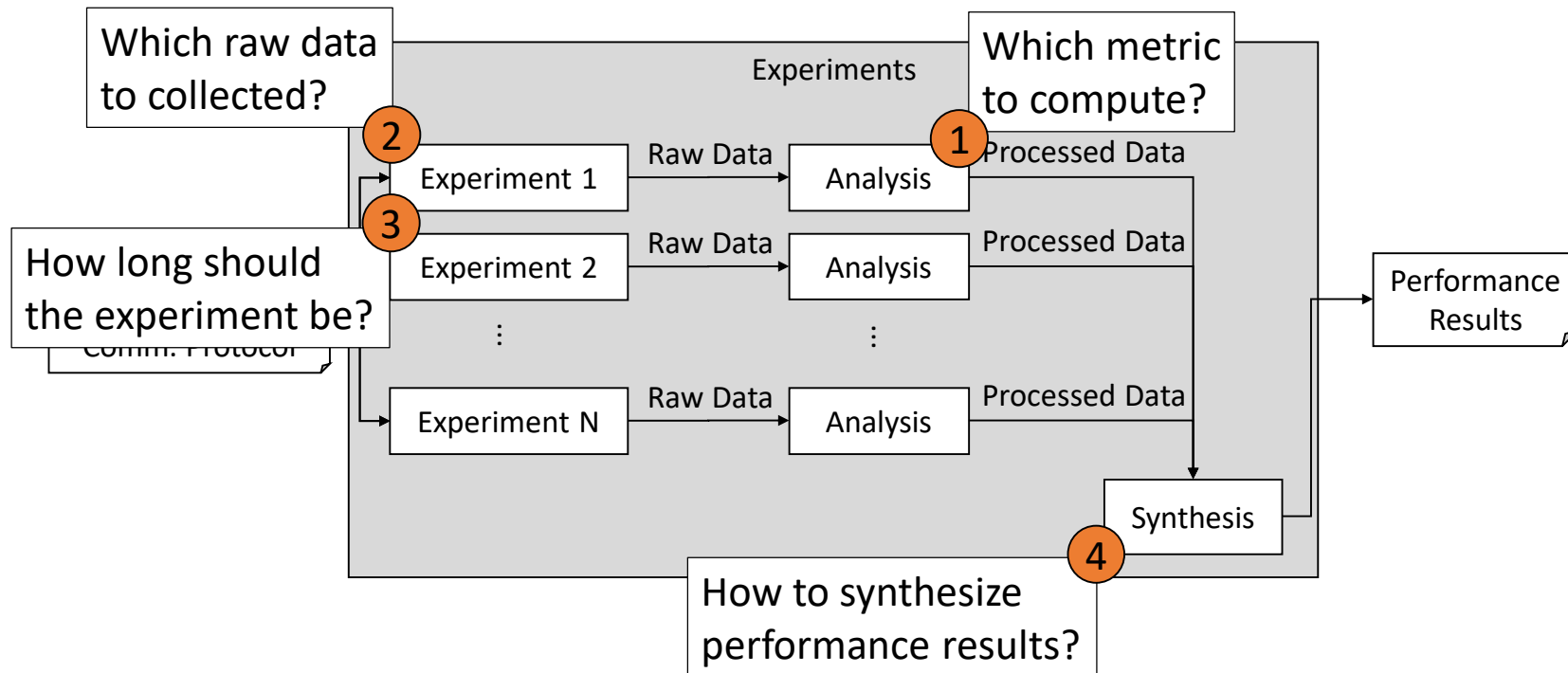


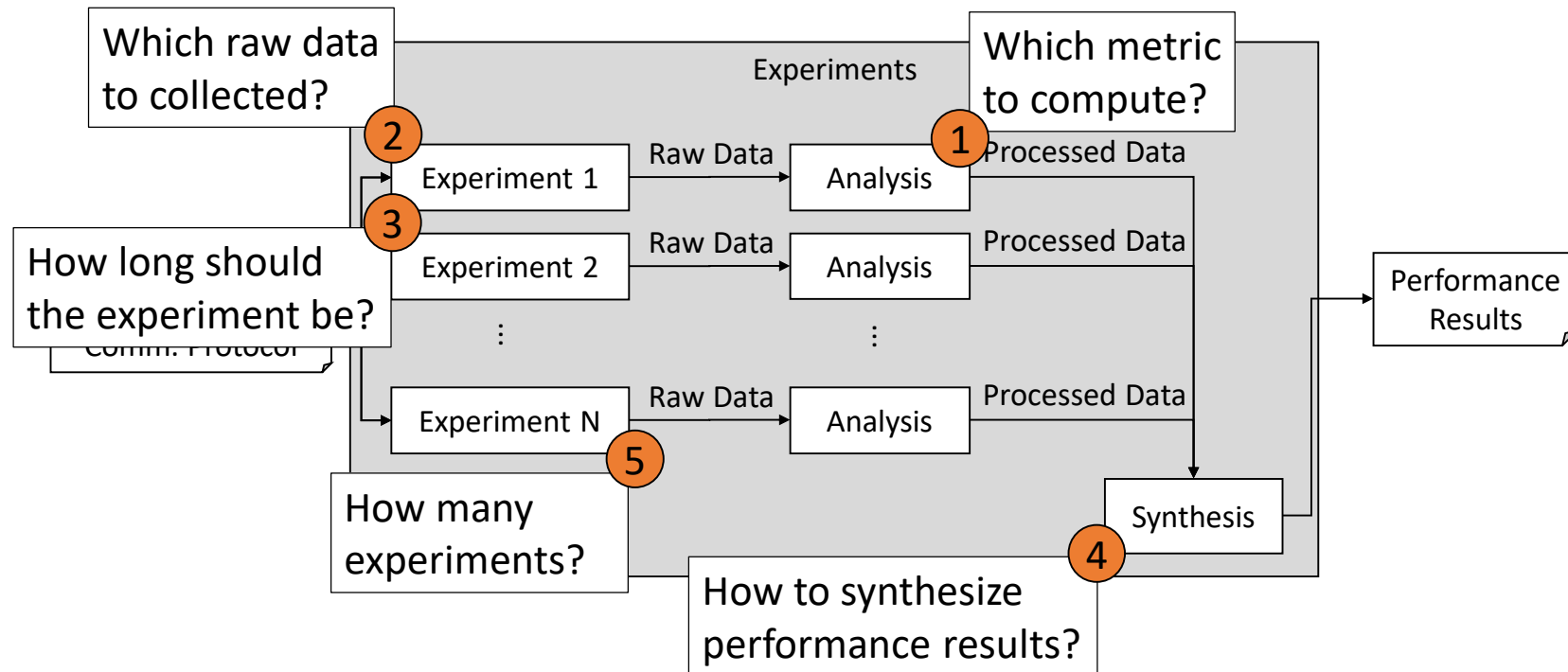








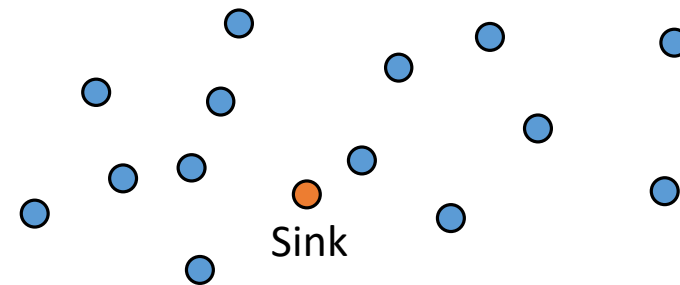




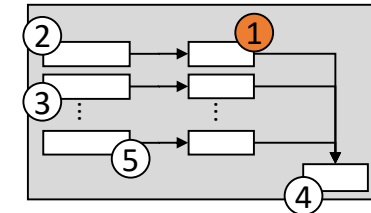
Case study – Periodic data collection

14 source nodes
200 payloads per source
2 Bytes per payload
10 payload per second
Periodic release, asynchronous

First payload released after 10s
Test stops 10s after last payload is released



Select the “metrics” based on the purpose of the evaluation



Performance
dimensions

Reliability
Latency
Energy efficiency

Average
Extremal

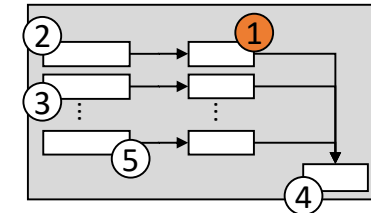
Metrics

Radio duty-cycle
Current draw

Measures

Mean
Median/percentiles
Max/Min

Select the “metrics” based on
the **purpose of the evaluation**



*How many application payloads can one expect to
successfully receive in one execution of the scenario?*

Dimension

Metric

Measure

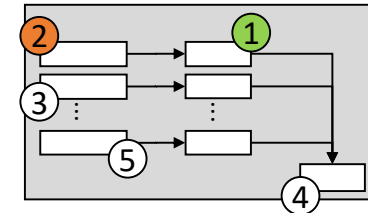
Average
reliability

PRR

~~Mean~~
Median

We are trying
to predict future
performance

Collect raw data
with the finest granularity possible



PRR

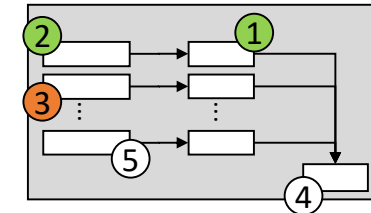
Log all received payloads at the sink

Current
draw

14400 samples/s
10 pA precision

1 every $\sim 7\mu s$

Define the length of the experiment
based on the scenario and the protocol



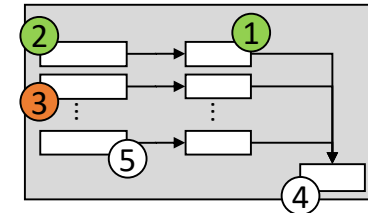
Generally
difficult

Correct approach **depends**
on the protocol under test

Easy case

If scenario is **terminating** and **short**,
then run it in full

Define the length of the experiment
based on the scenario and the protocol

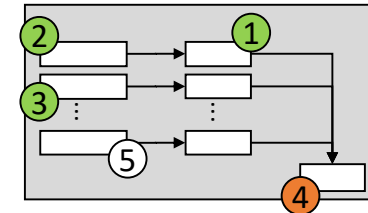


	Terminating	200 payload/source
+	Short	$10 + 200/10 + 10 = 40s$

⇒ Run in full

⇒ The uncertainty lies only in the
variability across experiments

Use performance indicators based on confidence intervals



Average
reliability

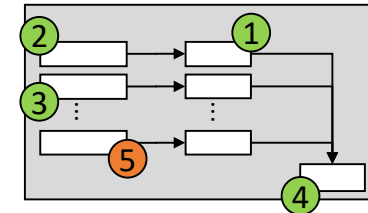
Overall PRR
 $\Rightarrow \forall \text{ experiment } j, x_j$

Received payloads
 $\frac{\quad}{200 * 14}$

95% CI on the median PRR for all exp.
 $\Rightarrow [x_m, x_{N-m+1}]$

Use conservative bound
 $\Rightarrow x_m$

Perform sufficiently many experiments
to obtain tight CI

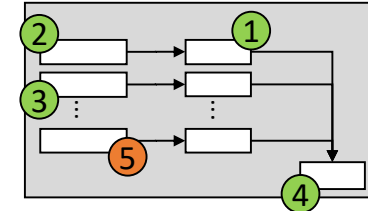


Intuition

Estimating **average** performance is easier
than **extremal** performance

Need more experiments to estimate a **95-th percentile**
than a **median**

Perform sufficiently many experiments
to obtain tight CI



Performance indicators based on 95% CI for the median

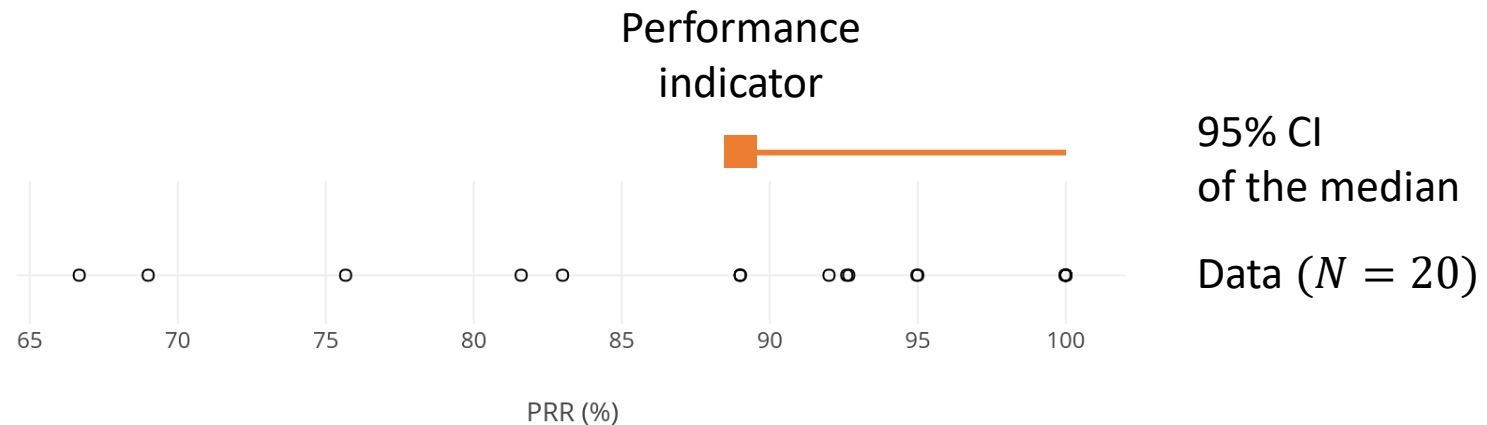
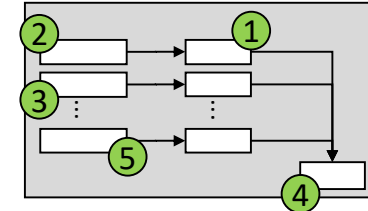
$$\Rightarrow N \geq 6$$

Aim for
tighter CI

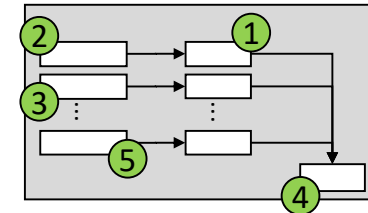
$$N = 20$$

$$\Rightarrow 95\% \text{ CI is } [x_6, x_{15}]$$

Following the methodology enables
unambiguous performance reports



Following the methodology enables
unambiguous performance reports



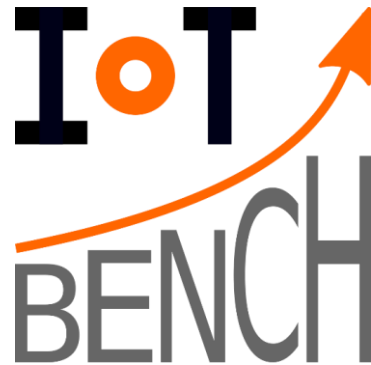
Protocol	A	B	C	D	E	F	G
Average Energy	0.82	0.83	0.89	0.86	0.90	0.43	0.25
Worst-case Energy	0.67	0.44	0.82	0.18	0.52	0.27	0.19
Reliability	0.40	0.41	0.89	0.06	0.48	0.27	0.25

We are on good way...

“ We need a benchmark for IoT networking. ”

- ⇔ Comparing performance
- ⇒ Repeatable experiments
- ⇒ Formalize the experimental methodology

Towards a Methodology for Experimental Evaluation in Low-Power Wireless Networking



Know your data

Use non-parametric statistics

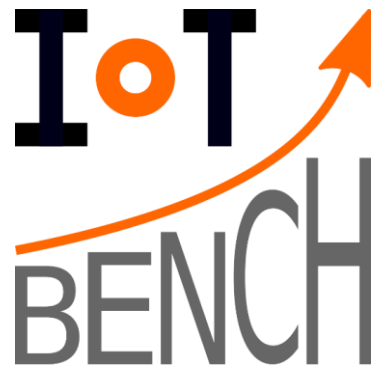
Formalizing low-power wireless
experimental evaluation

- Problem 1 Predictive statistics require **i.i.d. measurements**
Not a given. This must be checked, not assumed.
- Problem 2 What if the scenario is **not terminating**?
We still don't know how long one experiment should be.
- Problem 3 To be comparable, results must be **repeatable**
We still don't know how to formalize repeatability in our context.

These are **work-in-progress...**

72

Towards a Methodology for Experimental Evaluation in Low-Power Wireless Networking



Romain Jacob

Usman Raza

Lothar Thiele

Carlo Alberto Boano

Marco Zimmerling

ETH zürich

TOSHIBA



Includes material from Hanspeter Schmid and Alex Huber