

**École polytechnique de Louvain**

# **The $\Delta Q$ Oscilloscope: Real-Time Observation of Large Erlang Applications using $\Delta QSD$**

Author: **Francesco NIERI**

Supervisor: **Peter VAN ROY**

Readers: **Neil DAVIES, Tom BARBETTE, Peer STRITZINGER**

Academic year 2024–2025

Master [120] in Computer Science



# Abstract

It is difficult to study the detailed behaviour of large distributed systems while they are running. What happens when there is an overload? How can we feel something is wrong with the system before anything problematic can be observed? Current observability tools do not meet observability requirements when it comes to detecting problems early enough in running systems.

This thesis aim to provide further proof about how the  $\Delta$ QSD paradigm can be used to study the behaviour of running systems and to explore tradeoffs in system design, thanks to the implementation of the  **$\Delta$ Q oscilloscope**, a real time graphical dashboard that gives insights into a running Erlang system. The development of an Erlang wrapper, named `dqsd_otel`, allows the running system to communicate with the oscilloscope to receive real time insights about the execution of the former.

The oscilloscope performs statistical computations on the time series data it receives and displays the results in real time, thanks to the  $\Delta$ QSD paradigm. We provide a set of triggers which are set to capture rare events, like an oscilloscope would, and give a snapshot of the system under observation as if it was frozen in time. An implementation of a syntax to create outcome diagrams allows the creation of outcome diagrams which give an "observational view" of the system. Furthermore, the implementation of efficient algorithms allows for the computations to be done rapidly on precise representations of components.

We first provide an extensive summary of  $\Delta$ QSD concepts, which have been extended to allow the instrumentation of Erlang systems. Subsequently, we explain the user level concepts which are essential to understand how the oscilloscope works and understand what is displayed on the screen, delving later on into the mathematical foundations of the concepts. Lastly, we provide synthetic applications which prove the soundness of  $\Delta$ QSD and show how the oscilloscope is able to detect problems in a running system, diagnose it and explore design tradeoffs.

# Acknowledgments

This thesis is the culmination of my studies, I would like to thank the people who made this possible, those who supported me through the years and those who helped my with the thesis.

My family, especially **my mom, my dad, my brother** and **my sister**, for their help which was a crucial shoulder I could lay on while writing this thesis and most importantly throughout these five years.

My **friends**, to those who have taken time out of their lives to listen the thesis presentation and to those who through the years have been there for me.

**A-M.**, for the moments we shared these four years together in uni.

**My dad** and **Maurizio**, who nurtured the passion for coding in me.

**Peer Stritzinger, Stritzinger GmbH** and the **EEF Observability Working Group** (Bryan Naegele and Tristan Slougher) for their help in the EEF Slack, which helped me in understanding OpenTelemetry and gave me the `send_after` intuition.

**Neil Davies** for taking the time to proofread my thesis.

The **PNSol Ltd.** team for their extensive groundwork on  **$\Delta$ QSD** and its dissemination, which made this thesis possible.

Lastly, **Peter Van Roy**, for his year-long relentless interest, support and weekly and constant supervision which made sure the project would come to fruition.

# AI disclaimer

AI was used to aid the writing of the code. The **written master thesis** was written entirely without the aid of AI.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Context . . . . .	1
1.2	Objective . . . . .	1
1.3	Previous work . . . . .	2
1.4	Contributions . . . . .	2
1.5	Roadmap . . . . .	3
<b>2</b>	<b>Background</b>	<b>4</b>
2.1	An overview of $\Delta$ QSD . . . . .	4
2.1.1	Outcome . . . . .	4
2.1.2	Quality attenuation ( $\Delta$ Q) . . . . .	5
2.1.3	Failure semantics . . . . .	6
2.1.4	Partial ordering . . . . .	6
2.1.5	Timeliness . . . . .	6
2.1.6	QTA, required $\Delta$ Q . . . . .	6
2.1.7	Outcome diagram . . . . .	7
2.1.8	Independence hypothesis . . . . .	9
2.2	Observability . . . . .	10
2.2.1	erlang:trace . . . . .	10
2.2.2	OpenTelemetry . . . . .	11
2.3	Current observability problems . . . . .	13
2.3.1	Handling of long span . . . . .	14
<b>3</b>	<b>Design</b>	<b>15</b>
3.1	Probes . . . . .	15
3.2	Extending the notion of failure . . . . .	16
3.3	Time series - Oscilloscope outcome instances . . . . .	16
3.4	Erlang system . . . . .	17
3.4.1	System under test . . . . .	17
3.4.2	Wrapper/Adapter . . . . .	17
3.4.3	Inserting probes in Erlang - From spans to outcome instances . . . . .	18
3.5	Oscilloscope: C++ system . . . . .	18
3.5.1	Server . . . . .	18
3.5.2	Oscilloscope . . . . .	18
3.5.3	Inserting probes in the oscilloscope . . . . .	18

3.6	Oscilloscope system design diagram . . . . .	19
3.7	Triggers . . . . .	20
3.7.1	Snapshot . . . . .	20
3.8	Polling window . . . . .	20
<b>4</b>	<b>Oscilloscope: User level concepts</b>	<b>22</b>
4.1	$\Delta$ QSD implementation . . . . .	22
4.1.1	Internal representation of a $\Delta$ Q . . . . .	22
4.1.2	dMax . . . . .	23
4.1.3	QTA . . . . .	24
4.1.4	Confidence bounds . . . . .	24
4.2	$\Delta$ Q display . . . . .	24
4.3	Outcome diagram . . . . .	25
4.3.1	Causal link . . . . .	25
4.3.2	Sub-outcome diagrams . . . . .	25
4.3.3	Outcomes . . . . .	25
4.3.4	Operators . . . . .	25
4.3.5	Limitations . . . . .	26
4.4	Dashboard . . . . .	26
4.4.1	Sidebar . . . . .	26
4.4.2	Plots window . . . . .	28
4.5	Triggers . . . . .	28
4.5.1	Load . . . . .	28
4.5.2	QTA . . . . .	28
<b>5</b>	<b>Oscilloscope: implementation</b>	<b>29</b>
5.1	$\Delta$ QSD implementation . . . . .	29
5.1.1	dMax . . . . .	30
5.1.2	Operations . . . . .	30
5.1.3	Confidence bounds . . . . .	31
5.1.4	Rebinning . . . . .	31
5.2	Wrapper . . . . .	32
5.2.1	API . . . . .	32
5.2.2	Handling outcome instances . . . . .	34
5.2.3	TCP connection . . . . .	35
5.3	Parser . . . . .	36
5.3.1	ANTLR . . . . .	36
5.3.2	Grammar . . . . .	36
5.4	Oscilloscope GUI . . . . .	37
<b>6</b>	<b>Application on synthetic programs</b>	<b>39</b>
6.1	System with sequential composition . . . . .	39
6.1.1	System composition . . . . .	40
6.1.2	Determining parameters dynamically . . . . .	40
6.2	Detecting slower workers in workers . . . . .	45
6.2.1	First to finish application . . . . .	45
6.2.2	All to finish application . . . . .	47

<b>7</b>	<b>Performance study</b>	<b>49</b>
7.1	Convolution performance . . . . .	49
7.2	Wrapper performance . . . . .	50
7.3	GUI plotting performance . . . . .	50
<b>8</b>	<b>Conclusions and future work</b>	<b>52</b>
8.1	Future improvements . . . . .	52
8.1.1	Oscilloscope improvements . . . . .	53
8.1.2	Wrapper improvements . . . . .	53
8.1.3	Real applications . . . . .	54
8.1.4	Licensing limitations . . . . .	54



# Chapter 1

## Introduction

### 1.1 Context

$\Delta$ QSD is an industrial-strength approach for large-scale system design that can predict performance and feasibility early on in the design process. Developed over 30 years by a small group of people around Predictable Network Solutions Ltd, the paradigm has been applied in various industrial-scale problems with huge success and large savings in costs. [1] Moreover, it is the basis of Broadband forum's TR452 standard series, used in instrumenting data networks. [2]

Modern software development practices successfully fail to adequately consider essential quality requirements or even to consider properly whether a system can actually meet its intended outcomes, particularly when deployed at scale, the  $\Delta$ QSD paradigm addresses this problem! [3]

$\Delta$ QSD has important properties which make its application to distributed projects interesting, it supports:

- A compositional approach that considers performance and failure as first-class citizens.
- Stochastic approach to capture uncertainty throughout the design approach.
- Performance and feasibility can be predicted at high system load for partially defined systems.

While the paradigm has been successfully applied in **a posteriori** analysis, there is no way yet to analyse a distributed system which is running in real time with  $\Delta$ QSD! This is where the  $\Delta$ Q oscilloscope comes in.

### 1.2 Objective

This project will develop a practical tool, the  **$\Delta$ Q oscilloscope**, for the Erlang developer community.

The Erlang language and Erlang/OTP platform are widely used to develop distributed applications that must perform reliably under high load [4]. The tool will provide useful information for these applications both for understanding their behaviour, for diagnosing performance issues, and for optimizing performance over their lifetime.

The  $\Delta Q$  Oscilloscope will perform statistical computations to show real time graphs about the performance of system components. With the oscilloscope prototype we will present in this paper, we are aiming to show that the  $\Delta QSD$  paradigm is not only a theoretical paradigm, but it can be employed in a tool to diagnose large distributed systems.

The oscilloscope targets large distributed applications handling many independent tasks where performance and reliability are important.

### 1.3 Previous work

The  $\Delta QSD$  paradigm has been formalised across different papers [3] [5] and was brought to the attention of engineers via tutorials [1] and to students at Université Catholique de Louvain [6].

A Jupyter notebook workbench has been made available on GitHub [7], it shows real time  $\Delta Q$  graphs for typical outcome diagrams but is not adequate to be scaled to real time systems, it is meant as an interactive tool to show how the  $\Delta QSD$  paradigm can be applied to real life examples.

Observability tools such as Erlang tracing [8] and OpenTelemetry [9] lack the notions of failure as defined in  $\Delta QSD$ , which allows detecting performance problems early on, we base our program on OpenTelemetry to incorporate already existing notions of causality and observability to augment their capabilities and make them suitable to work with the  $\Delta QSD$  paradigm.

### 1.4 Contributions

There are a few contributions that make the master thesis and thus, the oscilloscope, possible:

- A graphical interface to display  $\Delta Q$  plots for outcomes.
- An Erlang OpenTelemetry wrapper to give OpenTelemetry spans a notion of failure and to communicate with the oscilloscope.
- An implementation of a syntax, derived from the original algebraic syntax to create outcome diagrams.
- The implementation of  $\Delta QSD$  concepts from theory to practice, allowing probes'  $\Delta Qs$  to be displayed and analysed on the oscilloscope.
- An efficient convolution algorithm based on the FFTW3 library.

- A system of triggers to catch rare events when system behaviour fails to meet quality requirements, giving a snapshot of the system, giving the user insights about their system's behaviour.
- Synthetic applications to test the effectiveness of  $\Delta$ QSD on diagnosing systems and their feasibility.

These contributions can show that the  $\Delta$ QSD has its practical applications and is not limited to a theoretical view of system design.

## 1.5 Roadmap

The following thesis will give the reader everything that is needed to use the oscilloscope and exploit it to its full potential.

We divided the thesis in multiple chapters, below is the roadmap of the content:

- The background chapter gives the reader an extensive background into the theoretical foundations of  $\Delta$ QSD, which are the basis of the oscilloscope and are fundamental to understand how to correctly use and analyse the output given by the oscilloscope. Secondly, an introduction to OpenTelemetry, the library we base our Erlang wrapper on, and the problems that are present in the observability library.
- The design chapter initially extends the concepts of  $\Delta$ QSD, introducing novel aspects which are key to understand how to properly instrument the different part of the systems, as they are the basis upon which we build the oscilloscope. We then delve how the parts of the system interact together and how to correctly apply the concepts we just introduced in the oscilloscope and the Erlang system.
- We then present the oscilloscope in two different chapters, first providing "user level concepts" of how  $\Delta$ QSD is used in the oscilloscope and what the user should expect graphically from the oscilloscope. Secondly, a more low level explanation, which goes into more technical details of the parts that compose the oscilloscope and the mathematical explanations of  $\Delta$ QSD concepts explained in the previous chapter.
- We then provide synthetic applications which have been tested with the oscilloscope that demonstrate the usefulness of the oscilloscope in a distributed setting. We also perform evaluations of the performance of the different parts we have developed to understand the overhead that are present.

We end by providing future possibilities which can be explored, and concepts which we believe ought to be implemented in observabilities tools. In the appendix, we provide a user manual to help users use the oscilloscope, along with C++ and Erlang source code of the oscilloscope and the wrapper.

The oscilloscope (<https://github.com/fnieri/DeltaQOscilloscope>) and wrapper([https://github.com/fnieri/dqsd\\_otel](https://github.com/fnieri/dqsd_otel)) can be found on GitHub as open source projects.

# Chapter 2

## Background

This chapter aims to provide firstly a complete background of the concepts key to understanding the  $\Delta$ QSD.

Secondly, we provide a comprehensive background into the observability solutions that have been explored for the oscilloscope, delving deeper into OpenTelemetry and its macros.

We finish by explaining the current limitations of OpenTelemetry and explaining where our oscilloscope comes in.

### 2.1 An overview of $\Delta$ QSD

$\Delta$ QSD is a metrics-based, quality-centric paradigm that uses formalised outcome diagrams to explore the performance consequences of design decisions. [5]

Key concepts of  $\Delta$ QSD are **quality attenuation ( $\Delta$ Q)** and **outcome diagram**

Outcome diagrams capture dependency and causality properties of the system. The  $\Delta$ QSD paradigm derives bounds on performance expressed as probability distribution, encompassing all possible executions of the system.

The following sections are a summary of multiple articles and presentation formalizing the paradigm. [1] [5] [3] [2]

#### 2.1.1 Outcome

An outcome  $O$  is a specific system behaviour that can be observed to start at some point in time and *may* be observed to complete at some later time. [2] Formally, what the system obtains by performing one of its tasks. One task corresponds to one outcome and viceversa. When an outcome is performed, it means that the task of an outcome is performed.

**Observables** Each outcome has two starting sets of events: the starting sets and the ending sets. Such sets are called the *observables*. Once an event from the starting set

occurs, there is no guarantee that a corresponding event in the terminating set will occur within the duration limit (required time to complete). An observable is *done* when it occurs during the time limit. [3]

**Outcome instance** An outcome instance is the result of an execution of an outcome given a starting event  $e_{in}$  and an end event  $e_{out}$ .

**Graphical Representation** Outcomes are represented as circles, with the starting and terminating set of events being represented by boxes.

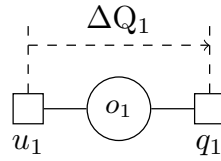


Figure 2.1: The outcome (circle) and the starting set (left) and terminating set (right) of events.

### 2.1.2 Quality attenuation ( $\Delta Q$ )

Assume a component  $C$  which receives a message  $m_{in}$  and outputs a message  $m_{out}$  after a delay  $d$ . Over multiple executions, we will have observed multiple delays which can be represented as a cumulative definition where  $p$  percent of delays have delay  $\leq d$ . [3]

$\Delta Q$  is a cumulative distribution function that defines both *latency* and *failure probability* between a start and end event [1]

In an ideal system, an outcome would deliver a desired behaviour without error, failure, delay, but this is not the case. The quality of an outcome response "attenuated to the relative ideal" (the cumulative distribution function) is called "quality attenuation" ( $\Delta Q$ ) and can depend on many factors (geographical, physical ...). Its distribution may be modelled by a random variable.

As  $\Delta Q$  captures deviation from ideal behavior and incorporates delay, which is a continuous random variable, and failures/timeouts, which are discrete variables, it can be described mathematically as an *Improper Random Variable*, where the probability of a finite or bounded delay  $< 1$ . Combining latency and failure together makes it easy to examine the tradeoffs between them.

$\Delta Q(x)$  is the probability that an outcome  $O$  occurs in time  $t \leq x$ . The *intangible mass*  $1 - \lim_{x \rightarrow \infty} \Delta Q(x)$  of a  $\Delta Q$  will encode the probability of failure/timeout/exception occurring. [5]

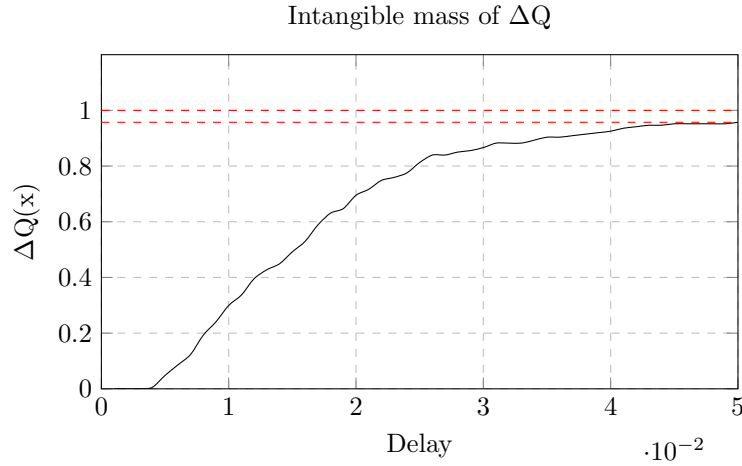


Figure 2.2: Intangible mass (red) of a  $\Delta Q$ , the observable had a failure rate of about 5%

### 2.1.3 Failure semantics

In the CDF representation of a  $\Delta Q$ , there is an  $f$  percent probability that the delay is infinite, this is what failure models. Concretely, it means that an input message  $m_{in}$  **has no output message**  $m_{out}$ . [3]

Combining delay and failure in a single quantity is what makes  $\Delta QSD$  a great choice to explore feasibility in system design. [1]

### 2.1.4 Partial ordering

A CDF of a  $\Delta Q$  is *less than* the other if its CDF is everywhere to the left and above the other. Mathematically, it is a partial order.

If two  $\Delta Q$ s intersect, they are not ordered. [1]

### 2.1.5 Timeliness

Timeliness is defined as a relation between an observed  $\Delta Q_{obs}$  and a required  $\Delta Q_{req}$ . Timeliness is delivering results within required time bounds (sufficiently often).

A system *satisfies timeliness* if  $\Delta Q_{obs} \leq \Delta Q_{req}$ . [3]

### 2.1.6 QTA, required $\Delta Q$

The *Quantitative Timeliness Agreement* (QTA) maps objective measurements to the subjective perception of application performance [2]. It specifies what the base system does and its limits.

**Slack** There is performance *slack* when a  $\Delta Q$  is strictly less than the requirement,.

**Hazard** There is performance *hazard* when a  $\Delta Q$  is strictly greater than the requirement.

**QTA example** : Imagine a system where 25% of the executions should take  $< 15$  ms, 50%  $< 25$  ms and 75%  $< 35$  ms, all queries have a maximum delay of 50ms and 5% of executions can timeout, the QTA can be represented as a step function.

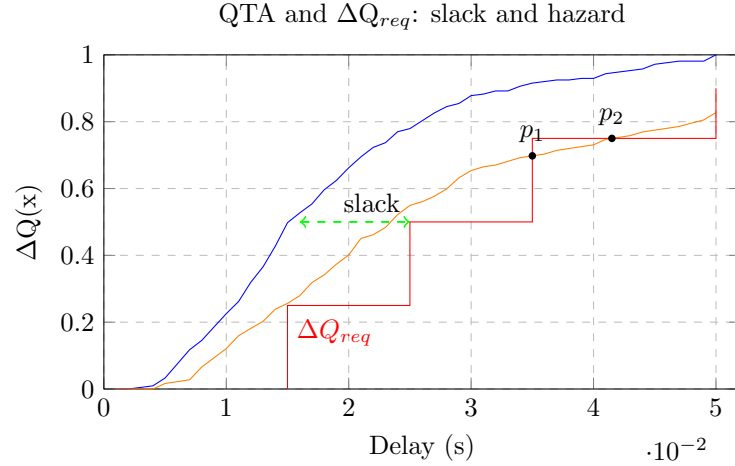


Figure 2.3: The system in blue is showing slack and satisfies the requirement, the system in orange is showing signs that it cannot handle the stress, it is not respecting the  $\Delta Q_{req}$ .

### 2.1.7 Outcome diagram

An outcome diagram is central to capture the causal relationships between the outcomes. It shows the causal connections between all the outcomes we are interested in, and it allows computing the  $\Delta Q$  for the whole system [1]. The outcome diagram captures the essential observational properties of a system. It maps a system's behaviour as seen from outside to concrete outcomes. [3]

There are four different operators that represent the relationships between outcomes. [1]

#### Sequential composition

If we assume two outcomes  $O_A$ ,  $O_B$  where end event of  $O_A$  is the start event of  $O_B$ , the two outcomes can be sequentially composed. The  $\Delta Q$  of  $O_{AB}$  is given by the convolution of the PDFs of  $O_A$  and  $O_B$ .

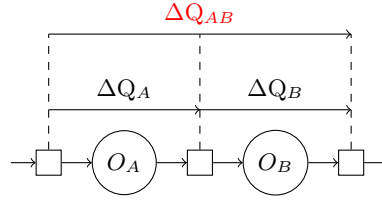


Figure 2.4: Sequential composition of  $O_A$  and  $O_B$ .

Where convolution ( $\circledast$ ) between two PDF is:

$$PDF_{AB}(t) = \int_0^t PDF_A(\delta) \cdot PDF_B(t - \delta) d\delta \quad (2.1)$$

and thus  $\Delta Q_{AB}$ :

$$\Delta Q_{AB} = \Delta Q_A \circledast \Delta Q_B \quad (2.2)$$

### First to finish (FTF)

If we assume two independent outcomes  $O_A, O_B$  with the same start event, first-to-finish occurs when at least one end event occurs, it can be calculated as:

$$\begin{aligned} \Delta Q_{FTF(A,B)} &= Pr[d_A > t \wedge d_B > t] \\ &= Pr[d_A > t] \cdot Pr[d_B > t] = (1 - \Delta Q_A) \cdot (1 - \Delta Q_B) \\ \Delta Q_{FTF(A,B)} &= \Delta Q_A + \Delta Q_B - \Delta Q_A \cdot \Delta Q_B \end{aligned} \quad (2.3)$$

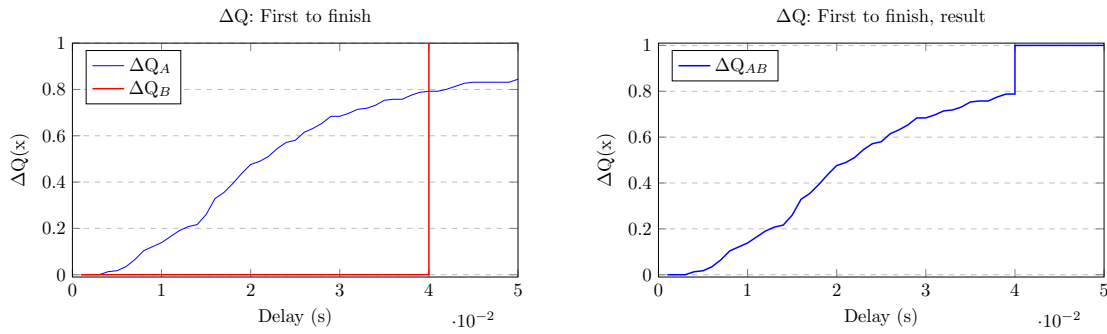


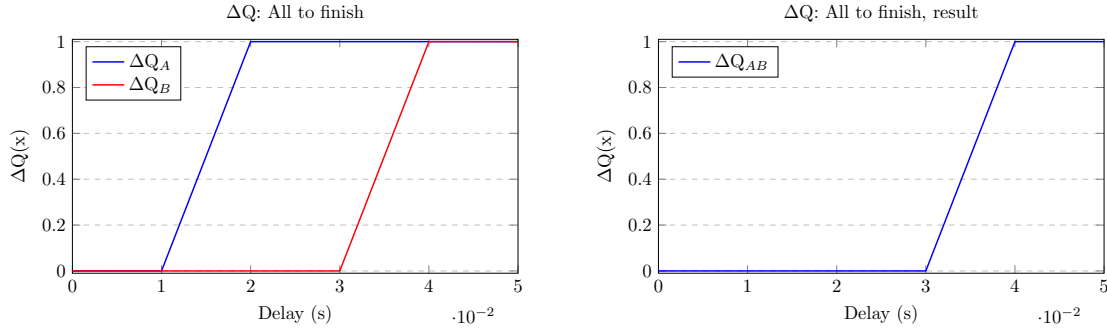
Figure 2.5: Left:  $\Delta Q_{(A,B)}$ . Right:  $FTF_{(A,B)} = \Delta Q_{AB}$

### All to finish (ATF)

If we assume two independent outcomes  $O_A, O_B$  with the same start event, all-to-finish occurs when both end events occur, it can be calculated as:

$$\begin{aligned} \Delta Q_{ATF(A,B)} &= Pr[d_A \leq t \wedge d_B \leq t] \\ &= Pr[d_A \leq t] \cdot Pr[d_B \leq t] = \Delta Q_A \cdot \Delta Q_B \\ \Delta Q_{ATF(A,B)} &= \Delta Q_A \cdot \Delta Q_B \end{aligned} \quad (2.4)$$




 Figure 2.6: Left:  $\Delta Q_{(A,B)}$ . Right:  $\text{ATF}_{(A,B)} = \Delta Q_{AB}$ 

### Probabilistic choice (PC)

If we assume two possible outcomes  $O_A$  and  $O_B$  and exactly one outcome is chosen during each occurrence of a start event and:

- $O_A$  happens with probability  $\frac{p}{p+q}$
- $O_B$  happens with probability  $\frac{q}{p+q}$

$$\Delta Q_{PC}(A, B) = \frac{p}{p+q} \Delta Q_A + \frac{q}{p+q} \Delta Q_B \quad (2.5)$$

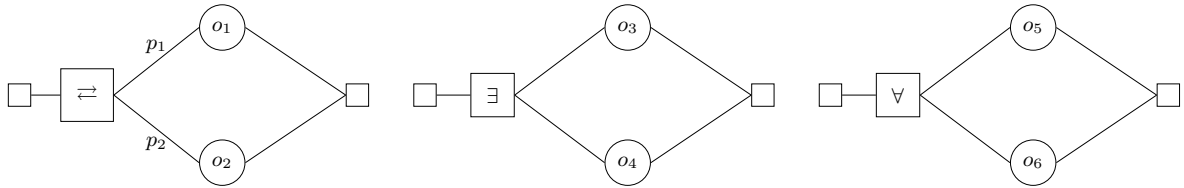


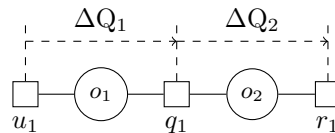
Figure 2.7: The possible operators in an outcome diagram: Probabilistic choice, first-to-finish, all-to-finish

First-to-finish, All-to-finish and probabilistic-choice are calculated on the CDF of the  $\Delta Q$  of their components.

These operators can be assembled together to create an outcome diagram, later on, we will see how one can go from the graphical representation to outcome diagrams which can be used in the  $\Delta Q$  oscilloscope.

#### 2.1.8 Independence hypothesis

Assume two sequentially composed outcomes  $o_1, o_2$  running on the same processor. We observe the delay of execution from the start event of  $o_1$  to the end event of  $o_2$ .



At low load, the two components behavior will be independent, the system will behave linearly, the observed total delay of execution will be equal to the result of convolution of  $o_1, o_2$  ( $o_1 \otimes o_2$ ).

When load increases, the two components will start to show dependent behaviour due to the processor utilisation increasing. The  $\Delta Q$  of the observed delay will then deviate from the  $\Delta Q$  which is the result of the convolution of  $o_1, o_2$

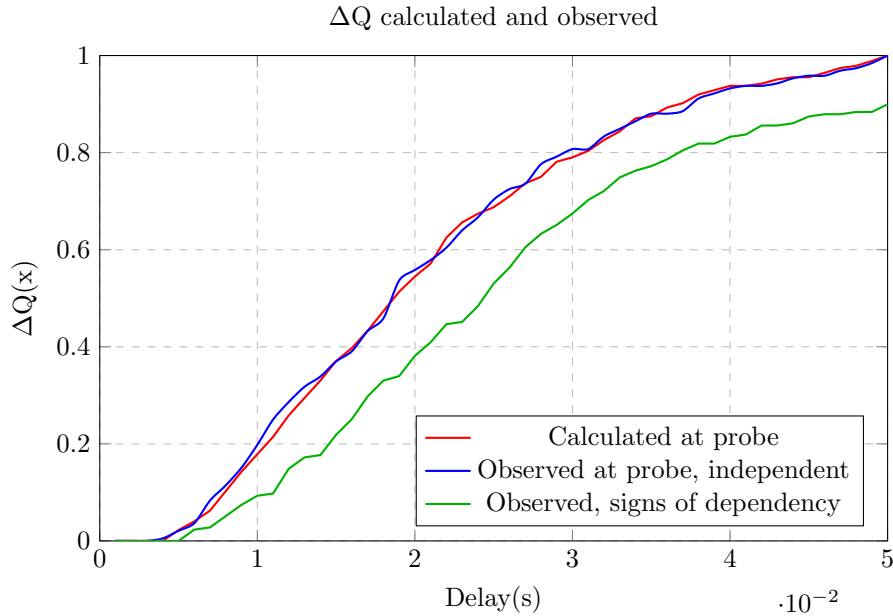


Figure 2.8: When the components are independent, what is observed (blue) and calculated (red) can be superposed, whilst when  $o_1$  and  $o_2$  show initial signs of dependency, what is observed (green) can be seen deviating from the  $\Delta Q_{o_1, o_2}$ .

When the system is far from being overloaded, the effect is noticeable thanks to  $\Delta QSD$  even if the system is far from being overloaded. As the cliff edge of overload is approached, the nonlinearity will increase.

## 2.2 Observability

Observability refers to the ability to understand the internal state by examining its output, in the context of a distributed system, being able to understand the internal state of the system by examining its telemetry data. [10]

In the case of the Erlang programming language, we explain below two tools that can be used to observe an Erlang program.

### 2.2.1 erlang:trace

The Erlang programming language gives the users different ways to observe the behaviour of a system, one of those is the function `erlang:trace/3`. The erlang run-time system

exposes several trace points that can be observed, observing the trace points allows users to be notified when they are triggered [8]. One can observe function calls, messages being sent and received, process being spawned, garbage collecting ....

```
-spec trace(PidPortSpec, How, FlagList) -> integer()
when
    PidPortSpec ::
        pid() |
        port() |
        all | processes | ports | existing | existing_processes |
        ↪ existing_ports | new |
        new_processes | new_ports,
    How :: boolean(),
    FlagList :: [trace_flag()].
```

Figure 2.9: erlang:trace/3 specification

Nevertheless, Erlang Tracing, according to our use case, has a major flaw: no notion of causality. If two messages  $a, b$  are sent and then received in disorder, the tracer has no default way of knowing which is which, this is a missing feature that is crucial for observing a program functioning and being able to connect an application to our oscilloscope. This is where the OpenTelemetry standard comes in.

## 2.2.2 OpenTelemetry

OpenTelemetry is an open-source, vendor agnostic observability framework and toolkit designed to generate, export and collect telemetry data, in particular traces, metrics and logs[10]. OpenTelemetry provides a standard protocol, a single set of API and conventions and lets you own the generated data, allowing to switch between observability backends freely.

OpenTelemetry is available for a plethora of languages, including Erlang, although, as of writing this, only traces are available in Erlang.

The Erlang Ecosystem Foundation has a working group focused on evolving the tools related to observability.

### Traces

Traces are why we are basing our program on top of OpenTelemetry, traces follow the whole "path" of a request in an application, traces are comprised of one or more spans.

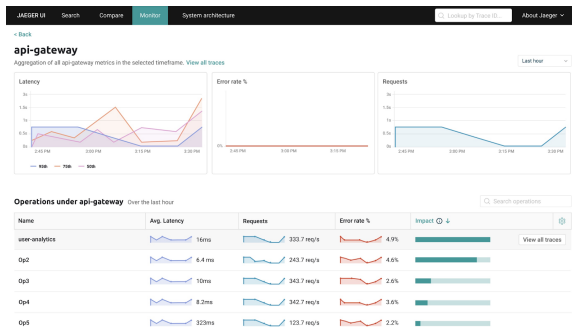
**Span** A span is a unit of work or operation. Spans can be correlated to each other and can be assembled into a trace. The notion of spans and traces allows us to follow the execution of a "request" and carry a context, allowing us to get the causal links of messages. [11]

```
{
  "name": "oscilloscope-span",
  "context": {
    "trace_id": "5b8aa5a2d2c872e8321cf37308d69df2",
    "span_id": "5fb397be34d26b51"
  },
  "parent_id": "0515505510cb55c13",
  "start_time": "2022-04-29T18:52:58.114304Z",
  "end_time": "2022-04-29T22:52:58.114561Z",
  "attributes": {
    "http.route": "some_route"
  },
}
```

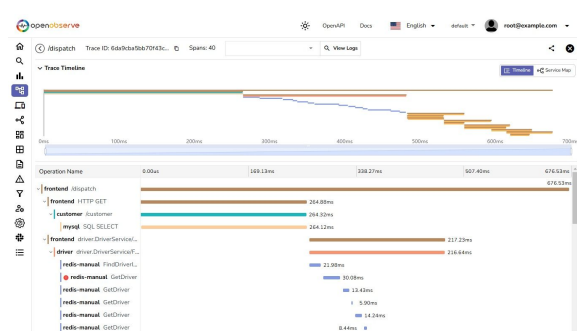
Figure 2.10: Example of span with a parent, indicating a causal link between parent and children span [11]

## Monitoring OpenTelemetry spans

OpenTelemetry gives the possibility to export traces to backends such as Jaeger, Zipkin, Datadog. A user can monitor their workflows, analyse dependencies, troubleshoot their programs by observing the flow of the requests in such backends[12]. These monitoring tools give extensive details about a running system, but may fail to capture essential requirements early enough.



(a) Jaeger interface [13].



(b) A span analysis on OpenObserve [14]

## Macros

OpenTelemetry provides macros to start, end and interact with spans in Erlang, the following code excerpts are taken from the instrumentation wiki. [15]

**?with\_span** `?with_span` creates active spans. An active span is the span that is currently set in the execution context and is considered the "current" span for the ongoing operation or thread. [16]

```
parent_function() ->
  ?with_span(parent, #{}, fun child_function/0).
```

```
child_function() ->
  %% this is the same process, so the span parent set as the active
  %% span in the with_span call above will be the active span in
  ↪ this function
  ?with_span(child, #{},
    fun() ->
      %% do work here. when this function returns, child
      ↪ will complete.
    end).
```

**?start\_span** ?start\_span creates a span which isn't connected to a particular process, it does not set the span as the current active span.

```
SpanCtx = ?start_span(child),
Ctx = otel_ctx:get_current(),
proc_lib:spawn_link(fun() ->
  otel_ctx:attach(Ctx),
  ?set_current_span(SpanCtx),
  %% do work here
  ?end_span(SpanCtx)
end),
```

**?end\_span** ?end\_span ends a span started with ?start\_span

## 2.3 Current observability problems

A legitimate question to pose would be why one would need an additional tool to observe their system, monitoring tools are already plenty and provide useful insights into an application's behaviour. While they may seem adequate to provide a global oversight of applications, they fail to diagnose real time problems like overload, dependent behaviour early enough and in a quick manner.

The problem we are trying to tackle can be described by the following situation: Imagine an Erlang application instrumented with OpenTelemetry, suddenly, the application starts slowing down, and the execution of a function takes 10 seconds instead of the usual 1 second. Between its start and its end, the user instrumenting the application sees nothing in their dashboard.

This is a big problem! One would like to know right away if something is wrong with their application, better! Even before problems are apparent. This is where the  $\Delta$ QSD paradigm and the  $\Delta$ Q oscilloscope come in handy.

By leveraging  $\Delta$ QSD notion of failure and QTAs, problems can be detected right away.

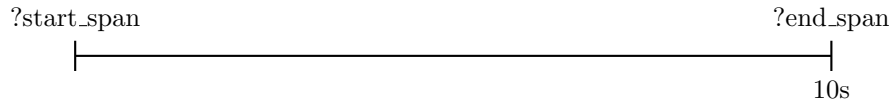


Figure 2.12: Execution of a long span in OpenTelemetry, the user will only be notified after 10 seconds that the function has ended (and taken too long).

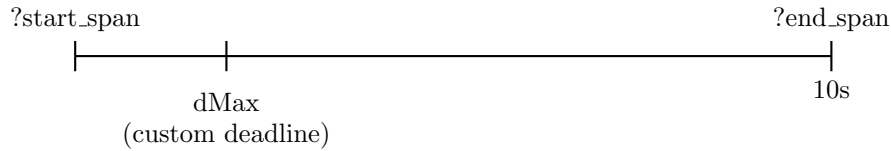


Figure 2.13: Execution of a long span in OpenTelemetry, the *dMax* deadline allows knowing that the span has taken too long.

### 2.3.1 Handling of long span

OpenTelemetry presents a bigger problem, what happens when there are long-running spans? Worse, what happens when spans are not actually terminated?

OpenTelemetry limits the length of its spans, moreover, those who are not terminated are lost and not exported. Why? Failed executions are those that tell more about a program's execution!

If the span is the parent/root span, its effect could trickle down to child spans. We can quickly see how this becomes problematic, all the information about an execution of your program ...lost. Moreover, a span could not be terminated for trivial reasons: refreshing a tab, network failures, crashes ... [17]. There are a few hacks that can be implemented, having shorter spans, carrying data in child spans, saving spans in a log to track spans which were not ended to manually set an end time; why the need to circumvent limitations when observing a system?

We believe that the wrapper we provide can be a great start to improve observability requirements surrounding OpenTelemetry. We will show in the evaluation on synthetic applications how  $\Delta$ QSD's notion of failure can help to detect overload problems in running systems right away. The oscilloscope can be an addition to be put next to monitoring platforms to give further statistical insights about a running system with high precision.

# Chapter 3

## Design

This chapter aims to first extend the concepts of  $\Delta$ QSD, giving more insights into how the systems need to be instrumented to correctly work together, and how the different parts need to be integrated to interact together.

- We first provide concepts of probes, we extend the  $\Delta$ QSD notion of failure and describe how time series will work in our oscilloscope, this part is crucial to understand how the different parts of the system work together.
- We then split the design of the oscilloscope in two. First explaining the Erlang side, where the system to be tested is. Secondly, we explain the C++ side. Both chapters explain how probes can be inserted and made to work together. We conclude the sections by showing the system diagram of the different parts.
- Lastly, we provide high level concepts of the key elements of the oscilloscope.

### 3.1 Probes

To observe a system, we must put probes in it. For each outcome of interest, a probe (observation point) is attached to measure the delay of the outcome, like one would in a true oscilloscope.

Consider the figure below, a probe is attached at every component to measure their  $\Delta$ Qs ( $c_2, c_3$ ), Another probe ( $p_1$ ) is inserted at the beginning and end of the system to measure the global execution delay. Thanks to this probe, the user can observe the  $\Delta$ Q "*observed at  $p_1$* ", which is the  $\Delta$ Q which was calculated from the data received by inserting probe  $p_1$ . The  $\Delta$ Q "*calculated at  $p_1$* " is the resulting  $\Delta$ Q from the convolution of the observed  $\Delta$ Qs at  $c_2$  and  $c_3$ .

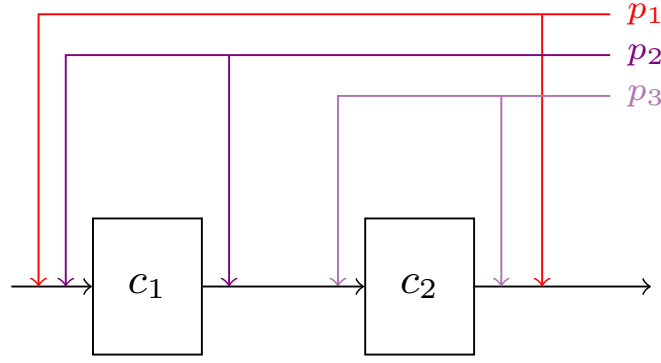


Figure 3.1: Probes inserted in a component diagram. In an applications instrumented with OpenTelemetry,  $p_1$  could be considered the root span,  $c_1$  and  $c_2$  its children spans sharing a causal link.

## 3.2 Extending the notion of failure

Recall the definition of failure: *"an input message  $m_{in}$  that has no output message  $m_{out}$ ".* If you recall the previous section ??, we introduced the notion of a maximum delay.

We extend the notion of failure to the following definition:

*"An input message  $m_{in}$  that has no output message  $m_{out}$  after  $x$  seconds"*

Where  $x$  is the  $dMax$  defined by the user. We can leverage this new definition to observe  $\Delta Qs$  in real time.

## 3.3 Time series - Oscilloscope outcome instances

Consider a probe  $p$  with two distinct sets of events, the starting set of events  $s$  and ending set of event  $e$ . The outcome instance of a message  $m_s \rightarrow m_e$ :

- The probe's  $p$  name
- The start time  $t_s$
- The end time  $t_e$
- Its status
- Its elapsed time of execution

The instance has three possible statuses: **success**, **timeout**, **failure**, it can thus be broken down in the representations, based on its status:

- $(t_s, t_e)$ : This representation indicates that the execution was successful ( $t < dMax$ ).
- $(t_s, \mathcal{T})$ : This representation indicates that the execution has timed out ( $t > dMax$ ). The end time and elapsed time is equal to  $t_s + \text{timeout}$



- $(t_s, \mathcal{F})$ : This representation indicates the execution has failed given a user defined requirement (i.e. a dropped message given buffer overload in a queue system). It must not be confused with a program failure (crash), if a program crashes during the execution of event  $e$ , it will time out since the wrapper will not receive an end message.

The **time series** of a probe is the sequence of  $n$  outcome instances and can then be easily modeled by  $\Delta Q$ .

**What can be considered a failed execution?** Imagine a queue with a buffer: the buffer queue being full and dropping incoming messages can be modeled as a failure.

More generally, the choice of what is considered a failed execution is left up to the user who is handling the spans and is program-dependent. Exceptions or errors can be kinds of failure.

On another note, the way of handling errored spans in OpenTelemetry can differ from user to user, so the wrapper will not handle ending and setting statuses for "failed" spans.

## 3.4 Erlang system

### 3.4.1 System under test

The system under test (**S.U.T**) is the Erlang system the engineer wishes to observe, it ideally is a system which already is instrumented with OpenTelemetry. The ideal system where  $\Delta QSD$  is more useful is a system that executes many independent instances of the same action.

### 3.4.2 Wrapper/Adapter

The adapter is the `dqsd_otel` Erlang application, a wrapper that starts and ends OpenTelemetry spans and translates them to outcome instances which are useful for the oscilloscope. This can be done thanks to probes being attached to the system under test, like an oscilloscope would! The outcome instances end normally like OpenTelemetry spans or, additionally, can timeout, given a custom timeout, and fail, according to user's definition of failure.

Handling of OpenTelemetry spans which goes beyond starting and ending them is delegated to the user, who may wish to do further operations with their spans. The wrapper is called from the system under test and communicates outcome instances data to the oscilloscope via TCP.

The wrapper can receive messages from the oscilloscope, the messages are about updating observable's *dMax* or starting and stopping the sending of data to the oscilloscope.

### 3.4.3 Inserting probes in Erlang - From spans to outcome instances

OpenTelemetry spans are useful to carry context, attributes and baggage in a program. The plethora of attributes they have is nevertheless too much for the oscilloscope.

To get the equivalent of spans for the oscilloscope, the wrapper needs to be called at the starting events of a probe to start an instance of a probe, and at the ending events to end the outcome instance and send the data to the oscilloscope. The name given with "start\_span" is the name of the probe.

```
% Start the outcome instance of worker_2
{WorkerCtx, WorkerPid} = dqsd_otel:start_span(<<"worker_2">>),

% Do work here ...

%End the outcome instance of worker_2
dqsd_otel:end_span(WorkerCtx, WorkerPid),
```

## 3.5 Oscilloscope: C++ system

### 3.5.1 Server

The server is responsible for receiving the messages containing the outcome instances from the wrapper. The server forwards the instances to the oscilloscope.

### 3.5.2 Oscilloscope

The oscilloscope is a C++ graphical application which implements a dashboard to observe  $\Delta Q$ s of probes inserted in the system. It receives the instances corresponding to probes from the server and adds them to the time series of the probes whose instance is being received. The oscilloscope has a graphical interface which allows the user to create an outcome diagram of the system under test, display real time graphs which show detail about the execution of the system and allow the user to set custom timeouts for probes. It can also display snapshots of the system as if it was frozen in time

### 3.5.3 Inserting probes in the oscilloscope

Probes are automatically inserted in the oscilloscope when creating an outcome diagram. They are inserted on the outcomes, operators and to the causal result of operations, we will see later on how they can be defined and how an outcome diagram can be created.

In the system below, which is equal to the one defined above, probes are automatically attached to outcomes  $o_1, o_2$ . The user who wants to observe the result of the sequential composition can insert probes at the start and end of the routine.

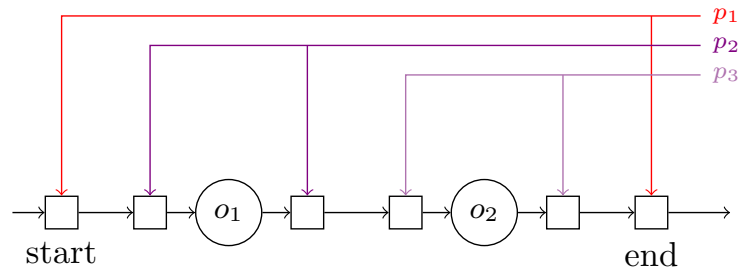


Figure 3.2: Probes inserted in the outcome diagram of the previous component diagram.  
Figure 3.1

As for operators, probes are automatically attached to the components inside them and to the start event and end events of the operators.

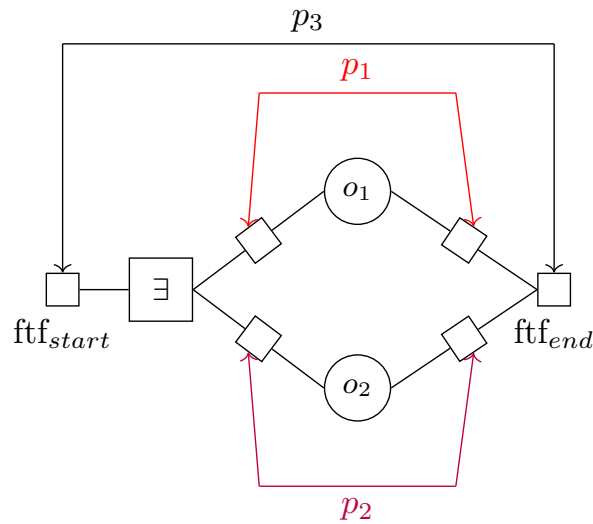


Figure 3.3: Probes inserted into an operator.

The **observed**  $\Delta Q$  for the first-to-finish operator is the  $\Delta Q$  from the instances (**start**, **end**). The **calculated**  $\Delta Q$  is the  $\Delta Q$  which is the result of the first-to-finish operator being applied on  $o_1, o_2$

### 3.6 Oscilloscope system design diagram

Now that we have an idea of how all the parts behave together and what they do, we can draw a complete diagram of the interactions of the different components.

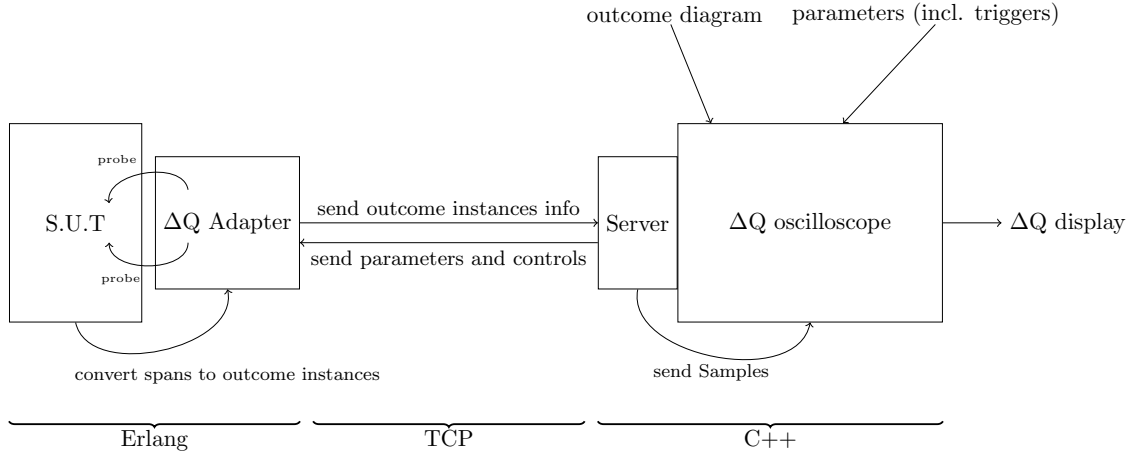


Figure 3.4: Global system design diagram.

## 3.7 Triggers

The concept of triggers is key to the oscilloscope, much like an oscilloscope that has a trigger mechanism to capture periodic signals or investigate a transient event [18], the  $\Delta Q$  oscilloscope has a similar mechanism that can recognize when an observed  $\Delta Q$  violates certain conditions regarding required behaviour and record snapshots of the system.

Each time an observed  $\Delta Q$  is calculated by the oscilloscope, it is checked against the requirements set by the user. If these requirements are not met, a trigger is fired and a snapshot of the system is saved to be shown to the user.

### 3.7.1 Snapshot

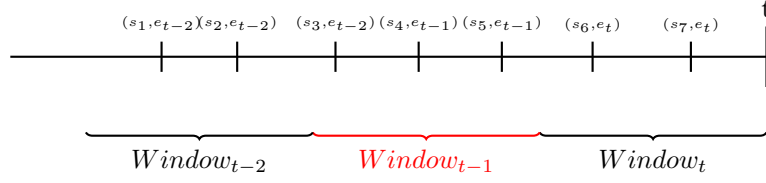
A snapshot of the system gives insights into the system before and after a trigger was fired. It gives the user a still of the system, as if it was frozen in time. All the  $\Delta Q$ s which are calculated during the system's execution are stored away. Then, if no trigger is fired, older  $\Delta Q$ s are removed. Otherwise, the oscilloscope keeps recording  $\Delta Q$ s without removing older ones, to allow the user to look at the state of the system before and after the trigger.

## 3.8 Polling window

To calculate a  $\Delta Q$ , we take all the outcome instances that ended within a window of time from  $t_l$  to  $t_u$ , a lower and upper time bound.

Suppose we are at time  $t$ , the window we will display is the **window of time**  $(t-1)_l - (t-1)_u$  with  $t-1$  equal to  $t-x$ , and  $x$  the polling rate. This is to account for various overheads that need to be taken into consideration, which could be network overhead, the wrapper overhead, C++ latency ... Imagine multiple outcome instances that are

ended at a time slightly lower but close to  $t$ , and due to the overheads the messages arrive at a time slightly higher but close to  $t$ , the outcome instance would not be taken into consideration for the calculation of a  $\Delta Q$ .



The polling window then advances every  $x$  seconds setting the new window:

$$\text{From: } (t-1)_l, (t-1)_u \xrightarrow{t+1} t_l, t_u.$$

$$\text{Where: } t_l = (t-1)_u \text{ and } t_u = (t-1)_u + x$$

# Chapter 4

## Oscilloscope: User level concepts

The following chapter gives insights on the user level concepts of  $\Delta QSD$  in the oscilloscope. They are the concepts needed by the user to understand how the oscilloscope works.

- We first provide insights into how  $\Delta QSD$  was implemented in the oscilloscope, the parameters that define a probe's  $\Delta Q$ , its representation and what can be done with  $\Delta Q$ s. We show how probe's  $\Delta Q(s)$  will be shown in the oscilloscope.
- We then provide a language to write outcome diagrams based on an already existing syntax.
- Lastly, we explain the different controls present on the oscilloscope dashboard.

### 4.1 $\Delta QSD$ implementation

Originally,  $\Delta Q(x)$  denotes the probability that an outcome occurs in a time  $t \leq x$ , defining then the "intangible mass" of such IRV as  $1 - \lim_{x \rightarrow \infty} \Delta Q(x)$ . We then extend the original definition to fit real time constraints, needing to calculate  $\Delta Q$ s continuously.

For a given probe,  $\Delta Q(t_l, t_u, dMax)$  is the probability that its instances with end time  $t_l \leq t_e \leq t_u$  occur in time  $t \leq dMax$ .

#### 4.1.1 Internal representation of a $\Delta Q$

We provide a class to calculate the  $\Delta Q$  of a probe between a lower time bound  $t_l$  and an upper time bound  $t_u$ . It can be calculated in two ways:

**Observed  $\Delta Q$**  The first way is by having  $n$  collected outcome instances between  $t_l$  and  $t_u$ , calculating its PDF and then calculating the *empirical cumulative distribution function* (ECDF) based on its PDF. This is called the **Observed  $\Delta Q$** .

**Calculated  $\Delta Q$**  A  $\Delta Q$  can also be calculated by performing operations which are the result of operations on two or more  $\Delta Q$ s, the notion of outcome instances is then

lost between calculations, as the interest shifts towards calculating the resulting PDFs and ECDFs. This is called the **Calculated  $\Delta Q$** .

### 4.1.2 dMax

The key concept of  $\Delta Q$ SD is having a maximum delay after which we consider that the execution is failed, this is represented in a prove as  $dMax$ . The user defines, for each prove the maximum delay its execution can have.

Setting a maximum delay for a probe is not a job that can be done one-off and blindly, it is something that is done with an underlying knowledge of the system inner-workings and must be thoroughly fine-tuned during the execution of the system by observing the resulting distributions of the obtained  $\Delta Q$ s.

We define in our oscilloscope a formula to dynamically define a maximum delay:

$$dMax = \Delta_T * N \quad (4.1)$$

Where  $\Delta_T$  is the bin width of the  $\Delta Q$  PDF and ECDF and  $N$  their number of bins.

The user must choose both via a slider.  $N$  is in the range  $[1, 1000]$ . This is a good enough bound to allow for finer grained representation, or less precision if needed.

Some tradeoffs must though be acknowledged when setting these parameters, a higher number of bins corresponds to a higher number of calculations and space complexity, a lower  $dMax$  may correspond to more failures. These are all tradeoffs that must be considered by the system engineer and set accordingly.

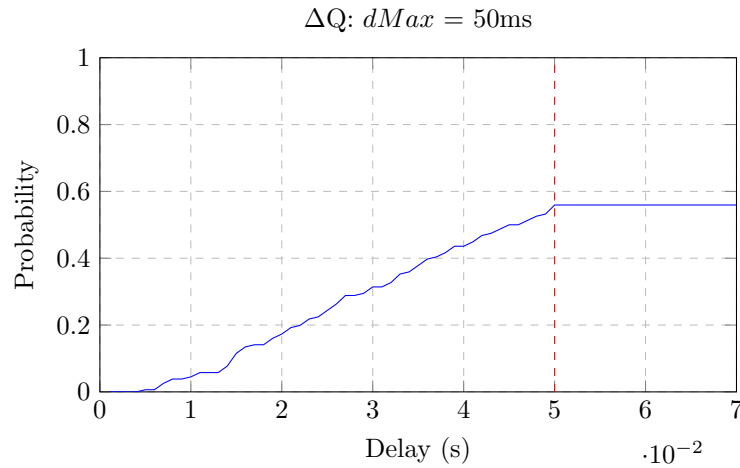


Figure 4.1:  $\Delta Q: dMax = 50ms$ , the CDF will stay constant when delay  $> dMax$

### dMax limitation

$dMax$  can **not** be lower than 1 millisecond and will be rounded to the **nearest** integer, this is a limitation of Erlang `send_after` function which only accepts integers and milliseconds values.

### 4.1.3 QTA

A simplified QTA is defined for probes. We define 4 points for the step function at 25, 50, 75 percentiles and the maximum amount of failures accepted for an observable. An observed  $\Delta Q$  will calculate that based on the samples collected.

### 4.1.4 Confidence bounds

To observe the stationarity of a system we must observe a window of  $\Delta Q$ s of a probe and calculate confidence bounds over said windows. The bounds can be updated dynamically by inserting or removing a  $\Delta Q$ , this allows us to consider a small window of execution rather than observing the whole execution.

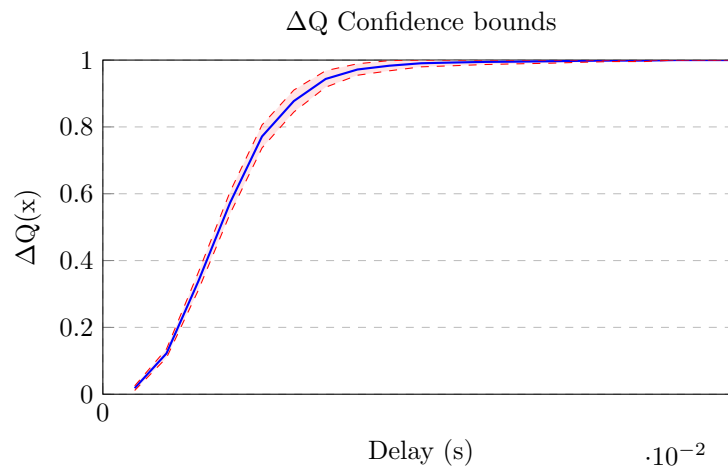


Figure 4.2: Upper and lower bounds (dashed, red) of the mean (blue) of multiple  $\Delta Q$ s. In a system that behaves linearly, the bounds will be close to the mean, once the overload is approaching, or a system is showing behaviour that diverges from a linear one, the bounds will be larger.

## 4.2 $\Delta Q$ display

A probe's displayed graph must contain the following functions:

- The mean and confidence bounds of a window of previous  $\Delta Q$ s.
- The observed  $\Delta Q$ .
- If applicable, the calculated  $\Delta Q$  from the components showing the causal links of a probe.
- Its QTA (if defined).

This allows for the user to observe if a  $\Delta Q$  has deviated from normal execution, analyse its stationarity, nonlinearity and observe its execution.



## 4.3 Outcome diagram

An abstract syntax for constructing outcome diagrams has already been defined in a previous paper [3], nevertheless, the oscilloscope needs a textual way to define an outcome diagram.

We define thus a grammar to create an outcome diagram in our oscilloscope, our grammar is a textual interpretation of the abstract syntax.

### 4.3.1 Causal link

A causal link between two components can be defined by a right arrow from `component_i` to `component_j`

```
component_i -> component_j
```

### 4.3.2 Sub-outcome diagrams

Multiple sub-outcome diagrams can be created for multiple parts of the system, these "sub-outcome diagrams" can then be linked together to form the global system outcome diagram. Recall Section 3.5.3, we defined a probe which observes the sequential composition of  $o_1, o_2$ . The probe (sub-outcome diagram)  $p_1$  can be defined as:

```
p_1 = o_1 -> o_2;
```

A probe is attached at the begin and end of  $p_1$ , it will observe the whole system and the calculated  $\Delta Q$  will be the convolution of  $o_1, o_2$ .

The lines defining these diagrams must be semicolon terminated. Outcomes and operators cannot be defined on their own, they must be observed in a sub-outcome diagram.

Sub-outcome diagrams can be reused in other diagrams by adding `s:` (sub-outcome diagram) before they are used.

```
p_3 = s:p_1 -> s:p_2;
```

This allows for easy composition and reuse of different parts of the system, allowing for independent refining of diagrams.

### 4.3.3 Outcomes

To attach a probe to an outcome observables, it is enough to declare an outcome with its name inside a diagram.

```
... = outcomeName;
```

### 4.3.4 Operators

First-to-finish, all-to-finish and probabilistic choice operators must contain at least two components.

**All-to-finish operator**

An all-to-finish operator needs to be defined as follows:

```
a:name(component1, component2...)
```

**First-to-finish operator**

A first-to-finish operator needs to be defined as follows.

```
f:name(component1, component2...)
```

**Probabilistic choice operator**

A probabilistic choice operator needs to be defined as follows:

```
p:name[probability_1, probability_2, ... probability_i](component_1,  
↪ component_2, ..., component_i)
```

In addition to being comma separated, the number of probabilities inside the brackets must match the number of components inside the parentheses. For  $n$  probabilities  $p_i$ ,  $0 < p_i < 1$ ,  $\sum_{i=0}^n p_i = 1$

**4.3.5 Limitations**

Our system has a few limitations compared to the theoretical applications of  $\Delta Q$ , namely, no cycles are allowed in the definition of outcome diagrams.

```
p_1 = s:p_2;  
p_2 = s:p_1;
```

The above example is not allowed and will raise an error when defined.

**4.4 Dashboard**

The dashboard is devised of multiple sections where the user can interact with the oscilloscope, create the system, observe the behaviour of its components, set triggers.

**4.4.1 Sidebar**

The sidebar has multiple tabs, we explain here the responsibility of each one.

**System/Handle plots tab**

**System creation** In this tab the user can create its system using the grammar defined before, he can save the text he used to define the system or load it, the system is saved to a file with any extension, we nevertheless define an extension to save the system to, the extension `.dq`. If the definition of the input is wrong, he will be warned with a pop up giving the error the parser generator encountered in the creation of a system.

**Adding a plot** Once the system is defined, the user can choose the probes he wants to plot. They can select multiple probes per plot and display multiple plots on the oscilloscope window.

**Polling rate** The user can choose the polling rate of the system: How often  $\Delta Q$ s are calculated and displayed in the oscilloscope.

**Editing a plot** By clicking onto a plot that is being shown, the user can choose to add or remove probes to and from it. Multiple probes can be selected to either be removed or added.

### Parameters tab

In this tab, the user can define parameters for the probes they have defined.

**Set a QTA** The user is given the choice to set a QTA for a given observable, they have 4 fields where they can fill in which correspond to the percentiles and the maximum amount of failures allowed, they can change this dynamically during execution.

**dMax, bins** The user has a slider which goes from -10 to 10, where they can set the parameters we explained previously,  $n$ , the exponent of  $\Delta_{tbase} \cdot 2^n$  and the bins  $N$ . When this information is saved by the user, the new  $dMax$  is transmitted to the wrapper and saved for the selected observable.

### Triggers tab

In the triggers tab the user can set triggers and observe the snapshots of the system.

**Set triggers** The user can set which triggers to fire for the probes they desire, they are given checkboxes to decide which ones to set as active or not (by default, the triggers are deactivated).

**Fired triggers** Once a trigger is fired, the oscilloscope starts a timer, during which all probes start recording the observed  $\Delta Q$ s (and the calculated ones if applicable) without discarding older ones. Once the timer expires, the snapshot is saved for the user in the triggers tab. In the dashboard, it indicates when the trigger was fired (timestamp) and the name of the probe which fired it.

### Connection controls

**Erlang controls** The user can set the IP and the port where the Erlang wrapper is listening from. Two additional buttons communicate with the wrapper by sending messages, they can start and stop the wrapper's sending of outcome instances.

**C++ server controls** The user can set the IP and the port for the oscilloscope's server.

### 4.4.2 Plots window

To the left, the main window shows the plots of the probes being updated in real time.

## 4.5 Triggers

There are two available triggers which can be selected by the user, the triggers are evaluated on the **observed**  $\Delta Q$ .

### 4.5.1 Load

A trigger on an observed  $\Delta Q$  can be fired if the amount of outcome instances received in a polling rate is greater than what the user defines:

$$\#instances(\Delta Q(t_l, t_u, dMax)) > \text{maxAllowedInstances}$$

### 4.5.2 QTA

A trigger on an observed  $\Delta Q$  can be fired if:

$$\Delta Q_{obs} > \text{observableQTA}$$

# Chapter 5

## Oscilloscope: implementation

The following chapter gives a more technical description of the oscilloscope.

- We provide a more in-depth look at the  $\Delta$ QSD concepts introduced in the previous chapter.
- We then explain how the wrapper works, its API and the underlying mechanism that let us export outcome instances to the oscilloscope.
- Next we give a technical explanation of the parser generator we used to parse the outcome diagram syntax.
- Lastly, we briefly talk about the dashboard graphical framework.

### 5.1 $\Delta$ QSD implementation

A probe's  $\Delta$ Q can be represented internally by a PDF and displayed as an ECDF. Here is how both can be calculated given  $n$  outcome instances.

#### PDF

We approximate the PDF of the observed  $\Delta$ Q via a histogram. We partition the values into  $N$  bins of equal width, this is required to ease future calculations. Given  $[x_i, x_{i+1}]$  the interval of a bin  $i$ , where  $x_i = i\Delta x$ , and  $\hat{p}(x_i)$  the value of the PDF at bin  $i$ , for  $n$  bins:

$$\begin{cases} \hat{p}(i) = \frac{s_i}{n}, & \text{if } i \leq n \\ \hat{p}(i) = 0, & \text{if } i > n \end{cases} \quad (5.1)$$

Where  $s_i$  the number of successful outcome instances whose elapsed time is contained in the bin  $i$ ,  $n$  the total number of instances.

## ECDF

The value  $\hat{f}(x_i)$  of the ECDF at bin  $i$  with  $n$  bins can be calculated as:

$$\begin{cases} \hat{f}(i) = \sum_{j=1}^i \hat{p}(j), & \text{if } i \leq n \\ \hat{f}(i) = \hat{f}(x_n), & \text{if } i > n \end{cases} \quad (5.2)$$

### 5.1.1 dMax

We introduced  $dMax$  in the previous section, we provide here the full equation that allows  $dMax$  to be calculated:

$$dMax = \Delta_{tbase} * 2^n * N \quad (5.3)$$

Where:

- $\Delta_{tbase}$  represents the base width of a bin, equal to 1ms.
- $N$  the number of bins.

### 5.1.2 Operations

In a previous section we talked about the possible operations that can be performed on and between  $\Delta$ Qs, the time complexity of FTF, ATF and PC is trivially  $\mathcal{O}(N)$  where  $N$  is the number of bins.

As to convolution, the naïve way of calculating convolution has a time complexity of  $\mathcal{O}(N^2)$ , this quickly becomes a problem as soon as the user wants to have a more fine-grained understanding of a component. Below we present two ways to perform convolution.

#### Convolution

**Naïve convolution** Given two  $\Delta$ Q binned PDFs  $f$  and  $g$ , the result of the convolution  $f \otimes g$  is given by [19]:

$$(f \otimes g)[n] = \sum_{m=0}^N f[m]g[n-m] \quad (5.4)$$

**Fast Fourier Transform Convolution** FFTW (Fastest Fourier Transform in the West) is a C subroutine library [20] for computing the discrete Fourier Transform in one or more dimensions, of arbitrary input size, and of both real and complex data. We use FFTW in our program to compute the convolution of  $\Delta$ Qs. We adapt our script from an already existing one found on GitHub. [21]

Whilst the previous algorithm is far too slow to handle a high number of bins, convolution leveraging Fast Fourier Transform (FFT) allows us to reduce the amount of calculations to  $\mathcal{O}(n \log n)$ . This is why the naïve convolution algorithm is not used. We will analyse the time gains in a later chapter.

FFT and naïve convolution produce the same results in our program barring  $\varepsilon$  differences (around  $10^{-18}$ ) in bins whose result should be 0.

FFTs algorithms are plenty, the choice of the one to use is left up to the subroutine via the parameter `FFTW_ESTIMATE` [22].

### Arithmetical operations

We can apply a set of arithmetical operations between  $\Delta$ Qs ECDFs, and on a  $\Delta$ Q.

**Scaling (multiplication)** A  $\Delta$ Q can be scaled w.r.t. a constant  $0 \leq j \leq 1$ . It is equal to binwise multiplication on ECDF bins.

$$\hat{f}_r(i) = \hat{f}(i) \cdot j \quad (5.5)$$

**Operations between  $\Delta$ Qs** Addition, subtraction and multiplication can be done between two  $\Delta$ Q of equal bin width (but not forcibly of equal length) by calculating the operation between the two ECDFs of the  $\Delta$ Qs:

$$\Delta Q_{AB}(i) = \hat{f}_A(i)[\cdot, +, -]\hat{f}_B(i) \quad (5.6)$$

### 5.1.3 Confidence bounds

To observe the stationarity of a system we must observe a window of  $\Delta$ Qs of an observable and calculate confidence bounds over said windows. We present here the formulae required to give such bounds with 68% confidence level.

For  $x_{ij}$  the value of an ECDF  $j$  at bin  $i$ , the mean of all ECDFs for the bin over a window is:

$$\mu_i = \frac{1}{n_i} \sum_{j=1}^{n_i} x_{ij} \quad (5.7)$$

Its variance:

$$\sigma_i^2 = \frac{1}{n_i} \sum_{j=1}^{n_i} x_{ij}^2 - \mu_i^2 \quad (5.8)$$

The confidence intervals  $CI_i$  for a bin  $i$  can then be calculated as:

$$CI_i = \mu_i \cdot \frac{\sigma_i}{\sqrt{n_i}} \quad (5.9)$$

### 5.1.4 Rebinning

Rebinning refers to the aggregation of multiple bins of a bin width  $i$  to another bin width  $j$ . Operations between  $\Delta$ Qs can be done on  $\Delta$ Qs that have the same bin width, this is why it is fundamental that all probes have a common  $\Delta_{tbase}$ . This allows for fast rebinning to a common bin width.

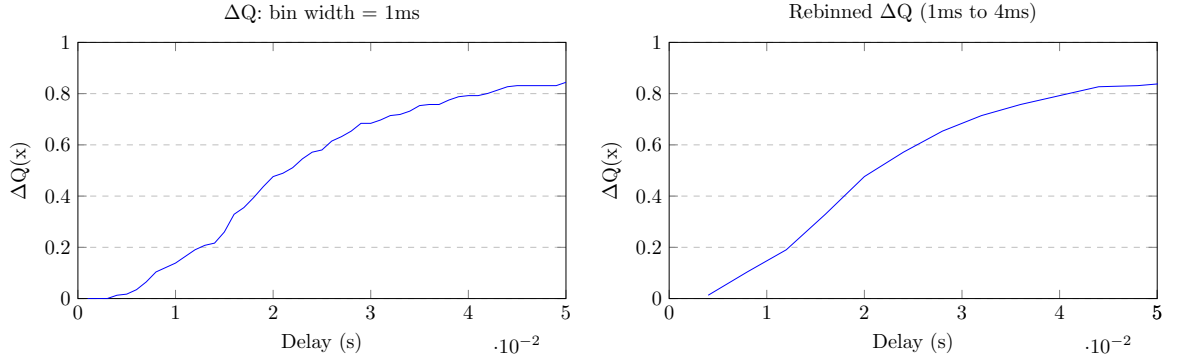
Given two  $\Delta$ Qs  $\Delta Q_i$ ,  $\Delta Q_j$ :

$$\Delta_{Tij} = \max \{ \Delta_{Ti}, \Delta_{Tj} \}$$

and the PDF of the rebinned  $\Delta Q$  at bin  $b$ , from the original PDF of  $n$  bins, where  $k = \frac{\Delta_{T_j}^i}{\Delta_{T_j}}$ .

$$p'_b = \sum_{n=b \cdot k}^{b+1 \cdot k - 1} p_n, \quad b = 0, 1, \dots, \lceil \frac{N}{k} \rceil \quad (5.10)$$

We perform rebinning to a higher bin width for a simple reason, while this leads to loss of information for the  $\Delta Q$  with the lowest bin width, rebinning to a lower bin width would imply inventing new values for the  $\Delta Q$  with the highest bin width.



(a) Sample  $\Delta Q$  with 1ms bins

(b)  $\Delta Q$  on the left after rebinning to 4ms bins

## 5.2 Wrapper

The wrapper, called `dqsd_otel` is a `rebar3` [23] application built to replace OpenTelemetry calls and create outcome instances, it is designed to be paired with the oscilloscope to observe an Erlang application.

### 5.2.1 API

The wrapper functions to be used by the user are made to replace OpenTelemetry calls to macros as for `?start_span` and `?with_span` and `?end_span`. This is to make the wrapper less of an encumbrance for the user.

Moreover, the wrapper will always start OpenTelemetry spans but only start outcome instances if the wrapper has been activated. The wrapper can be activated by the oscilloscope by pressing the "start wrapper" button and can be stopped via the "stop wrapper" button.

`start_span/1`, `start_span/2`

```
start_span/1: -spec start_span(binary()) -> {opentelemetry:span_ctx(),
  ↳ pid() | ignore}.
start_span/2: -spec start_span(binary(), map()) ->
  ↳ {opentelemetry:span_ctx(), pid() | ignore}.
```



**Parameters:**

- Name: Binary name of the probe.
- Attributes: The OpenTelemetry span attributes (Only for `start_span/2`).

`start_span` incorporates OpenTelemetry `?start_span(Name)` macro.

**Return:** The function returns either:

- `{SpanCtx, span_process_PID}` if the wrapper is active and the probe's *dMax* has been set.
- `{SpanCtx, ignore}` if one of the two previous conditions was not respected.

With `SpanCtx` being the context of the span created by OpenTelemetry.

**with\_span/1, with\_span/2**

```
with_span/1: -spec with_span(binary(), fun(() -> any())) -> any().
with_span/2: -spec with_span(binary(), map(), fun(() -> any())) ->
  ↳ any().
```

**Parameters:**

- Name: Binary name of the probe.
- Fun: Zero-arity function representing the code of block that should run inside the `?with_span` macro.
- Attributes: The OpenTelemetry span attributes (Only for `with_span/3`).

`with_span` incorporates OpenTelemetry `with_span` macro.

**Return:** `with_span` returns what `Fun` returns (`any()`).

**end\_span**

```
-spec end_span(opentelemetry:span_ctx(), pid() | ignore) -> ok |
  ↳ term().
```

**Parameters:**

- `SpanCtx`: The context of the span returned by `start_span`.
- `Pid`: `span_process_PID` || `ignore`.

As is the case for `start_span`, `end_span` incorporates an OpenTelemetry macro, in this case `?end_span(Ctx)`.

### **fail\_span**

```
-spec fail_span( pid() | ignore) -> ok | term().
```

#### **Parameter:**

- Pid: `ignore` || `span_process_PID`.

`fail_span` does not incorporate any `OpenTelemetry` macro, it is let up to the user to decide how to handle failures in execution.

### **span\_process**

`span_process` is the process, spawned by `start_span`, responsible for handling the `end_span`, `fail_span`, `timeout` messages.

Upon being spawned, the process starts a timer with time equal to the `dMax` set by an user for the probe being observed, thanks to `erlang:send_after`. When the timer runs out, it sends a `timeout` message to the process.

The process can receive three kinds of messages:

- `{end_span, end_time}`: This will send a custom span to the oscilloscope with the start and end time of the execution of the probe.
- `{fail_span, end_time}`: This will send a custom span to the oscilloscope indicating that an execution of a probe has failed.
- `{timeout, end_time(StartTime + dMax)}`: If the program hasn't ended the span before `dMax`, the timer will send a `timeout` message and it will send an outcome instance to the oscilloscope indicating that an execution of a probe has timed out.

The process is able to receive one and only message, if the execution times out and subsequently the span is ended, the oscilloscope will not be notified as the process is defunct. This is assured by Erlang documentation:

*If the message signal was sent using a process alias that is no longer active, the message signal will be dropped. [24]*

## **5.2.2 Handling outcome instances**

To create outcome instances of a probe we must obtain three important informations:

- Its name.
- The time when the span was started.
- Its `dMax`.

They start time and end time are supplied by calling this function:

```
StartTime/EndTime = erlang:system_time(nanosecond).
```

The name is given when starting a span and the *dMax* is stored in a dictionary in the wrapper.

The outcome instance is created only if two conditions are met: the wrapper has been set as active and the user set a timeout for the probe, the functions will spawn a `span_process` process, passing along all the necessary informations.

Once the span is subsequently ended/timed out/failed, the function `send_span` creates a message carrying all the informations and sends it to the C++ server. The formatting of the messages is the following:

```
n:Observed name, b: Start time (beginning), e: End time (end or  
↪ deadline), s: The status
```

### 5.2.3 TCP connection

The wrapper is composed of two `gen_server` which handle communication to and from the oscilloscope. This `gen_server` behaviour allows the wrapper to send spans asynchronously to the oscilloscope.

#### TCP server

The TCP server is responsible for receiving commands from the oscilloscope. It can be run by setting its IP and port via:

```
-spec start_server(string() | binary() | tuple(), integer()) -> ok |  
↪ {error, Reason}
```

The oscilloscope can send commands to the wrapper, these commands are:

- `start_stub`: This command sets the wrapper as active, it can now send outcome instances to the oscilloscope if the probe's *dMax*s are defined.
- `stop_stub`: This commands sets the wrapper as inactive, it will no longer send outcome instances to the oscilloscope.
- `set_timeout;probeName;timeout`: This command indicates to the wrapper to set the *dMax* = timeout for a probe, a limit of the wrapper is that erlang:send\_after does not accept floats as timeouts, so the timeout will be rounded to the nearest integer.

#### TCP client

The TCP client allows the wrapper to send the spans to the oscilloscope. The client connects over TCP to the oscilloscope by connecting to the oscilloscope server's address and opens a socket where it can send the outcome instances.

```
-spec try_connect(string() | binary(), integer()) -> ok.
```

## 5.3 Parser

To parse the system, we use the C++ ANTLR4 (ANother Tool for Language Recognition) library.

### 5.3.1 ANTLR

ANTLR is a parser generator for reading, processing, executing or translating structured text files. ANTLR generates a parser that can build and walk parse trees [25].

ANTLR is just one of the many parsers generators available in C++ (flex/bison, lex, yacc), although it presents certain limitations, its generated code is simpler to handle and less convoluted with respect to the other possibilities.

ANTLR uses Adaptive LL(\*) (*ALL*(\*)) parser, namely, it will move grammar analysis to parse-time, without the use of static grammar analysis. [26]

### 5.3.2 Grammar

ANTLR provides a yacc-like metalanguage [26] to write grammars. Below, is the grammar for our system:

```
grammar DQGrammar;

PROBE_ID: 's';
BEHAVIOR_TYPE: 'f' | 'a' | 'p';
NUMBER: [0-9]+('.'[0-9]+)?;
IDENTIFIER: [a-zA-Z_][a-zA-Z0-9_]*;
WS: [ \t\r\n]+ -> skip;

start: definition* system? EOF;

definition: IDENTIFIER '=' component_chain ' ';

component_chain : component ('->' component)*;

component : behaviorComponent | probeComponent | outcome;

behaviorComponent : BEHAVIOR_TYPE ':' IDENTIFIER ('[' probability_list
    ↪ ']')? '(' component_list ')';

probeComponent : PROBE_ID ':' IDENTIFIER;

probability_list: NUMBER (',' NUMBER)+;
component_list: component_chain (',' component_chain)+;
outcome: IDENTIFIER;
```

## Limitations

A previous version was implemented in Lark [27], a python parsing toolkit. The python version was quickly discarded due to a more complicated integration between Python and C++. Lark provided Earley(SPPF) strategy which allowed for ambiguities to be resolved, which is not possible in ANTLR.

For example the following system definition presents a few errors:

```
probe = s -> a -> f -> p;
```

While Lark could correctly guess that everything inside was an outcome, ANTLR expects ":" after "s, a, f" and "p", thus, one can not name an outcome by these characters, as the parser generator thinks that an operator or a probe will be next.

## 5.4 Oscilloscope GUI

Our oscilloscope graphical interface has been built using the QT framework for C++. Qt is a cross-platform application development framework for creating graphical user interfaces. [28] We chose Qt as we believe that it is the most documented and practical library for GUI development in C++, using Qt allows us to create usable interfaces quickly, while being able to easily pair the backend code of C++ to the frontend.

The interface is composed of a main window, where widgets can be attached to it easily. Everything that can be seen is customisable widgets. This allows for easy reusability, modification and removal without great refactoring due in other parts of the system.

In the photo below we can see each top level widget (a QWidget that contains other widgets) in the main window, the widgets could easily be switched to other places of the window or rearranged.

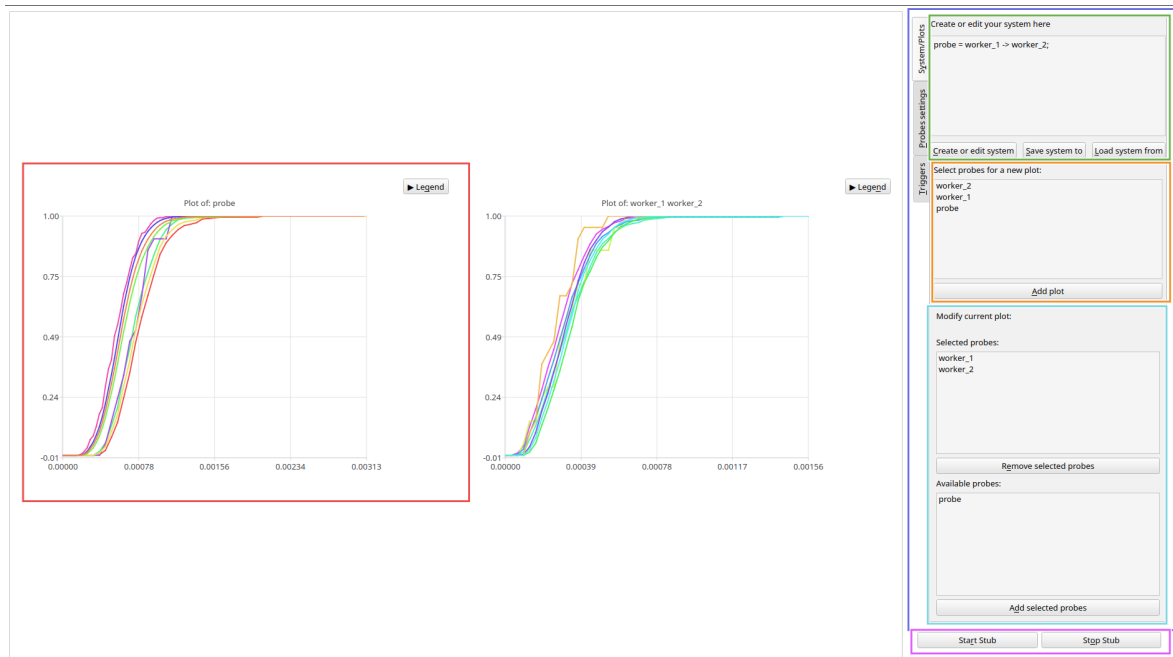


Figure 5.2: The oscilloscope displaying probes, the boxes represent a top level widget, which may contain other widgets inside.

# Chapter 6

## Application on synthetic programs

This section aims to provide an example of how the oscilloscope could be used to instrument an application, in this case, a synthetic one. We explain how the  $\Delta$ QSD paradigm can be applied to explore tradeoffs in design and to gain more insights into a running system.

### 6.1 System with sequential composition

We model a first system with two sequentially composed component. We choose two model the two components as M/M/1/K queues.

**Why M/M/1/K queues?** An average component in a distributed system can be modeled as an M/M/1/K, due to the exponential inter-arrival rate of messages  $\lambda$ , the exponential distribution of the execution delay  $\mu$ , the buffer size of messages  $K$  of a component and the failure rate  $f$ . [1]

Let us first provide a refresher about M/M/1/K queues:

- $\lambda$ : The arrival rate.
- $s$ : The service time, is the time it takes to serve a message.
- $\mu$  The service **rate** and  $E[s] = \frac{1}{\mu}$
- Offered load  $\rho = \frac{\lambda}{\mu}$

We will control  $\lambda$  to show its effects on the offered load. The offered load can tell much about the system:

- At low load ( $\rho < 0.8$ ) the failure will tend to 0, the system is behaving correctly and the  $\Delta$ Q will show that, the delay will tend to 1.
- Once  $\rho$  is approaching high load ( $\rho > 0.8$ ) we can observe the failure increasing quickly, but we can observe the system starting to get bad after  $\rho > 0.5$ ! [1]

### 6.1.1 System composition

The system has two components **worker\_1**, **worker\_2**. Each individual component is made of a buffer queue of size  $K = 1000$  and a worker process.

The system sends  $n$  messages per second following a Poisson distribution to **worker\_1**'s queue, the queue then reduces its available buffer size.

The buffer notifies its worker, which then does  $N$  loops, which are defined upon start, of fictional work. The worker then passes a message to **worker\_2**'s queue, which has another queue of same size, who passes the message to **worker\_2**'s worker, which does the same amount of loops. When a worker completes its work, it notifies the queue, freeing one "message" from its buffer size.

If the queue's buffer is overloaded, it will drop the incoming message and consider the execution a failure.

A probe  $p$  is defined, which observes the execution from when the first message up until **worker\_2** is done.

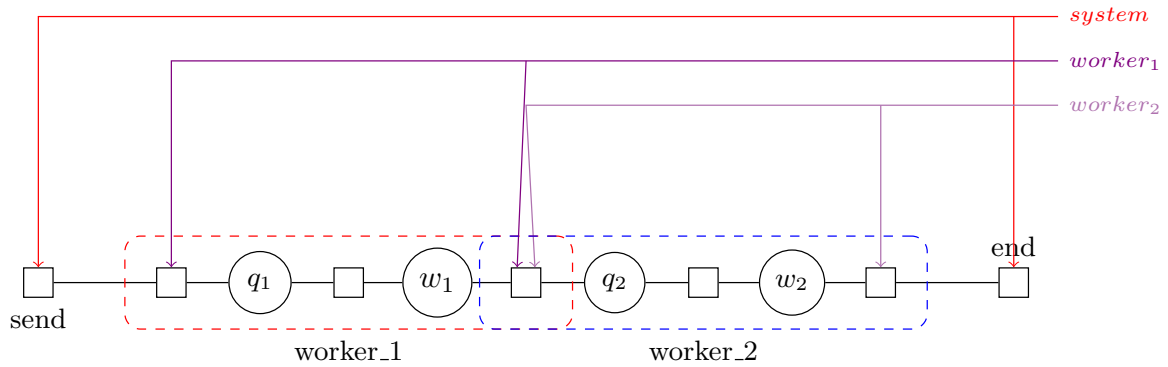


Figure 6.1: Outcome diagram of the M/M/1/K queue with the colored lines representing the probes that were inserted.

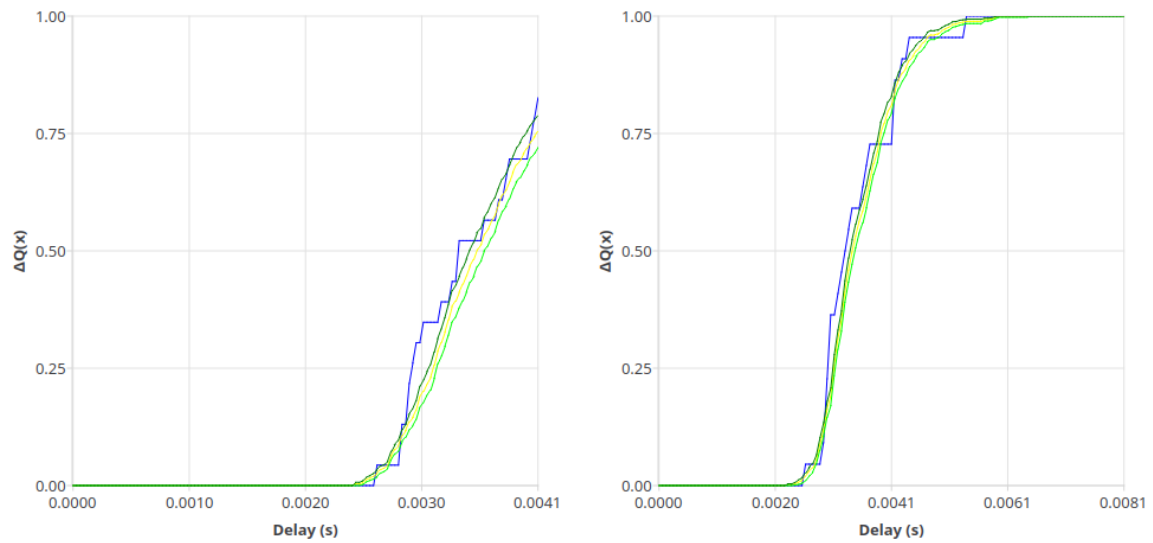
### 6.1.2 Determining parameters dynamically

We stated previously that determining parameters is something that must be done with an underlying knowledge of the system. The oscilloscope can provide knowledge of the system, here is an example of **worker\_1** and **worker\_2** as observed in the oscilloscope.

Imagine the engineer supposes the workers executions should a maximum of 4 ms to complete, but doesn't actually know how long the executions should take. The engineer, after having set the required parameters observes in the following graph in the oscilloscope Section 6.1.2.

The oscilloscope shows the engineer that their assumptions do not correspond to the actual system  $\Delta Q$ , the user can then modify the parameters to observe the actual system's behaviour. By setting  $dMax$  to 8 ms, they can observe the worker's  $\Delta Q$ s failure approaching 0.





(a) worker\_1  $\Delta Q$  with confidence bounds plot with 4 ms  $dMax$ . (b) worker\_1  $\Delta Q$  with confidence bounds plot with 8 ms  $dMax$ .

On the other hand, the engineer's assumption could have been what he truly expected from the system, in this case, the oscilloscope tells him that the system is not behaving as expected.

### Low Load

Let's first observe the system at low load, we will send 50 messages per second to observe the system under test to get key properties.

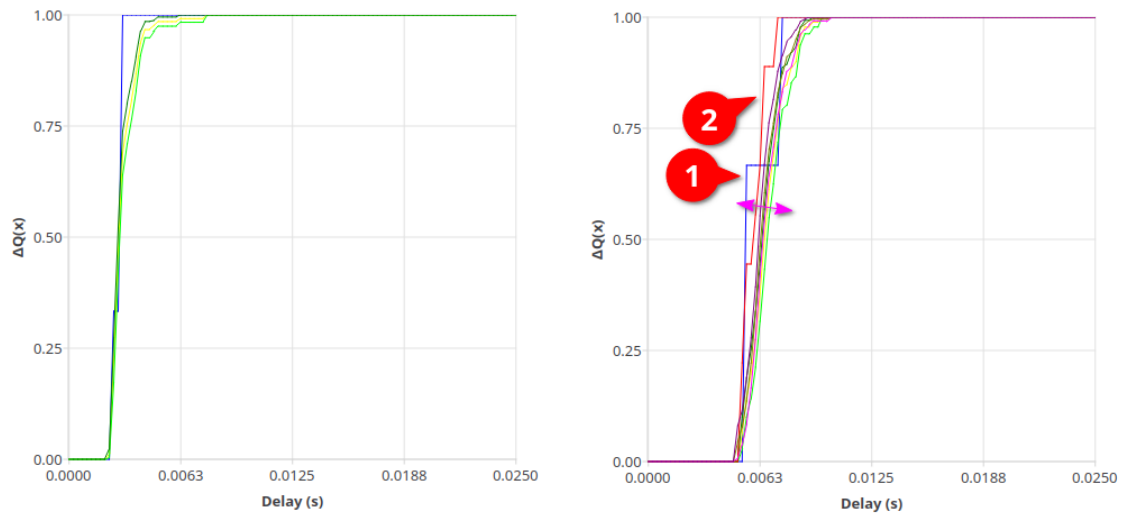


Figure 6.3: Left: worker\_1  $\Delta Q$ , in blue, the observed  $\Delta Q$ , in green, the confidence bounds. Right: probe  $\Delta Q$ , in blue, the observed  $\Delta Q$ , in red the calculated  $\Delta Q$ . In magenta, the two  $\Delta Q$ s confidence bounds overlapping.

We first observe the worker's  $\Delta Q$ , we can observe that the average execution takes  $\approx 30\text{ms}$ . We then have  $\mu_{worker} = \frac{1}{0.0033} \approx 300 \text{ req/s}$ . Thus  $\rho = \frac{50}{322} = 0.16$ , we are in nice grounds!

At low load, we can observe in the oscilloscope the probe **observed  $\Delta Q$**  and **calculated  $\Delta Q$**  confidence bounds overlap.

### Early signs of overload

We can see how even at load = 0.5 the system is starting to show bad behaviour. Let us observe what happens when  $\lambda = 150 \rightarrow \rho = 0.5$ .

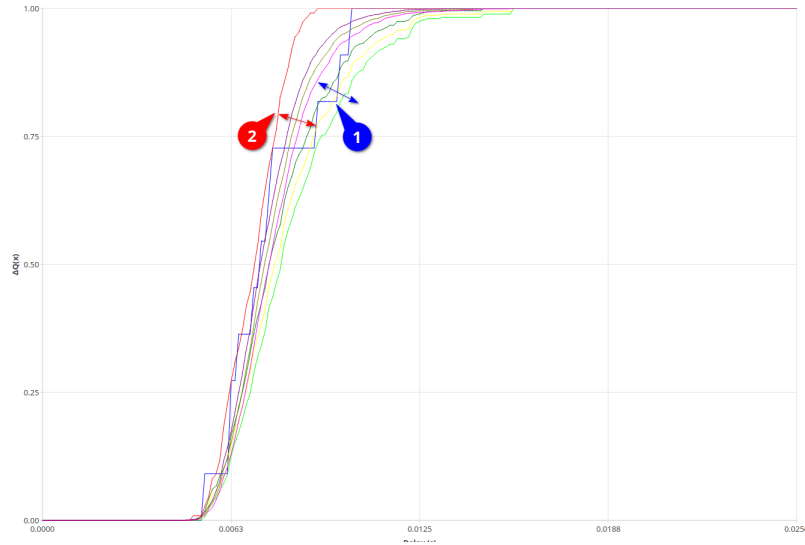
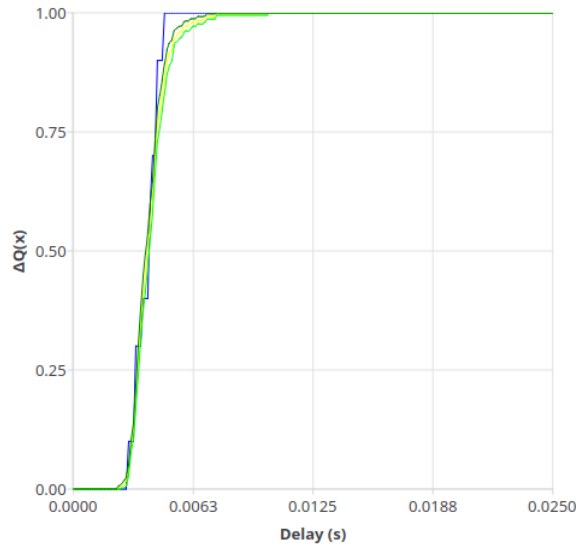


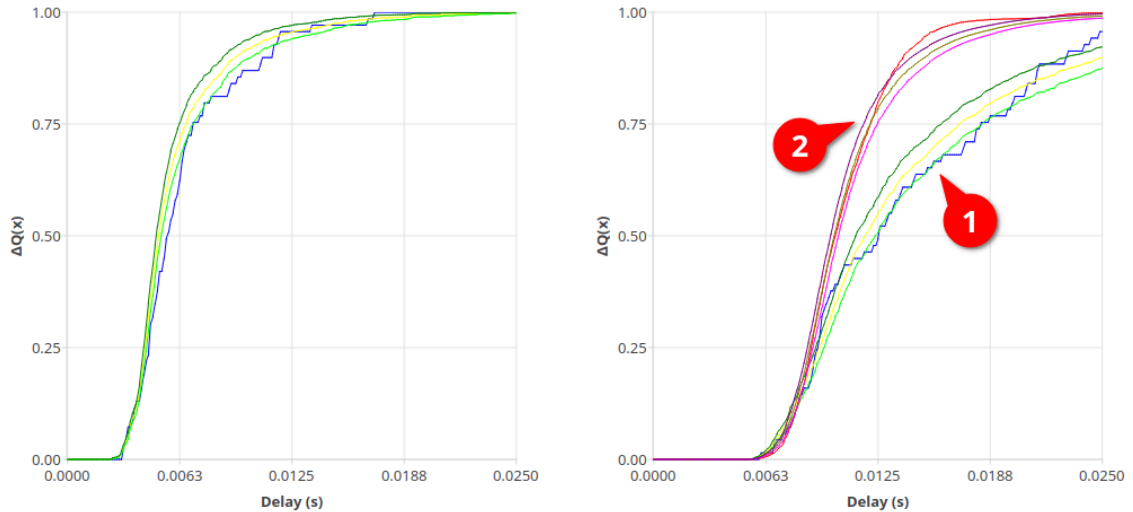
Figure 6.4: probe  $\Delta Q$ s, in blue (1), the observed  $\Delta Q$ , in red (2) the calculated  $\Delta Q$ . Arrow, above, blue: Observed  $\Delta Q$  confidence bounds. Arrow, below, red: Calculated  $\Delta Q$  confidence bounds.

Figure 6.5: worker\_1  $\Delta Q$  (blue) and its confidence bounds

Recall 2.8, we can start to observe early signs of dependency! At load 0.5 the calculated  $\Delta Q$  is deviating from the observed one. This is a sign that the performance is degrading. Worker\_1 is slowing down, but we nevertheless do not observe failures in probe's  $\Delta Q$ s.

### High load

Performance at 0.5 offered load are already remarkable, the  $\Delta Q$ s are shifting. We can go even further and observe the system under high load situations. We set  $\lambda = 200 \rightarrow \rho 0.83$ , just above the high load threshold.

Figure 6.6: Left: worker\_1  $\Delta Q$ . Right: probe  $\Delta Q$ , in blue (1), the observed  $\Delta Q$  with its confidence bounds, in red (2) the calculated  $\Delta Q$  with its confidence bounds.

This is what we expected previously, and confirms what is expected by queueing theory,  $\Delta Q$  is capable of observing the basic observation requirements and capable of recognising dependency. While what is expected by the execution of the queue (observed  $\Delta Q$ ) is a nice normally distributed CDF with little to no failure. What we can actually observe is a degraded performance in both workers and the probe observed execution.

The workers CDF has completely degraded, with the average request taking almost double the time as under normal queueing conditions.

Further degradation can be observed by increasing  $\lambda = 300, 350 \rightarrow rho \geq 1$ .

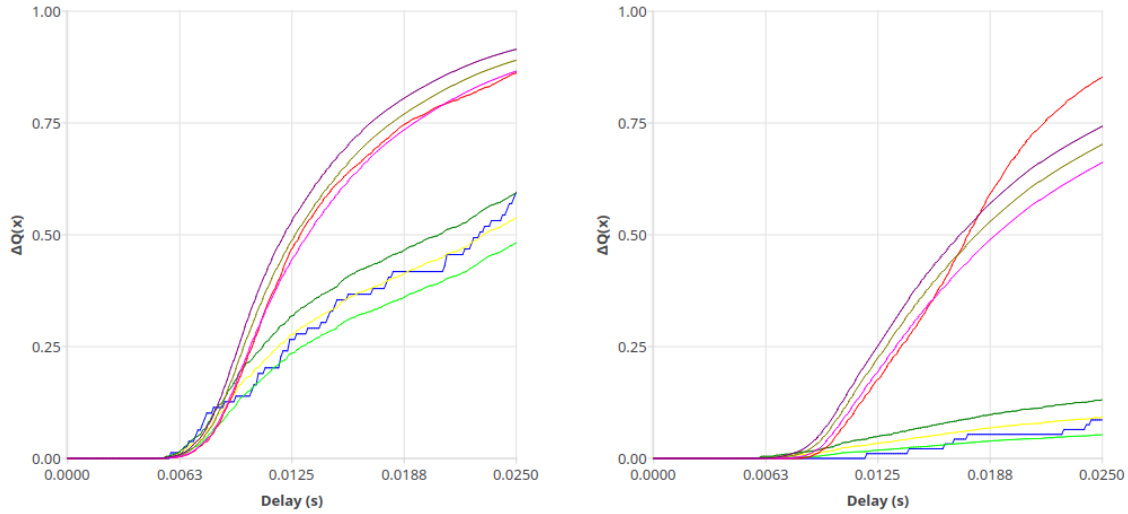


Figure 6.7: Left: probe  $\Delta Q$ s at  $\lambda = 300$ . Right: probe  $\Delta Q$ s at  $\lambda = 350$

The system degrading clear, the  $\Delta Q$ s show how almost all messages are being dropped or take  $> dMax$ . Let us look at triggers and how they can be useful to diagnose such cases.

### Triggers

By observing the system under test in high load cases, we can set the load trigger by setting the sampling window to 1 second and trigger when outcome instances  $\gtrsim 150$ . We can also set a trigger based on observation of the running system.

**QTA trigger** By observing the system, we create a QTA for the probe with: 25% = 0.0075 s, 50% = 0.0125 s, 75% = 0.015s and minimum intangible mass = 0.9.

By setting the trigger to fire for  $\Delta Q_{obs} < QTA$ . We captured a handful of snapshots. Here,  $\lambda = 150$ .

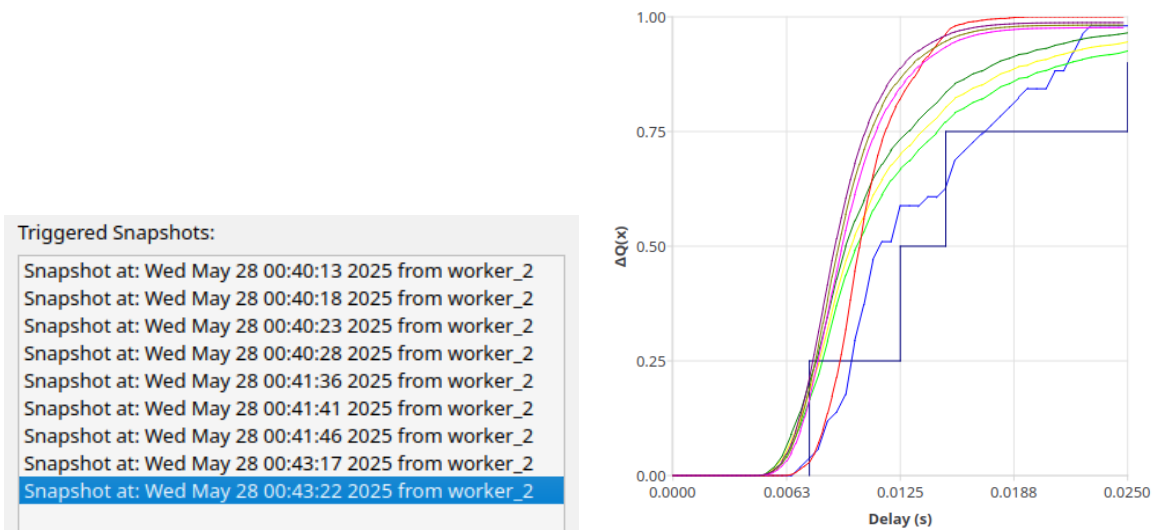
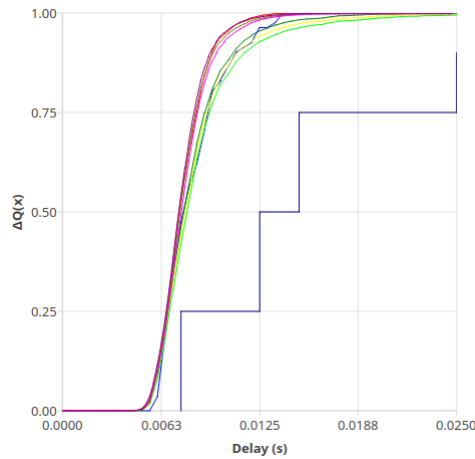


Figure 6.8: Left: Fired triggers. Right: Sample QTA violation.

QTA triggers can help to detect overhead even before high load becomes evident!

**Instances trigger** By knowing the inner details of the system, setting a QTA on the number of instances can be useful. Here is an example of a fired trigger on the number of instances.

Even though the QTA requirement isn't being violated, the number of instances fires a trigger, where the user can observe that the system is showing early signs of overload.



## 6.2 Detecting slower workers in workers

### 6.2.1 First to finish application

Next, we provide a synthetic application modeling an application that can be modeled by a first to finish operator

**Why first to finish?** Recall the previous FTF graph Figure 2.5. Assume a send request to "the cloud" that waits for a response or a timeout, it is modeled by a FTF operator.

### Using the wrong operator

What happens if the wrong operator is chosen to represent the causal relationships between the outcomes? What if the user believes that the system diagram is the one we presented before Figure 6.1? The result on the oscilloscope will clearly show that something is wrong!

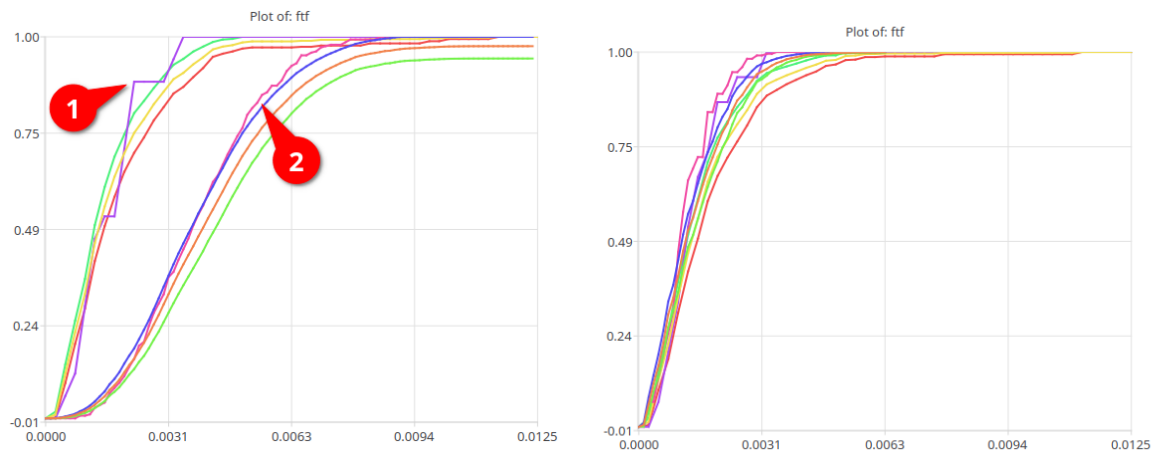


Figure 6.9: *(Left)* FTF plot **with wrong outcome diagram definition** as shown in the oscilloscope. (1) Observed  $\Delta Q$ . (2) Calculated  $\Delta Q$ . *(Right)* FTF plot **with correct outcome diagram definition** as shown in the oscilloscope. Observed  $\Delta Q$  and calculated  $\Delta Q$  overlapping.

On the left, we can observe how the **calculated  $\Delta Q$**  (2) is clearly greater than the **observed  $\Delta Q$**  (1). A difference this drastic tells us that the proposed outcome diagram does not correctly represent the actual system. On the right, if no dependencies are present and the correct operator is chosen, the two graphs will overlap.

### Introducing a slower component

Let us introduce a slower worker into the system, we introduce an artificial delay into worker\_2 (about 20ms). If the oscilloscope works correctly, the paradigm operations are sound and no dependencies are present in the system, we should not see any difference in the observed and calculated  $\Delta Q$ s of the FTF operator.

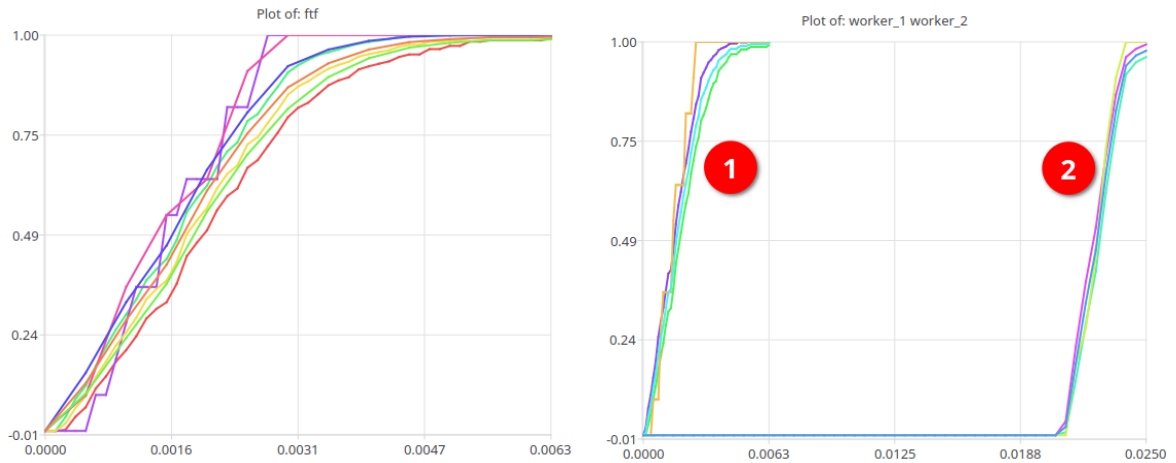


Figure 6.10: (Left) FTF plot of worker\_1 and worker\_2, observed and calculated  $\Delta Q$  overlapping.

(Right) worker\_1 (1) and worker\_2 (2)  $\Delta Q$ s.

The FTF plot correctly displays how worker\_2 does not have an effect on the ftf plot.

## 6.2.2 All to finish application

We can extend the previous application to an all-to-finish operator, this operator can for instance parallel work, a task that requires a lot of computation and can be done in separate pieces by separate workers. [1]

### Introducing a slower component

Like we did for the FTF operator, let's introduce a slower work into the mix. We introduce a slight delay to show how even a few milliseconds can be noticeable right away by a keen eye (or by triggers, which avoids having to look constantly at the graphs). The delay is a 2ms sleep on worker\_2.

**Worker's performance** We present here the worker's graphs as shown on the oscilloscope.

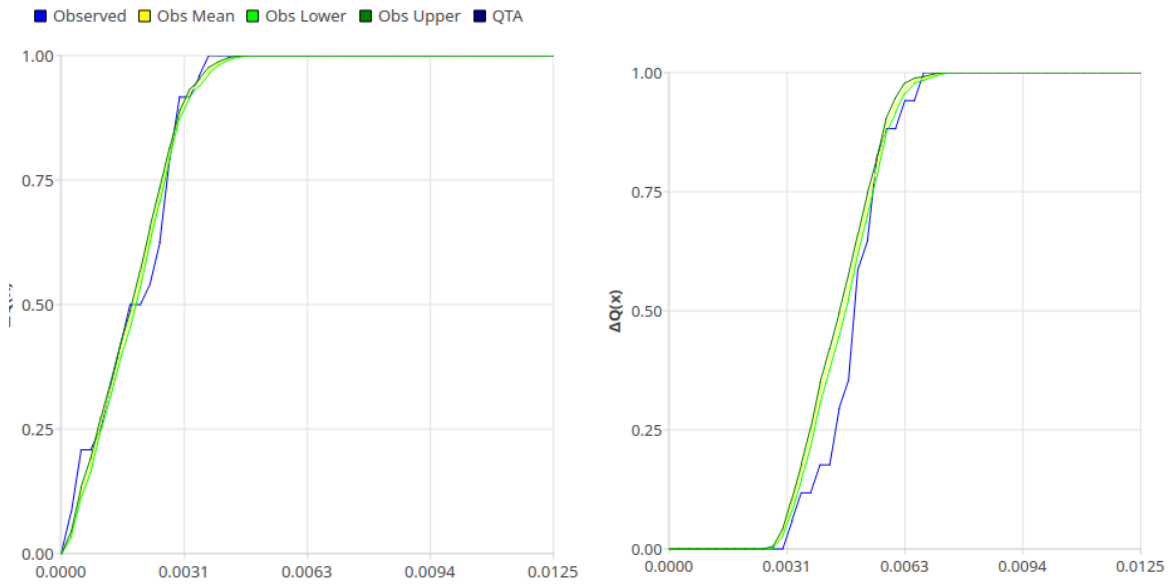
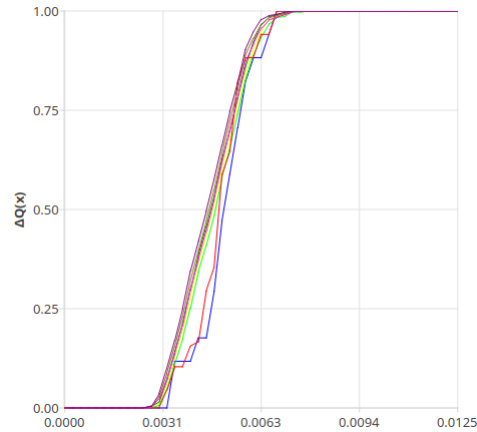


Figure 6.11: Left: worker\_1 plot. Right: worker\_2 plot.

The difference in the worker's  $\Delta Q$  can be noticed with  $\Delta Q_w2 > \Delta Q_w1$ . The difference can then be observed in the all-to-finish plot, where the operator's  $\Delta Q$ s (both observed and calculated) can be overlaid on top of worker\_2  $\Delta Q$ , showing once again that the  $\Delta QSD$  foundation is sound.



These plots show the usefulness of  $\Delta QSD$ , the system can be decomposed to understand which part of the system is showing hazards, furthermore, the causal relationships can be observed to determine the behaviour of a part down to the single component.



# Chapter 7

## Performance study

This chapter evaluates the components and operations we introduced in previous sections, analysing their performances

### 7.1 Convolution performance

We implemented two versions of the convolution algorithm as described before, the naïve version and the FFT version. We compared their performance when performing convolution on two  $\Delta Q$ s of equal bins.

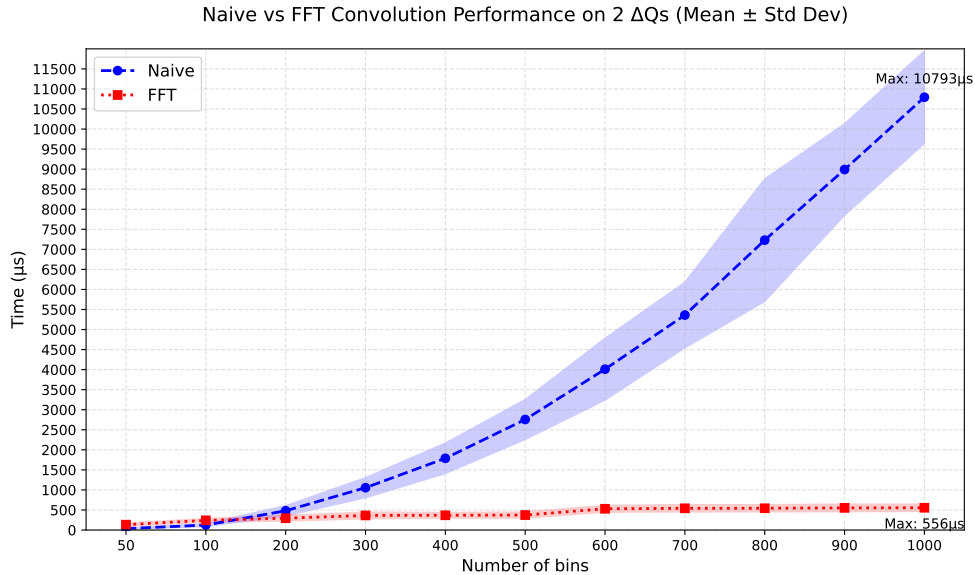


Figure 7.1: Performance comparison of two convolution algorithms

As expected, the naïve version has a time complexity of  $\mathcal{O}(n^2)$  and quickly scales with the number of bins, this is clearly inefficient, as a more precise  $\Delta Q$  will result in a much slower program.

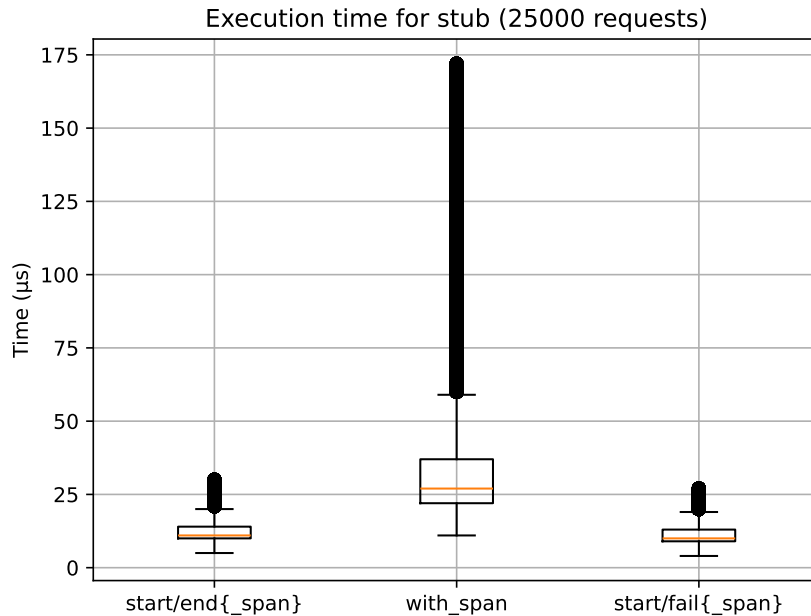
As for the FFT algorithm, it is slightly slower when the number of bins is lower than 100. This is due to the FFTW3 routine having slightly higher overhead.

## 7.2 Wrapper performance

We evaluated the performance of the wrapper to measure its impact in a normal execution, namely we tested the following calls which represent a normal usage of the wrapper.

- `start_span`  $\rightarrow$  `end_span`.
- `with_span` with the following function: `fun()`  $\rightarrow$  `ok`.
- `start_span`  $\rightarrow$  `fail_span`.

We ran the simulation for 25000 subsequent iterations, these are the results.

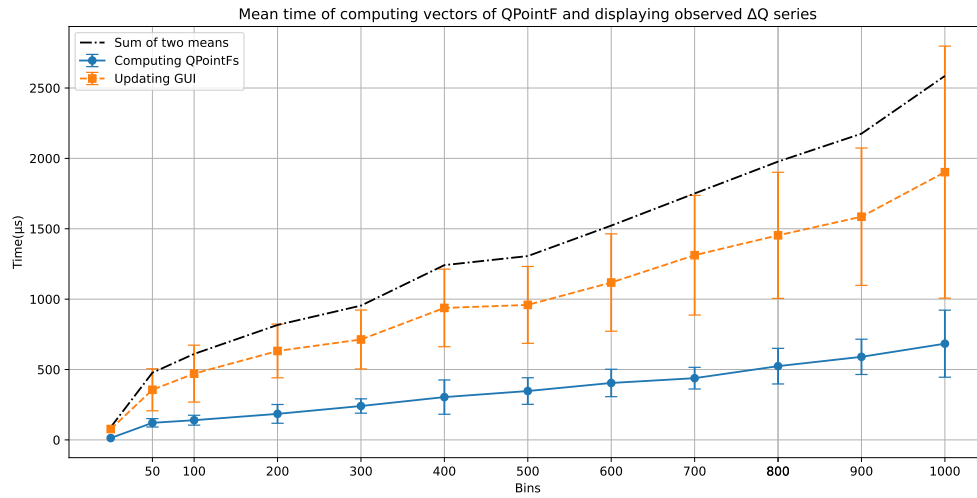


The overhead is minimal, around 10 microseconds on average to start and end/fail a span. The same cannot be said about with span, the increased overhead is nevertheless due to a function needing to be called inside it for it to record a span.

## 7.3 GUI plotting performance

We evaluated the performance of the GUI plotting routine for an observed  $\Delta Q$ . The routine first prepares the  $\Delta Q$ , creating vectors `QPointF` (a Qt class representing a point for a `QtChart`), representing the x and y values of the  $\Delta Q$ s CDF. The vectors are created for the lower bound, the upper bound the mean of the window of  $\Delta Q$ s and the observed  $\Delta Q$ .

Then, once the vectors are prepared, Qt replaces the old points with the new points for every element being plotted.



The result scales up to 2 ms for 1000 bins, since multiple plots have to be shown, if the probe observes the causal relationships of outcome and operators, the calculated  $\Delta Q$  will need to be plotted too, the time increase will be twofold. If a user decides to have finer grained representations and multiple plots at once, decreasing (decreasing or increasing?) the polling rate will avoid the plots frame skipping.

# Chapter 8

## Conclusions and future work

The following project is the beginning of the  $\Delta Q$  oscilloscope, our initial goal was to create an application to observe running distributed applications, namely, Erlang ones. A prototype was successfully created thanks to the following feats:

- The graphical dashboard for the  $\Delta Q$  oscilloscope, built in C++, which allows real time display of  $\Delta Q$ s for the probes inserted in the system.
- Fast convolution algorithms to perform statistical analysis on probes.
- The creation of a textual syntax to create outcome diagrams.
- The `dqsd_otel` Erlang wrapper to connect an OpenTelemetry instrumented Erlang app to the oscilloscope.

The user has full control over the outcome diagrams and can update them dynamically to add or remove probes, this allows full control of what the user decides to include, allowing a finer grained representation or a more general view of the system.

The oscilloscope and the Erlang can communicate via TCP socket connections to exchange outcome instances and probe parameters,

We showed how it can be useful in detecting early signs of overload many crucial features are still missing from the dashboard, and it could require less code modifications in the Erlang side. The next important step of the oscilloscope is its trial in a true distributed application. This would further reinforce the solidity of the paradigm in detecting problems in design of large systems.

### 8.1 Future improvements

We believe the oscilloscope and the Erlang application can be drastically improved, the size of the project and its intended goal is too big to be encompassed in a single master thesis. We list here some improvements which could be made to both the oscilloscope and the wrapper.

### 8.1.1 Oscilloscope improvements

- The oscilloscope could be turned into a **web app**, we feel that a C++ oscilloscope is a good prototype and proof of concept, but its usability would be greater in a browser context. It would be great as a plugin for already existing observability platforms like Grafana.
- A wider selection of **triggers**, as of writing this thesis, only the QTA trigger and load are available, this is a limitation due to time constraints. Nevertheless, triggers can be easily implemented in the available codebase.
- **Better communication between stub - server - oscilloscope.** The current way of sending outcome instances may be a limiting factor under high load, if hundred of thousands of spans were to be sent, the current way the server and oscilloscope are tied together may throttle communications. TCP socket connections could quickly become the chokepoint which makes the oscilloscope temporarily unusable.

Future improvements on the server side could implement epoll system server calls to make the server more efficient; **Detaching server from client**, as of right now, the oscilloscope and the server are tied together, using ZeroMQ to assure real time server-client communications could be an interesting solution to explore.

- **Improve real time graphs.** The class QtCharts does not perform correctly with high frequencies update. Moreover, since we are plotting multiple series (from a minimum 4 to a maximum of 9) per probe, which allows up to 1000 bins per probe, the performance quickly degrades with more probes being displayed. A better graphing class for Qt could definitely improve the experience.
- **Saving probe parameters:** As of writing this thesis, there is no way to save the parameters one may have set.
- **Deconvolution:** An important aspect of  $\Delta$ QSD, which was not introduced in this paper is deconvolution. It is used to check for infeasability in system desing. Since convolution has already been implemented, this could be integrated using the FFTW3 library.
- **Exporting graphs:** The graphs can only be observed in the oscilloscope and have no way to be exported to other programs via standard formats.
- **Many more:** This oscilloscope is just a start, if we were to list everything we may want to add, it would take many pages. What we provide is a sufficient enough basis to provide possibilities to observe a running system and understand the power of  $\Delta$ QSD in analysing its behaviour.

### 8.1.2 Wrapper improvements

- As suggested by Bryan Naegele, a member of the observability group of Erlang, the wrapper, instead of working on top of OpenTelemetry, could be directly included inside the context of a span by using the ctx library [29], which provides deadlines for contexts, propagating the value in `otel_ctx`, making it available

to the OpenTelemetry span processor. Leveraging `erlang:send_after` as we already do, we could create outcome instances with telemetry events to handle successful executions and timeouts. The span processor will then be responsible for creating outcome instances, without creating the need for custom functions in the wrapper, like we have now.

### 8.1.3 Real applications

A flaw of the oscilloscope and wrapper is that they have not been tested on real applications, while their usefulness has been proven on synthetic applications, the lack of real life applications is a weakness.

### 8.1.4 Licensing limitations

Lastly, a notable limitation is created by **Qt**, namely, QtCharts. The usage of Qt does not allow us to release our project under BSD/MIT licenses, but rather a GPLv3 one (we cannot release it under LGPL due to QtCharts). [30]

# Bibliography

- [1] Peter Van Roy and Seyed Hossein Haeri. *The  $\Delta$ QSD Paradigm: Designing Systems with Predictable Performance at High Load. Full-day tutorial*. 15th ACM/SPEC International Conference on Performance Engineering. 2024. URL: <https://webperso.info.ucl.ac.be/~pvr/ICPE-2024-deltaQSD-full.pdf>.
- [2] Peter Thompson and Rudy Hernandez. *Quality Attenuation Measurement Architecture and Requirements*. Tech. rep. MSU-CSE-06-2. Sept. 2020. URL: <https://www.broadband-forum.org/pdfs/tr-452.1-1-0-0.pdf>.
- [3] Seyed Hossein Haeri et al. “Algebraic Reasoning About Timeliness”. In: *Proceedings 16th Interaction and Concurrency Experience, ICE 2023, Lisbon, Portugal, 19th June 2023*. Ed. by Clément Aubert et al. Vol. 383. EPTCS. 2023, pp. 35–54. DOI: 10.4204/EPTCS.383.3. URL: <https://doi.org/10.4204/EPTCS.383.3>.
- [4] Erlang/OTP. *What is Erlang*. Accessed: (26/05/2025). 2025. URL: <https://www.erlang.org/faq/introduction.html>.
- [5] Seyed H. Haeri et al. “Mind Your Outcomes: The  $\Delta$ QSD Paradigm for Quality-Centric Systems Development and Its Application to a Blockchain Case Study”. In: *Comput.* 11.3 (2022), p. 45. DOI: 10.3390/COMPUTERS11030045. URL: <https://doi.org/10.3390/computers11030045>.
- [6] Peter Van Roy. *LINFO2345 lessons on  $\Delta$ QSD*. Accessed: (19/05/2025). UCLouvain, 2023. URL: <https://www.youtube.com/watch?v=tF7fbU9Gce8>.
- [7] Neil J. Davies and Peter W. Thompson.  *$\Delta$ QSD workbench - GitHub*. Accessed: (19/05/2025). 2022. URL: <https://github.com/DeltaQ-SD/dqsd-workbench>.
- [8] Erlang programming language. *Erlang tracing*. Accessed: (19/05/2025). 2024. URL: <https://www.erlang.org/doc/apps/erts/tracing.html>.
- [9] OpenTelemetry. *OpenTelemetry in Erlang/Elixir*. Accessed: (19/05/2025). 2025. URL: <https://opentelemetry.io/docs/languages/erlang/>.
- [10] OpenTelemetry. *What is OpenTelemetry?* Accessed: (19/05/2025). 2025. URL: <https://opentelemetry.io/docs/what-is-opentelemetry/>.
- [11] OpenTelemetry. *OpenTelemetry - Traces*. Accessed: (19/05/2025). 2025. URL: <https://opentelemetry.io/docs/concepts/signals/traces/>.
- [12] The Jaeger Authors. *Jaeger*. Accessed: (19/05/2025). 2025. URL: <https://www.jaegertracing.io/>.
- [13] Dotan Horovits. *From Distributed Tracing to APM: Taking OpenTelemetry & Jaeger Up a Level*. Accessed: (19/05/2025). 2021. URL: <https://logz.io/blog/monitoring-microservices-opentelemetry-jaeger/>.
- [14] Sampath Siva Kumar Boddeti. *Tracing Made Easy: A Beginner’s Guide to Jaeger and Distributed Systems*. Accessed: (19/05/2025). 2024. URL: <https://>

- openobserve.ai/blog/tracing-made-easy-a-beginners-guide-to-jaeger-and-distributed-systems/.
- [15] OpenTelemetry. *Instrumentation for OpenTelemetry Erlang/Elixir*. Accessed: (19/05/2025). 2025. URL: <https://opentelemetry.io/docs/languages/erlang/instrumentation/>.
  - [16] OpenTelemetry. *Active spans, C++ Instrumentation*. Accessed: (19/05/2025). 2025. URL: <https://opentelemetry.io/docs/languages/cpp/instrumentation/>.
  - [17] Hazel Weakly. *OpenTelemetry Challenges: Handling Long-Running Spans*. Accessed: (21/05/2025). 2024. URL: <https://thenewstack.io/opentelemetry-challenges-handling-long-running-spans/>.
  - [18] KeySight. *What is an Oscilloscope Trigger?* Accessed: (23/05/2025). 2022. URL: <https://www.keysight.com/used/id/en/knowledge/glossary/oscilloscopes/what-is-an-oscilloscope-trigger>.
  - [19] Steven B. Damelin and Willard Miller Jr. *The Mathematics of Signal Processing*. USA: Cambridge University Press, 2012. ISBN: 1107601045.
  - [20] FFTW3. *Fastest Fourier Transform in The West*. Accessed: (19/05/2025). 2025. URL: <https://www.fftw.org/>.
  - [21] Jeremy Fix. *FFTConvolution*. Accessed: (21/05/2025). 2013. URL: [https://github.com/jeremyfix/FFTConvolution/blob/master/Convolution/src/convolution\\_fftw.h](https://github.com/jeremyfix/FFTConvolution/blob/master/Convolution/src/convolution_fftw.h).
  - [22] *Planning-rigor flags*. Accessed: (23/05/2025). URL: <https://www.fftw.org/doc/Planner-Flags.html>.
  - [23] Rebar3. *Rebar3 Basic Usage*. Accessed: (25/05/2025). 2025. URL: [https://rebar3.org/docs/basic\\_usage/](https://rebar3.org/docs/basic_usage/).
  - [24] Erlang System Documentation. *Signals/Sending signals*. Accessed: (24/05/2025). 2025. URL: [https://www.erlang.org/doc/system/ref\\_man\\_processes](https://www.erlang.org/doc/system/ref_man_processes).
  - [25] ANTLR. *What is ANTLR4?* Accessed: (19/05/2025). 2025. URL: <https://www.antlr.org/>.
  - [26] Terence Parr and Kathleen Fisher. “LL(\*): the foundation of the ANTLR parser generator”. In: *Proceedings of the 32nd ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2011, San Jose, CA, USA, June 4-8, 2011*. Ed. by Mary W. Hall and David A. Padua. ACM, 2011, pp. 425–436. DOI: 10.1145/1993498.1993548. URL: <https://doi.org/10.1145/1993498.1993548>.
  - [27] lark-parser. *Lark - A parsing toolkit for Python*. Accessed: (24/05/2025). 2025. URL: <https://github.com/lark-parser/lark>.
  - [28] Wikipedia. *Qt (software) Wikipedia, The Free Encyclopedia*. Accessed: (24/05/2025). 2025. URL: [https://en.wikipedia.org/wiki/Qt\\_\(software\)](https://en.wikipedia.org/wiki/Qt_(software)).
  - [29] Tristan Sloughter. *ctx*. Accessed: (21/05/2025). 2023. URL: <https://github.com/tsloughter/ctx>.
  - [30] The Qt Company. *Add-ons available under Commercial Licenses, or GNU General Public License v3*. Accessed: (23/05/2025). 2025. URL: <https://doc.qt.io/qt-5/qtmodules.html>.



UNIVERSITÉ CATHOLIQUE DE LOUVAIN  
École polytechnique de Louvain

Rue Archimède, 1 bte L6.11.01, 1348 Louvain-la-Neuve, Belgique | [www.uclouvain.be/epl](http://www.uclouvain.be/epl)