



**NATIONAL AND KAPODISTRIAN UNIVERSITY OF ATHENS**

**FACULTY OF EXACT SCIENCES  
DEPARTMENT OF INFORMATICS AND TELECOMMUNICATIONS**

**BACHELOR THESIS**

**RHEA: A Reactive, Heterogeneous, Extensible and Abstract  
Framework for Dataflow Programming**

**Orestis Melkonian**

**Supervisors: Panos Rondogiannis, Professor EKPA  
Angelos Charalambidis, Researcher NCSR**

**ATHENS**

**APRIL 2016**

## **BACHELOR THESIS**

RHEA: A Reactive, Heterogeneous, Extensible and Abstract Framework for Dataflow  
Programming

**Orestis Melkonian**

**A.M.: 1115201000128**

**SUPERVISORS:** **Panos Rondogiannis**, Professor EKPA  
**Angelos Charalambidis**, Researcher NCSR

## ABSTRACT

The dataflow computational model enables writing highly parallel programs, which will be deployed on a heterogeneous network, in a concise and readable way. The main advantage is the fact that the system can be conceptually separated into several independent components that can be run in parallel and deployed on different machines. Therefore, concurrency and distribution is implicit and little or no responsibility is given to the programmer. The framework proposed in this thesis constitutes the underlying system that make this style of programming possible in JVM-based languages (e.g. Java, Scala, Closure), while at the same time making it easy to integrate other technologies that rely on the PubSub model, in order to move away from imperative languages and enter a higher level of abstraction. Particular emphasis was put on three domains, namely *Big Data*, *Robotics* and *IoT*. A detailed specification of the framework's characteristics and capabilities is given, followed by some robotic applications that demonstrate these attributes. Finally, related work and recent technologies are discussed.

**SUBJECT AREA:** Programming Languages

**KEYWORDS:** dataflow programming, frp, stream processing, distributed systems, declarative languages, implicit concurrency, node placement

*”τὰ ὄντα ἰέναι τε πάντα καὶ μένειν ουδέν”*  
*(all entities move and nothing remains still)*  
*- Heraclitus*

## **ACKNOWLEDGEMENTS**

I would like to thank Angelos Charalambidis for his immensely helpful supervision and guidance throughout the whole period of 6 months that I was present in NCSR.

I would also like to thank Professor Panos Randoiannis for being a major influence in my current research interests through the undergraduate courses "Theory of Computation" and "Principles of Programming Languages", where I was introduced to a much more declarative way of programming and started appreciating mathematical sensibility in computing.

# CONTENTS

<b>PROLOGUE</b>	<b>11</b>
<b>1. INTRODUCTION</b>	<b>12</b>
1.1 Main concept	12
1.2 Motivation	12
1.2.1 Declarative languages	12
1.2.2 Data versus Computation	12
1.2.3 Dataflows in Robotics	13
1.2.4 Dataflows in Big Data	13
1.3 Structure of thesis	13
<b>2. BACKGROUND</b>	<b>15</b>
2.1 The dataflow computational model	15
2.2 Functional reactive programming	16
2.3 Publish-Subscribe model	17
2.4 ROS: Robot Operating System	17
2.5 Internet of things - MQTT	18
2.6 The Reactive Streams Standard	18
<b>3. APPROACH</b>	<b>20</b>
3.1 Reactive	20
3.2 Heterogeneous	21
3.3 Extensible	21
3.4 Abstract	21
<b>4. IMPLEMENTATION</b>	<b>22</b>
4.1 The Strategy design pattern	22
4.2 Notifications	22
4.3 External Input-Output	22

<b>4.4</b>	<b>Software structure</b>	<b>23</b>
<b>4.5</b>	<b>Internal representation</b>	<b>23</b>
<b>4.6</b>	<b>Stream variables</b>	<b>23</b>
<b>4.7</b>	<b>Stream operators</b>	<b>24</b>
4.7.1	Creation	24
4.7.2	Combining	25
4.7.3	Filtering	25
4.7.4	Conditional	26
4.7.5	Transformational	27
4.7.6	Feedback	28
4.7.7	Error-handling	28
4.7.8	Backpressure	29
4.7.9	Utility	30
<b>4.8</b>	<b>Evaluation</b>	<b>31</b>
<b>4.9</b>	<b>Distribution</b>	<b>31</b>
4.9.1	Hazelcast	32
4.9.2	Machine configuration	32
<b>4.10</b>	<b>Optimization</b>	<b>32</b>
4.10.1	Proactive filtering	33
4.10.2	Granularity adjustment	33
4.10.3	Node placement	34
<b>4.11</b>	<b>Serialization</b>	<b>35</b>
<b>5.</b>	<b>APPLICATIONS</b>	<b>36</b>
<b>5.1</b>	<b>Hamming numbers</b>	<b>36</b>
<b>5.2</b>	<b>Camera surveillance</b>	<b>36</b>
<b>5.3</b>	<b>Robot control panel</b>	<b>37</b>
<b>5.4</b>	<b>Robot hospital guide</b>	<b>38</b>
<b>6.</b>	<b>RELATED WORK</b>	<b>41</b>
<b>6.1</b>	<b>Big Data</b>	<b>41</b>
6.1.1	GoogleDataflow	41
6.1.2	TensorFlow	42

6.1.3	Akka . . . . .	42
6.1.4	Spark . . . . .	43
6.1.5	Naiad . . . . .	44
6.1.6	dispel4py . . . . .	44
<b>6.2</b>	<b>Robotics . . . . .</b>	<b>45</b>
6.2.1	roshask . . . . .	45
6.2.2	Yampa . . . . .	45
6.2.3	Flowstone . . . . .	46
<b>6.3</b>	<b>Internet of Things . . . . .</b>	<b>47</b>
6.3.1	NoFlo . . . . .	47
6.3.2	Node-RED . . . . .	47
<b>7.</b>	<b>FUTURE WORK . . . . .</b>	<b>49</b>
7.1	More strategies . . . . .	49
7.2	Dynamic reconfiguration . . . . .	49
7.3	Advanced network profiling . . . . .	50
7.4	Advanced fault-tolerance . . . . .	50
7.5	Integration with other technologies . . . . .	51
7.6	Visual language . . . . .	51
7.7	Stream reasoning . . . . .	51
<b>8.</b>	<b>CONCLUSIONS . . . . .</b>	<b>53</b>
	<b>ABBREVIATIONS, INITIALS AND ACRONYMS . . . . .</b>	<b>54</b>



## LIST OF FIGURES

Figure 1:	Natural numbers . . . . .	15
Figure 2:	PubSub typical layout . . . . .	17
Figure 3:	Reactive properties . . . . .	20
Figure 4:	Heterogeneity pipeline . . . . .	21
Figure 5:	Strategy design pattern . . . . .	22
Figure 6:	Simple stream example . . . . .	24
Figure 7:	Split example . . . . .	24
Figure 8:	Partitioning . . . . .	32
Figure 9:	Optimization stages . . . . .	33
Figure 10:	Take/skip/distinct before map . . . . .	33
Figure 11:	Filter before map . . . . .	33
Figure 12:	Filter/distinct before concat/merge . . . . .	33
Figure 13:	Merge maps . . . . .	34
Figure 14:	Embed map in creation . . . . .	34
Figure 15:	Embed repeat in creation . . . . .	34
Figure 16:	Combine map with filter . . . . .	34
Figure 17:	Combine filter with exists . . . . .	34
Figure 18:	Combine map with exists . . . . .	34
Figure 19:	Combine map with zip . . . . .	34
Figure 20:	Combine zip with map . . . . .	34
Figure 21:	Meaningless nevers . . . . .	35
Figure 22:	Task fusion . . . . .	35
Figure 23:	Serialization process . . . . .	36

Figure 24:	Hamming numbers . . . . .	37
Figure 25:	Camera surveillance . . . . .	38
Figure 26:	Robot control panel . . . . .	39
Figure 27:	Robot hospital guide . . . . .	40
Figure 28:	TensorFlow graph . . . . .	42
Figure 29:	Akka DSL . . . . .	43
Figure 30:	Naiad application . . . . .	44
Figure 31:	Yampa operators . . . . .	46
Figure 32:	Flowstone screenshot . . . . .	46
Figure 33:	NoFlo graph editor . . . . .	47
Figure 34:	Node-RED graph editor . . . . .	48
Figure 35:	DyKnow architecture . . . . .	52
Figure 36:	The <i>RHEA</i> ecosystem . . . . .	53

## PROLOGUE

This bachelor thesis is the continuation of my internship at the National Centre for Scientific Research "Demokritos", particularly in the Software and Knowledge Engineering Laboratory (SKEL).

The main task I was assigned was the implementation of a framework for robot programming using a dataflow approach. During that internship, I came to realize that my work could be easily generalized to cover a much broader application area than just robot software.

The name of the framework stems from the ancient Greek Titaness Rhea(*Ρέα*), daughter of earth goddess Gaia(*Γαία*) and sky god Uranus(*Ουρανός*) and etymologically derives from the verb *ρέω*(to flow).

## 1. INTRODUCTION

### 1.1 Main concept

My main contribution is the design and implementation of a framework for dataflow programming to be deployed anywhere, ranging from low-performance robots and sensors to clusters of computer and even the Cloud.

The main idea is to provide the programmer with a different execution model, the dataflow model, which allows for a more abstract way of thinking and has the advantage of exposing opportunities for parallelism (amongst CPU cores) and distribution (amongst computational machines), which can then be automatically realised by the "intelligent" underlying system.

Therefore, the programmer will be able to utilize available computational resources without any effort, while at the same time reducing development time/cost and maintaining a much cleaner and easier-to-refactor software system. Resource utilization may appear in the form of faster execution or more robust error-handling.

Chain of Gain. Introduce running example?

### 1.2 Motivation

#### 1.2.1 Declarative languages

Software is becoming increasingly more complex each year, as computing capabilities are strengthened and user needs become more and more demanding. Thus the need for higher abstraction becomes imperative, as it provides a more structured, easier to debug and maintainable way of developing software. In other words, abstraction in computer science acts as a mean to overcome complexity.

In programming languages, the level of the aforementioned abstraction is measured regarding the amount of low-level details a programmer has to specify. Therefore, languages can be divided in two categories: the imperative ones, in which the programmer specifies what needs to be done and how to do it, and the declarative ones, where the programmer only specifies what needs to be done and rely on the underlying compiler/interpreter to produce the exact commands that will realize the desired behaviour. The most well-known declarative programming paradigms are functional and logic programming, each providing higher abstraction in different aspects. My approach was greatly influenced by the functional paradigm.

#### 1.2.2 Data versus Computation

A common problem in heterogeneous systems is that different representations of the same entities/data-types coexist in the same software and, as a consequence, pure computational tasks are intermingled with data-converting tasks. This makes the code less read-

able and harder to maintain and understand. In the dataflow execution model, where the program is modelled as directed graph of data flowing between operations, there is a clear separation of these two aspects as data (edges) are completely decoupled from computation (nodes). This motivation is strengthened even more, when cross-machine communication is included, and apart from converting data from one representation to another, serialization(i.e. conversion to bytes) is also mandatory.

### 1.2.3 Dataflows in Robotics

In control theory, which is the main background theory used in robotics, most architectures and/or algorithms are represented as dataflow diagrams for the sake of clarity and intuition. Translating these diagrams into common "imperative" software is not an easy task and is usually the source of bugs. Thus, having a dataflow execution model will nullify the need for such a translation.

Specifically, most robotic applications follow the *Robot Perception Architecture (RPA)*, where inputs to system are the robot's sensors, which are then processed by a dataflow graph, whose output is given as commands to the robot actuators.

Moreover, robotics typically involve several different robotic systems, whose combination is even more challenging. If each individual system is represented as a dataflow graph, composing them together is as trivial as connecting inputs with outputs, which is not the case in a traditional architecture, which is not component-based.

### 1.2.4 Dataflows in Big Data

Another reason for following a dataflow approach is the attention that it recently has drawn in the *Big Data* field. As data size is growing exponentially and distribution is not a luxury but a necessity, a more scalable and decentralized architecture was destined to be examined in more depth. As we will discuss in the *Related Work* chapter, there are many recent frameworks that became famous for their scalability due to the fact that they rely on a dataflow approach.

## 1.3 Structure of thesis

There are nine chapters which compose this thesis: *Introduction, Background, Approach, Implementation, Applications, Related Work, Future Work, Conclusions*.

*Background* introduces the reader to basic background knowledge, necessary for complete comprehension.

*Approach* presents the main characteristics of my approach.

*Implementation* gives a more detailed specification of the framework.

*Applications* present some use-cases, ranging from general mathematical problems to

real-life robot scenarios.

*Related Work* discusses relevant concepts and technologies, which influenced major decisions concerning the design and implementation of the framework.

*Future Work* suggests some interesting topics for future research, whose embedding in the framework is meaningful.

*Conclusions* sums up.

## 2. BACKGROUND

### 2.1 The dataflow computational model

The increased interest in parallelism during the 70's gave rise to the dataflow execution model, which is an alternative to the classical "von-Neumann" model. In the dataflow model, everything is represented in a dataflow graph, where nodes are independent computational units and edges are communication channels between these units. A node/unit is fired immediately when its required inputs are available and therefore no explicit control commands are needed for execution. Figure 1 shows a dataflow graph enumerating the set  $\mathbb{N}$  of natural numbers.

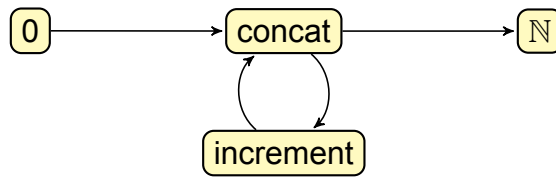


Figure 1: Natural numbers

In the dataflow graph above, we can discern three types of nodes: sources, which do not have any incoming edge and act as value generators to initiate computation, sinks, which do not have any outgoing edges and inner nodes, which transform one or more incoming streams and redirect their output to other nodes. The *zero* node just produces a stream with a single value 0 and then terminates. *Concat* produces a single stream by concatenating the stream produced by *zero* and *increment*, while *increment* transforms its input stream by adding one to its values. Finally, the sink node displays the result, which is the stream of natural numbers.

Streams can be infinite, such as the stream produced by *concat* because it is the concatenation of a single-value stream and an infinite one. Moreover, the graph is cyclic as *concat* feeds input to *increment* and vice versa. The most interesting fact is that there nodes are independent and therefore can run in parallel. For instance, while *increment* is processing value 5 (i.e. to produce value 6), the previous result (i.e. value 5) passes through *concat* to reach the sink node, which can concurrently process it to display it.

The main advantage of the dataflow model is its implicit parallelism, deriving from the fact that the computational units are totally independent and therefore can be executed in parallel. A possible single-machine implementation could represent edges as in-memory data storage, whereas a multi-machine one could represent them as channels between TCP sockets, allowing communication across the network. Its great flexibility and composability makes it a good candidate for the underlying architecture of a framework with a high level of abstraction.

## 2.2 Functional reactive programming

A relatively recent programming paradigm is *Functional Reactive Programming (FRP)*, which provides a conceptual framework for implementing reactive (i.e. time-varying and responding to external stimuli) behaviour in *hybrid systems* (i.e. containing both continuous and discrete components), such as robots, in functional programming languages.

To implement such systems in conventional imperative languages, one must use asynchronous *callbacks* (i.e. each change is handled by a registered *callback* function). Although this solution is satisfactory for simple schemes, more complex scenarios eventually lead to highly incoherent code structure, often called *spaghetti code*, in the sense that control rapidly moves between disconnected parts of the system, similar to the notorious *GOTO* command. This phenomenon stems from the unary nature of *callback* functions, which requires some kind of "internal plumbing" in order to achieve mechanisms for handling combination of changes (e.g. when multiple changes occur simultaneously). *FRP* provides a solution to this shortcoming of *callback* functions, because changes are represented as variables (*signals*), which can be passed as parameters to arbitrary functions, called *signal functions*.

*FRP* first appeared as a composable library for graphic animations [1], but quickly evolved into a generic paradigm [2, 3, 4]. Moreover, extensive research has investigated *FRP* as a framework for robotics [5, 6].

Although appealing at first, *FRP* was not appropriate for systems with real-time constraints, due to uncontrollable time- and space- leaks [7]. The solution was a generalization of monads called *arrows* [8], which provided the necessary guarantees that the aforementioned common errors do not occur. Let's see the example of calculating a robot's x-coordinate. Here is the mathematical formula drawn from control theory:

$$x = 1/2 \int (vr + vl) \cos \theta$$

Below is the *FRP* code corresponding to the formula above:

---

```

1 x = let
2   v = (vrSF &&& vlSF) >>> lift (+)
3   t = thetaSF >>> arr cos
4   in (v &&& t) >>> lift (*) >>> integral >>> lift (/2)

```

---

As the above may seem counter-intuitive and difficult to understand, a new notation was conceived, notably the *arrow notation* [9]:

---



```

1 x = proc inp -> do
2   vr <- vrSF -< inp
3   vl <- vlSF -< inp
4   theta <- thetaSF -< inp
5   i <- integral-< (vr+vl) * cos(theta)
6   returnA -< (i/2)

```

The main advantages of *FRP* are its close correspondence to mathematics[2], which make it an ideal framework for modelling real-time systems, and its concise representation of time-varying values via *signals*.

### 2.3 Publish-Subscribe model

*Publish/Subscribe (PubSub)* is a messaging pattern that became popular due to the loose coupling of its components, suited for the most recent large-scale distributed applications.

There is no point-to-point communication and no synchronization. *Publishers* advertise messages of a given type to a specific message class or *topic* that is identified by a *key-word*, whereas *subscribers* listen on a specific *topic* without any knowledge of who the publishers are. The component responsible for relaying the messages between machines and/or processes and finding the cheaper dissemination method is called the *message broker*. Figure 2 illustrates an abstract representation of the *PubSub* model.

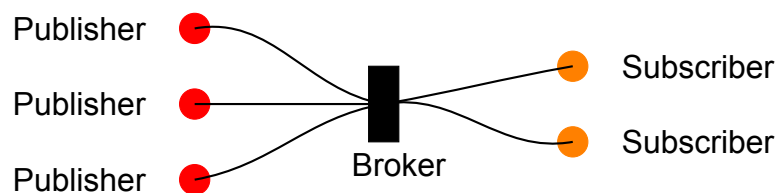


Figure 2: PubSub typical layout

### 2.4 ROS: Robot Operating System

*ROS* is an open-source middleware for robot software, which emphasizes large-scale integrative robotics research [10]. It provides a *thin* communication layer between heterogeneous computers, from robots to mainframes and it has been widely adopted by the research community around the world, due to its flexibility and maximal support of reusability through packaging and composability. It provides a compact solution to the development complexity introduced by complex robot applications that consist of several modules and require different device drivers for each individual robot.

It follows a peer-to-peer network topology, implemented using a topic-based *PubSub* messaging protocol and its architecture reflects many sound design principles. Another great property of *ROS* is that it is language-agnostic, meaning that only a minimal specification language for message transaction has been defined, so contributors are free to imple-

ment small-size clients in different programming languages, with *roscpp* and *rospy* being the most widely used ones.

A typical development scenario is to write several *nodes*, that subscribe to some topics and, after doing some computation, publish their results on other topics. The main architectural issue here is that subscribing is realized through asynchronous callback functions, so complicated schemes easily lead to unstructured code, which obviously lead to unreadable and hard-to-maintain code. My approach gives a solution to the aforementioned problem.

## 2.5 Internet of things - MQTT

The birth of the Internet gave rise to a concept called *Internet of Things (IoT)*, which is essentially the ability of many heterogeneous devices, ranging from low-cost sensors to vehicles with embedded electronics, to collect data and exchange it amongst themselves using the Internet. This gave rise to smart grids, smart homes and eventually smart cities.

The development of such systems though, due to their heterogeneity, is rather complex and costly. Typical software architectures were not meant to be used in such environments and therefore new tools and concepts needed to be invented. Recent development of a variety of middleware frameworks, showed that a standard protocol of communication is imperative along with supporting tools[?]. The most widely spread protocol is *MQTT*, which follows the *PubSub* messaging pattern and provides a very minimal communication layer in order not to put a strain on the resource-bounded system[?].

For instance, an *IoT* application could connect to some sensors by subscribing to their corresponding topics, taking decisions that would result in some commands to some actuators, by publishing to their corresponding topics.

Fortunately, the dataflow model seems to be rather fitting for these scenarios[?], as every node in the graph is completely independent, and consequently can be any "*thing*". This useful property of the model makes it a good architectural choice for such applications. The only thing to consider is how these things will communicate in a standard way, so as to be able to add new types of *things* and integrate it in an effortless way to an existing dataflow network.

## 2.6 The Reactive Streams Standard

*PubSub* is widely used by different frameworks but still lacks standardization. The *Reactive Streams Standard (RSS)* is an initiative to provide a standard for asynchronous stream processing with non-blocking back pressure. This encompasses efforts aimed at runtime environments (JVM and JavaScript) as well as network protocols [13].

*RSS* defines two minimal interfaces for the roles of *Subscriber*<sup>1</sup> and *Publisher*<sup>2</sup>. A *Subscriber* implementation should define reactions to observed values, including normal and erroneous termination, whereas a *Publisher* implementation should accept requests from *Subscribers* and start emitting values to them.

Below we see a minimal example of using *RSS* to define a publisher that emits values 1..10 and a subscriber that prints all observed values and finally connect them together.

---

```
1 Publisher<Integer> pub = new Publisher {
2     void subscribe(Subscriber<Integer> sub) {
3         for (int i = 1; i <= 10; i++)
4             sub.onNext(i);
5         sub.onComplete();
6     }
7 };
8
9 Subscriber<Integer> sub = new Subscriber {
10     void onNext(Integer i) {
11         println(i);
12     }
13     void onComplete() {
14         println("Complete");
15     }
16     void onError(Throwable t) {
17         t.printStackTrace();
18     }
19 };
20
21 pub.subscribe(sub);
```

---

<sup>1</sup><http://www.reactive-streams.org/reactive-streams-1.0.0-javadoc/org/reactivestreams/Subscriber.html>

<sup>2</sup><http://www.reactive-streams.org/reactive-streams-1.0.0-javadoc/org/reactivestreams/Publisher.html>

### 3. APPROACH

The design was heavily influenced by principles set out by the FRP and dataflow models.

#### 3.1 Reactive

The system is *reactive*, as close as possible to the definition of the Reactive Manifesto [12].

The system is *responsive*, meaning it should be able to handle time-sensitive scenarios if at all possible. This is the cornerstone of usability and utility, but more than that, it enables quick error-detection and error-handling.

The system is *resilient*, meaning it is able to recover robustly and gracefully after a failure, due to the fact that nodes in the dataflow graph are completely independent and recovery of each one can be done in isolation. Another thing to note here is that special error messages are built-in and make it very easy to propagate errors between *components*, in case the error-handling part of a component is decoupled from the computational logic. This leads to much more robust architectures for large-scale systems, where fault-tolerance is mission-critical.

The system is *elastic*, meaning it will adjust itself depending on the available resources and demanded workload. For instance, the granularity of the graph (i.e. number of nodes) is adjusted so as to match a heuristic-based value (e.g. total number of threads).

The system is *message-driven*, meaning it relies solely on asynchronous message-passing for inter-component communication leading to loose coupling, isolation, location transparency and the error propagation mentioned above. Location transparency is critical to preserve the semantics whether on a single host or a machine cluster.

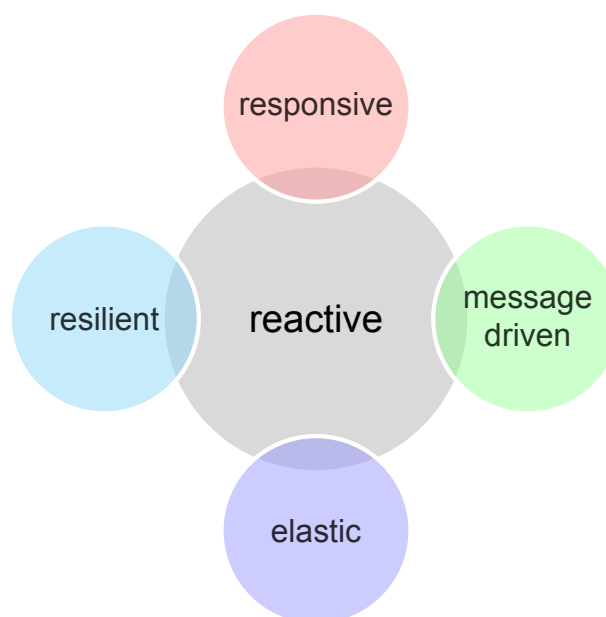


Figure 3: Reactive properties

### 3.2 Heterogeneous

One of the major concerns while designing the framework was the ability to deploy it anywhere, from low-cost robots to mainframes. Obviously, such attribute would require a very flexible runtime environment. To satisfy this requirement, the strategy design pattern was used for evaluation, meaning that the core system only builds the internal representation of the dataflow graph and partitions it across the available computational resources. From there onwards, each partial graph can be evaluated by a different *EvaluationStrategy* (see Implementation chapter), which could interpret it using the Java 8 Streams library or even compile into CUDA code for execution on a GPU.

Figure 5 illustrates a simple example of a robot application pipeline, where input to the dataflow graph is what the robot's camera senses and, after some image processing and some computation-heavy decision making, a command to an actuator of the robot is executed. Orange nodes are deployed on the robot's on-board computer, the green node is deployed on an off-board GPU and the red node is deployed on the main server.

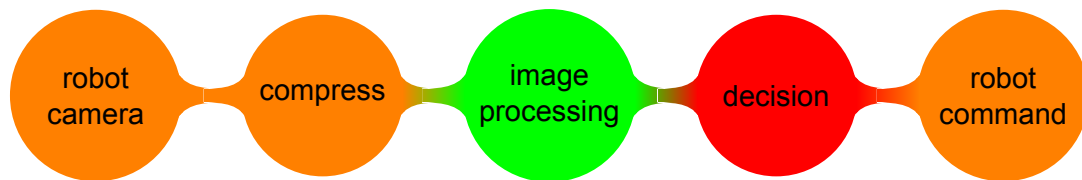


Figure 4: Heterogeneity pipeline

### 3.3 Extensible

While doing my internship at *NCSR*, I realized that the goal I was pursuing was impossible to reach closure in a 6-month period by a mere undergraduate student. Therefore, I decided that everything should be implemented in a way that will allow future contribution.

With that concept in mind, I started generalizing and abstracting away everything I had done so far with the hope that the framework will raise attention later on. I can now say I am satisfied with the level of abstraction the core system has reached and I hope the stressful refactoring that the framework went through will blossom in the form of future contributions.

### 3.4 Abstract

The framework is *abstract* in terms of implementation details, as it is completely agnostic of any machine-specific requirements. It is designed as a unifying conceptual base for further extensions and careful consideration was taken not to restrict in any aspect, architectural or not. This was achieved by making many parts of the core system pluggable, allowing for easy refactoring on most of its internal functionality. Moreover, the internal graph representation does not include information on how a node is executed, but only on its semantics.

## 4. IMPLEMENTATION

This section examines the major characteristics of the framework's implementation.

### 4.1 The Strategy design pattern

In software engineering, and especially in *object-oriented programming (OOP)*, a *design pattern* is a general repeatable solution to a commonly occurring problem in software design [11].

One such solution, the *Strategy* design pattern is used when a particular algorithm can be implemented by a variety of behaviours/classes [11]. In such a case, a good idea is to isolate the algorithm in a separate *interface* and allow the system to select the appropriate instantiating classes at runtime.

Figure 3 illustrates the basic UML diagram of the strategy design pattern.

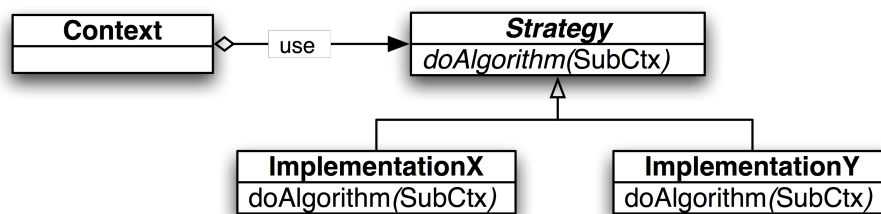


Figure 5: Strategy design pattern

### 4.2 Notifications

Every value passed through the framework's *streams* is wrapped inside a *Notification* object, which discriminated stream values into three categories: *onNext* (when the stream provides a regular value), *onError* (when an error occurs) and *onComplete* (when the stream completes its output).

### 4.3 External Input-Output

In order to make the framework easy to integrate with other stream and/or dataflow technologies, I decided that every input/output node (i.e. publisher/subscriber in the *PubSub* terminology or source/sink in the dataflow terminology) should implement the interfaces that *RSS* defines.

A sink node (output) should implement the Subscriber interface<sup>3</sup>, which essentially defines three methods corresponding to reactions to a *Notification*, one for each of the categories mentioned above.

A source node (input) should implement the Publisher interface<sup>4</sup>, which defines a single method *subscribe(Subscriber)*, where a Subscriber requests the Publisher to start emitting

<sup>3</sup><http://www.reactive-streams.org/reactive-streams-1.0.0-javadoc/org/reactivestreams/Subscriber.html>

<sup>4</sup><http://www.reactive-streams.org/reactive-streams-1.0.0-javadoc/org/reactivestreams/Publisher.html>

values.

Many existing technologies provide these interfaces, or at least adapters from their internal representations, and therefore they are very easy to be integrated to the framework.

#### 4.4 Software structure

The core system<sup>5</sup> is organized in the following top-level packages:

<i>org.rhea_core.internal</i>	all internal functionalities such as graphs, expressions and notifications
<i>org.rhea_core.evaluation</i>	everything associated with the evaluation of the constructed dataflow graph
<i>org.rhea_core.distribution</i>	everything associated with distributing the evaluation across the available computational resources
<i>org.rhea_core.optimization</i>	includes some built-in optimizers for adjusting the granularity of the graph
<i>org.rhea_core.network</i>	everything associated with networking
<i>org.rhea_core.io</i>	defines the interfaces that sources/sinks should implement
<i>org.rhea_core.util</i>	helpful utilities needed throughout the project

#### 4.5 Internal representation

For representing the internal structure of the dataflow graph, the *JGrapht* open-source Java library was used, which provides many graph data structures and common graph-theory algorithms<sup>6</sup>. The main class representing the internal dataflow graph is *FlowGraph*, which is located in the *org.rhea\_core.internal* package.

#### 4.6 Stream variables

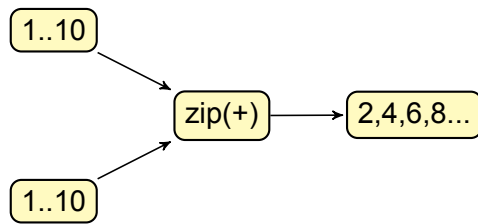
The construction of the aforementioned internal graph is always implicit, through a rich set of stream operators. The only data type handled by the programmer is *org.rhea\_core.Stream*, which contains the parametric type of its single output. Each *Stream* object contains internally a *FlowGraph*, which is only to be accessed and manipulated by the internal module, evaluation strategies and optimizers.

Figure 6 illustrates the dataflow graph(left) produced by the framework code(right) using the *Stream* data type.

These variables can be used, preferably together with their parametric type, and reused in different parts of the graph. This is necessary, for instance, when you wish to split a node's output to different inputs. Figure 7 shows such an example.

<sup>5</sup><https://github.com/rhea-flow/rhea-core>

<sup>6</sup><http://jgrapht.org/>



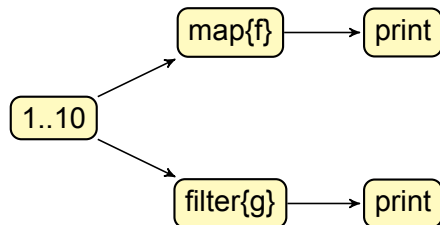

---

```

1 Stream.zip(
2     Stream.range(1, 10),
3     Stream.range(1, 10),
4     (x, y) -> x + y);
  
```

---

Figure 6: Simple stream example




---

```

1 Stream<Int> st = Stream.range(1, 10);
2 st.map(f).print();
3 st.filter(g).print();
  
```

---

Figure 7: Split example

## 4.7 Stream operators

This section displays the set of primitive operators, from which all the available stream operators are derived. Type information is not shown for the sake of readability, but all operations are type-safe.

In the marble diagrams on the right, circles represent *onNext* notifications, green bars represent *onComplete* and red x signals an *onError*.

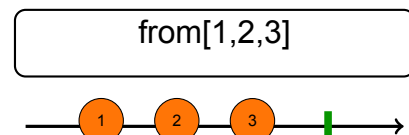
### 4.7.1 Creation

The following operators act as source nodes in the dataflow graph.

**from** (Iterable i)

*inputs:* none

*output:* i as a stream



**fromSource** (Source s)

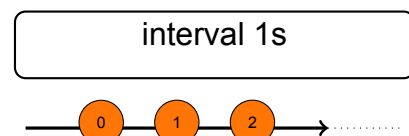
*inputs:* none

*output:* the values emitted by s

**interval** (TimeInterval t)

*inputs:* none

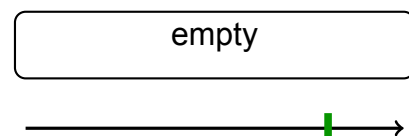
*output:* the natural numbers emitted every t



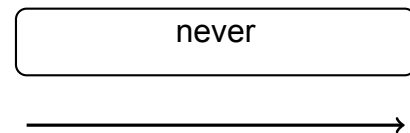
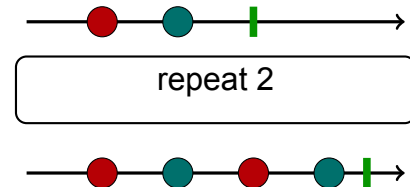
**empty**

*inputs:* none

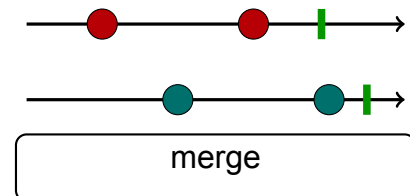
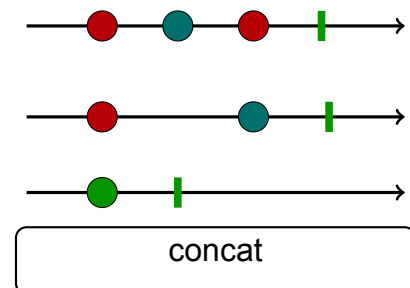
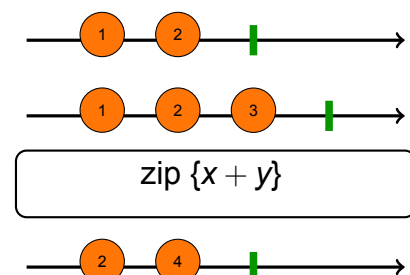
*output:* an empty stream



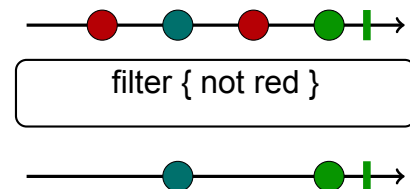
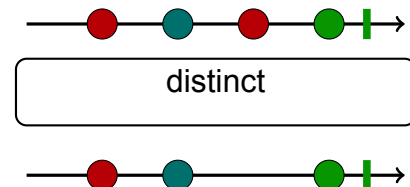
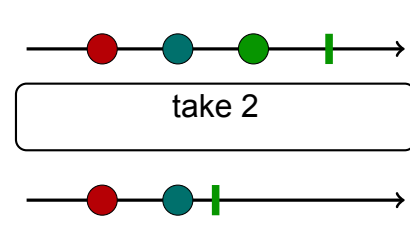
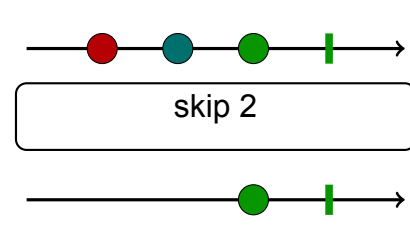


**never***inputs:* none*output:* a stream that emits no notification**repeat** (int *n*)*inputs:* single*output:* repeats the values of the input stream *n* times, or infinitely if *n* < 0**defer** (*f* : () → *Stream*)*inputs:* none*output:* the stream generated by the given stream factory *f***4.7.2 Combining**

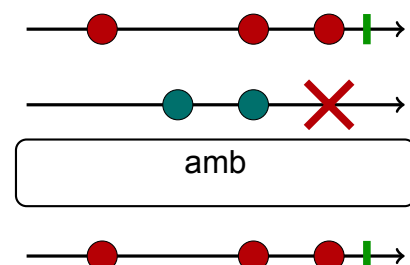
The following operators combine multiple streams to produce another one.

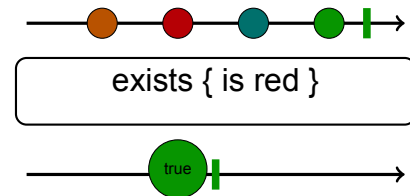
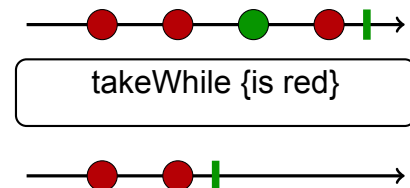
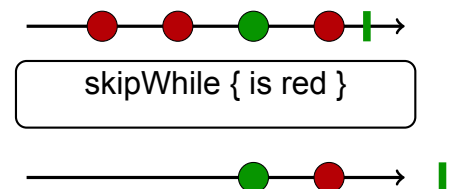
**merge***inputs:* many*output:* merges the input streams and completes as soon as any of them completes**concat***inputs:* many*output:* concatenates the input streams in the given order**zip** (*f* : *A*<sub>1</sub>...*A*<sub>*n*</sub> → *B*)*inputs:* many*output:* zips the input streams with the given function *f***4.7.3 Filtering**

The following operators filter the values emitted by another stream.

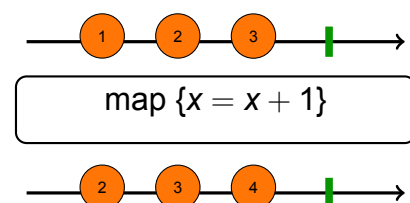
**filter** ( $f : A \rightarrow \text{Boolean}$ )*inputs*: single*output*: emits only values  $i$  of the input stream, where  $f(i) = \text{true}$ **filterMap** ( $f : A \rightarrow B, g : B \rightarrow \text{Boolean}$ )*inputs*: single*output*: emits only values  $f(i)$  of the input stream the, where  $(f \circ g)(i) = \text{true}$ **distinct***inputs*: single*output*: removes all duplicate values of the input stream**take** (int  $n$ )*inputs*: single*output*: extracts the first (last)  $n$  values of the input stream, if  $n > 0 (< 0)$ **skip** (int  $n$ )*inputs*: single*output*: skips the first (last)  $n$  values of the input stream, if  $n > 0 (< 0)$ **4.7.4 Conditional**

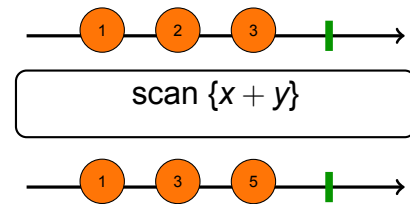
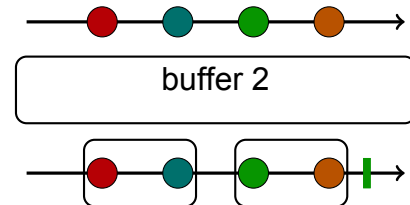
The following operators behave depending on some conditions on their input stream

**amb***inputs*: many*output*: emits the values of one its input streams, whichever emits a value or terminates first

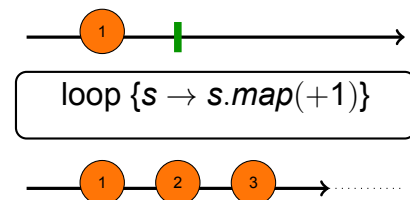
**exists** ( $f : A \rightarrow \text{Boolean}$ )*inputs*: single*output*: returns True, if the input stream contains a value  $i$ , where  $f(i) = \text{true}$  and False otherwise**takeUntil** (Stream  $s$ )*inputs*: single*output*: extracts values of the input stream, until stream  $s$  emits a value**skipUntil** (Stream  $s$ )*inputs*: single*output*: skips values of the input stream, until stream  $s$  emits a value**takeWhile** ( $f : A \rightarrow \text{Boolean}$ )*inputs*: single*output*: extracts values of the input stream until a value  $i$  is emitted, where  $f(i) = \text{false}$ **skipWhile** ( $f : A \rightarrow \text{Boolean}$ )*inputs*: single*output*: skips values of the input stream until a value  $i$  is emitted, where  $f(i) = \text{false}$ **4.7.5 Transformational**

The following operators transform their input stream

**map** ( $f : A \rightarrow B$ )*inputs*: single*output*: transforms the input stream by applying function  $f$  to every value emitted

**scan** ( $B$  seed,  $f : B \times A \rightarrow B$ )*inputs:* single*output:* transforms the input stream by sequentially applying function  $f$  to every value emitted and emitting each result along the way**buffer** (int  $n$ )*inputs:* single*output:* packs together every  $n$  values of the input stream into a single List item**buffer** (TimeInterval  $t$ )*inputs:* single*output:* packs together values of the input stream emitted every  $t$  into a single List item**4.7.6 Feedback**

This operator enables cycles in the dataflow graph.

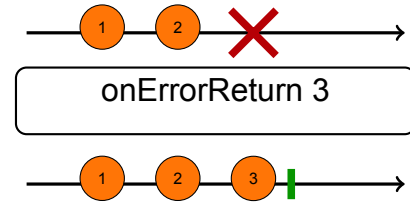
**loop** ( $f : Stream \rightarrow Stream$ )*inputs:* single*output:* attaches a sub-graph to the input stream's output, whose result acts as feedback to the attachment point**4.7.7 Error-handling**

The following operators are a mean to handle errors.

**onErrorResume** (Stream  $s$ )*inputs:* single*output:* mirrors the input stream, but instead of emitting an `onError` Notification when an error occurs, continues emitting values of the given stream  $s$

**onErrorReturn** ( $f : Throwable \rightarrow A$ )*inputs*: single

*output*: mirrors the input stream, but instead of emitting an *onError* Notification when an error *e* occurs, emits the value  $f(e)$  followed by a *onComplete* Notification

**retry** (int *n*)*inputs*: single

*output*: mirrors the input stream, but instead of emitting an *onError* Notification when an error occurs, resubscribes to it *n* times if  $n > 0$ , infinitely otherwise

**4.7.8 Backpressure**

The following operators specify how a node behaves when the requests are too intense to handle, computationally or memory-wise. *Backpressure* is the mechanism that handles fast publishers that interact with slow subscribers.

**onBackpressureBuffer***inputs*: single

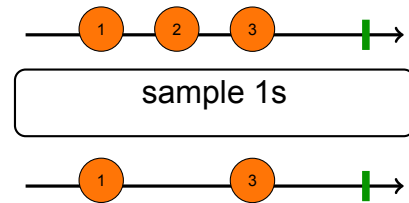
*output*: buffers values that cannot be handled by the subscriber to emit them later on

**onBackpressureDrop***inputs*: single

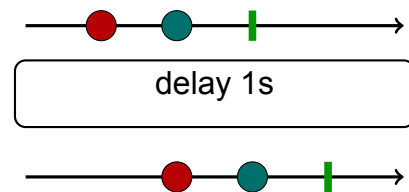
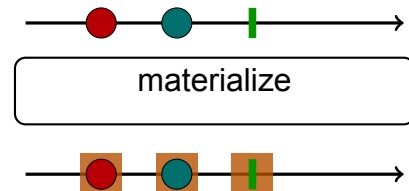
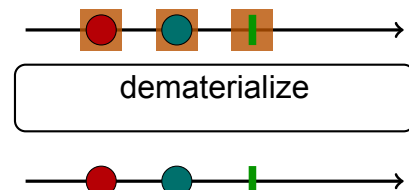
*output*: drops values that cannot be handled by the subscriber, instead of emitting them

**onBackpressureLatest***inputs*: single

*output*: drops values that cannot be handled by the subscriber and always buffers the last one, instead of emitting them

**sample** (TimeInterval *t*)*inputs*: single*output*: emits only the most-recent emitted value from the input stream within intervals of *t***timeout** (TimeInterval *t*)*inputs*: single*output*: mirrors the input stream, but emits *onError* if there is no emission within windows of *t***4.7.9 Utility**

The following operators provide some helpful utilities.

**doOnNext** (Action *a*)*inputs*: single*output*: execute Action *a* whenever `onNext(Complete/Error)` is called**cache***inputs*: single*output*: caches values emitted by the input stream for future subscribers**delay** (TimeInterval *t*)*inputs*: single*output*: emits the values of the input stream shifted forward in time by *t***materialize***inputs*: single*output*: wraps all values of input stream as Notifications**dematerialize***inputs*: single*output*: reverses the effect of **materialize**

All other operators can be produced by combining the above primitive ones (e.g. *flatMap*  $\equiv$  *map*  $\circ$  *merge*).

## 4.8 Evaluation

Every primitive operator corresponds to an expression implementing the *Transformer* interface, defined in the *org.rhea\_core.internal.expressions* package.

A complete dataflow is defined by a *Stream* variable and an object implementing the *Output* interface, which can be either an *Action*, a *Sink* or a list of these.

In order to evaluate a constructed dataflow graph the strategy design pattern is used, therefore a class implementing the *EvaluationStrategy* interface, found in the *org.rhea\_core.evaluation* package, needs to be provided. An *EvaluationStrategy* just takes the *Stream* variable and its corresponding *Output* and executes it, however desired.

The strategies I have implemented so far follow:

### RxJavaEvaluationStrategy<sup>7</sup>

Uses rxjava<sup>8</sup>, which is a famous and well-maintained library for asynchronous programming using the *Observable* type, which is very close, semantically, to my *Stream* type.

### RosEvaluationStrategy<sup>9</sup>

Integrates the ROS middleware into the framework, by providing the *RosTopic* class, which implements the *AbstractTopic* interface defined in the *org.rhea\_core.io* package. This strategy's job is to set up a ROS client and configure every *RosTopic* used within the dataflow that needs to be evaluated to use this client. After that, evaluation is propagated to a generic strategy (e.g. rxjava).

### MqttEvaluationStrategy<sup>10</sup>

Integrates the MQTT middleware into the framework, in the same way ROS is integrated.

## 4.9 Distribution

An evaluation strategy executes the requested dataflow graph in a single machine, without concern about distribution and resource utilization.

For distribution and cluster management, the strategy design pattern is used again, specifically the *DistributionStrategy* interface, which is defined in the *org.rhea\_core.internal.distribution* package. Its responsibility is to take the whole initial graph that we need to evaluate and, after adjusting its granularity (i.e. size) to fit the available resources (see *Optimization*

<sup>7</sup><https://github.com/rhea-flow/rx-eval>

<sup>8</sup><https://github.com/ReactiveX/RxJava>

<sup>9</sup><https://github.com/rhea-flow/ros-eval>

<sup>10</sup><https://github.com/rhea-flow/mqtt-eval>

section), partition it across all computational resources, maybe using different evaluation strategies.

#### 4.9.1 Hazelcast

Due to the RSS being in its infant stage, no working implementation exists for superimposing a network protocol onto it (e.g. RSS over TCP). For this reason, I relied upon the open-source *Hazelcast* library<sup>11</sup> to discover and manage multiple machines and used its internal decentralized *PubSub* model to communicate intermediate results across the network. Figure 8 illustrates a dataflow graph on the left and the same graph partitioned over several machines on the right.

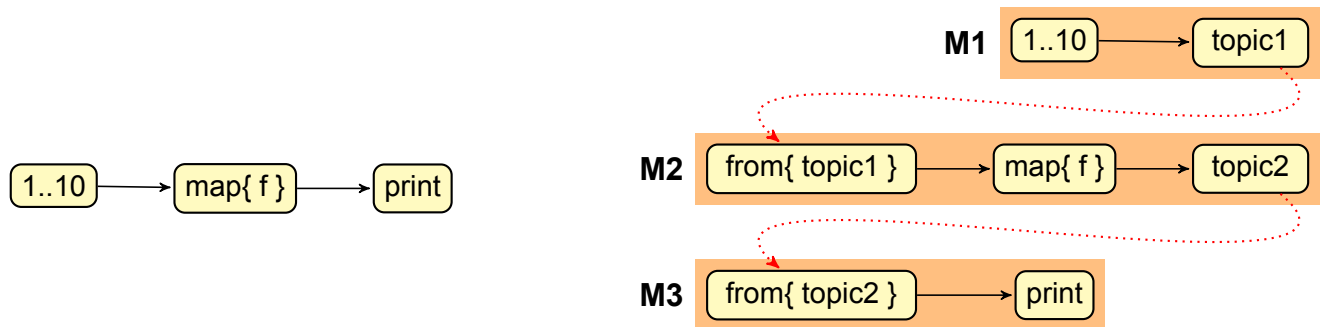


Figure 8: Partitioning

#### 4.9.2 Machine configuration

According to the distribution strategy being used, the available machines will require a certain initial configuration. For the *Hazelcast* case, a little piece of setup code needs to be executed on every member of the cluster, which is together with the main *Strategy* class. Moreover, helpful information can also be added at this step, such as number of CPU cores. It is the distribution strategy's responsibility to ensure that this information is properly distributed and handled.

Apart from this initial configuration, the distribution strategy needs to enable members to declare certain skills that they possess, which are required by specialized nodes. For instance, a source node emitting values from a *ROS* topic must be executed on a machine having *ROS* installed, in order to set up a *ROS* client. In the *Hazelcast* case, these skills are just *strings* and are declared in the initialization code of each machine separately.

#### 4.10 Optimization

This section describes three stages of optimization the dataflow graph goes through before being evaluated. To aid extensibility the strategy design pattern is again used, whose corresponding interface *OptimizationStrategy* resides in the *org.rhea\_core.optimization*

<sup>11</sup><http://hazelcast.org/>



package. Figure 9 illustrates the optimization stages.

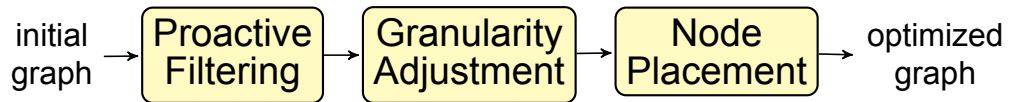


Figure 9: Optimization stages

#### 4.10.1 Proactive filtering

The first optimization stage is a heuristic one, based on the fact that if a filter operation can be moved earlier (i.e. closer to source nodes) while preserving the original semantics, then there will be benefit concerning computational cost and cross-machine communication overhead. The figures below show the corresponding graph transformations.



Figure 10: Take/skip/distinct before map



Figure 11: Filter before map

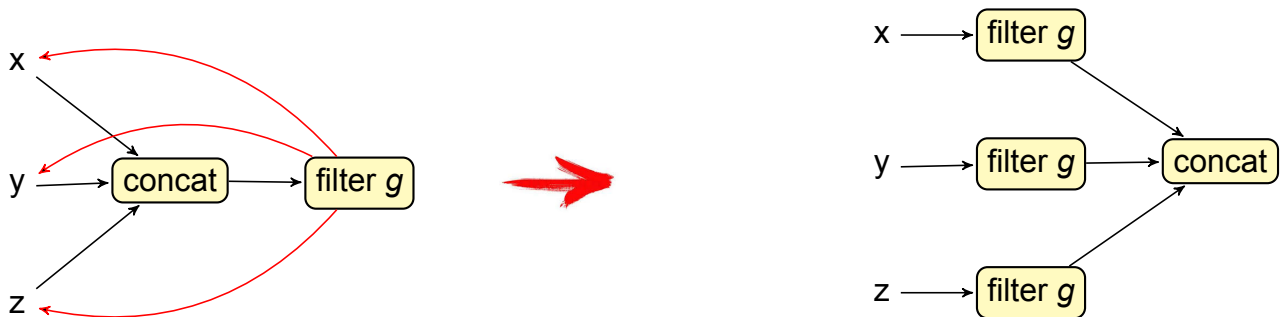


Figure 12: Filter/distinct before concat/merge

#### 4.10.2 Granularity adjustment

Different nodes of the dataflow graph will be executed on a separate thread/process. The fact that graphs can grow very big poses a problem when available computational resources are limited. For this reason, the second optimization stage tries to adjust the granularity of the dataflow graph to a desired value, which is the number of available threads amongst all machines.

To achieve this, the optimizer applies some semantic-preserving transformation, as shown in the figures below (for simplicity, only a single example of each general case is demonstrated).



Figure 13: Merge maps



Figure 14: Embed map in creation



Figure 15: Embed repeat in creation



Figure 16: Combine map with filter



Figure 17: Combine filter with exists



Figure 18: Combine map with exists

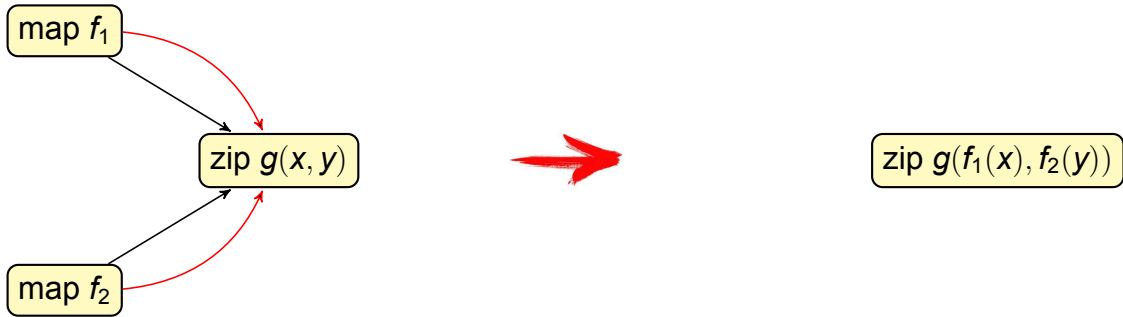


Figure 19: Combine map with zip



Figure 20: Combine zip with map

#### 4.10.3 Node placement

After the first two passes, we have an optimized dataflow graph with fine-tuned granularity. At this stage, nodes are mapped to tasks and are deployed across the available machines.

If the desired granularity has not been reached yet, the *DistributionStrategy* applies fusion to pairs of tasks until it reaches it, as shown in Figure 22.

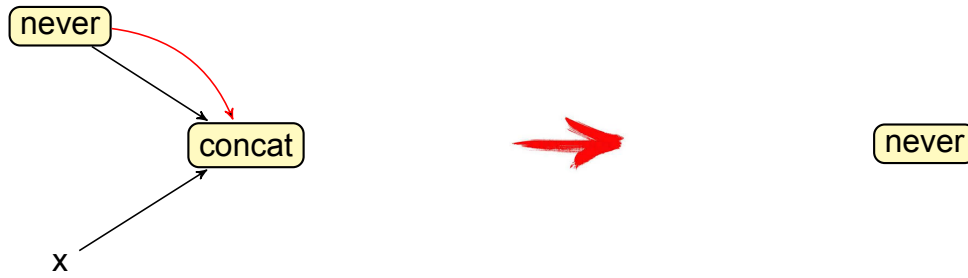


Figure 21: Meaningless nevers



Figure 22: Task fusion

The final decision to be made is where each of these newly constructed tasks will be executed, although some of them need to necessarily be placed on specific machines with certain skills.

Apart from these hard constraints, we need to minimize communication overhead. For this purpose, the strategy design pattern is again used, namely the *NetworkProfileStrategy* that is defined in the *org.rhea\_core.distribution* package. Its responsibility is to calculate a network distance between each pair of available machines, which is fed as input to the *NodePlacement* optimizer.

At this stage, we identify the aforementioned network distance as cost and apply brute-force to find the optimal placement of the (groups of) tasks that minimize that cost.

#### 4.11 Serialization

As communication between machines across a network is mandatory, data types emitted through the streams must be serialized on departure and de-serialized on arrival at each machine. For this reason, each *DistributionStrategy* must be configured with a class implementing the *Serializer* interface, define in the *org.rhea\_core.serialization* package. The byte representation of the objects is parametric for maximum flexibility.

A default *Serializer* is provided with the core system, which can serialize every class implementing the *Serializable* interface. In addition to that, the *JsonIO* library<sup>12</sup> is used which allows serialization of many types of classes, but still does not cover every possible one. Figure 23 depicts the serialization process in more detail.

<sup>12</sup><https://github.com/jdereg/json-io>

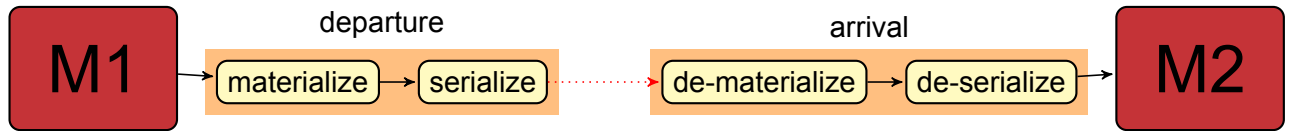


Figure 23: Serialization process

## 5. APPLICATIONS

This chapter will demonstrate some use-cases of the framework, which can be found in the framework repository<sup>13</sup>. At first, a mathematical example is shown, which is followed by three robotic applications.

### 5.1 Hamming numbers

Consider the problem of enumerating the *Hamming numbers*, which are generated by the mathematical formula  $\mathbb{H} = 2^i 3^j 5^k$ , where  $i, j, k \in \mathbb{N}$ . There is an intuitive dataflow solution to the above problem, taken from the book of *Lucid*, which is the first functional dataflow language[?]. The figure below shows the dataflow graph on the left and the corresponding *RHEA* code on the right.

The code is written in Scala to utilize the *Pimp my library* design pattern[?], which is used to easily add new functions to already existing libraries, using Scala's *implicit conversions* (line 28). In the example above, we define two new Stream operators, namely *multiply* (line 10), which just multiplies the stream with a constant, and *mergeSort* (line 13), which produces an ordered stream given two ordered streams as input. We also see the power of the *loop* operator (line 2), which allows us to define cycles in an effortless manner.

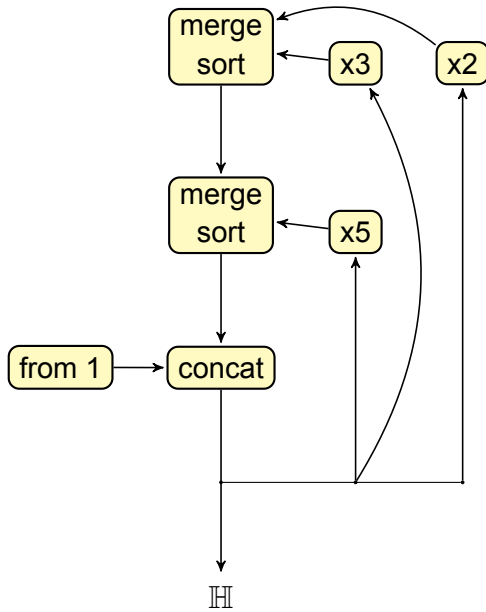
### 5.2 Camera surveillance

Moving to a more realistic, but still minimal, use-case, this example demonstrates how cleanly we can express a surveillance application from a robot's camera.

The camera should send frames to be displayed on the screen, only when motion is detected. Figure 25 shows the dataflow graph on the left and the corresponding *RHEA* code on the right.

In the code above, we can see how easy it is to view a *ROS* topic as a stream, using the *RosEvaluationStrategy* (line 3). Moreover, there is a nice separation between program logic (stream declaration in lines 6-17) and implementation details (*motionDetected* function in line 19).

<sup>13</sup><https://github.com/rhea-flow>



```

1 Stream.just(1)
2 .loop((entry: Stream[Int]) =>
3   (entry.multiply(2) mergeSort entry.multiply(3))
4     mergeSort
5     entry.multiply(5) : Stream[Int])
6 .distinct
7 .print
8
9 class IntStream(stream: Stream[Int]) {
10   def multiply(constant: Int): Stream[Int] =
11     stream.map(i => i * constant)
12
13   def mergeSort(other: Stream[Int]): Stream[Int] = {
14     val queue = new PriorityQueue[Int]()
15     Stream.zip(stream, other, (x, y) => {
16       val min: Int = Math.min(x, y)
17       val max: Int = Math.max(x, y)
18       queue.add(max)
19       if (min < queue.peek())
20         min
21       else {
22         queue.add(min)
23         queue.poll()
24       }
25     }).concatWith(Stream.from(queue))
26   }
27 }
28 implicit def enrichStream(st: Stream[Int]): IntStream =
29   new IntStream(st)

```

Figure 24: Hamming numbers

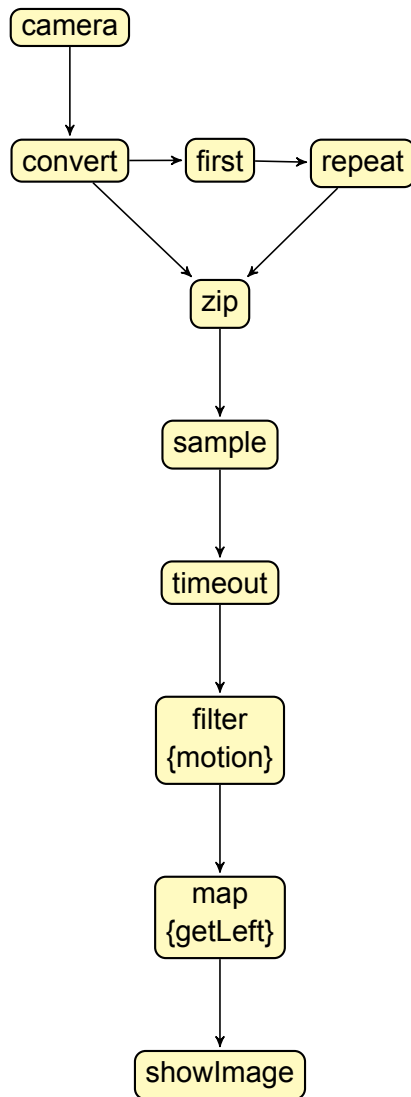
### 5.3 Robot control panel

This application concerns real-time monitoring of a robot, that is publishing its information and sensor-data to *ROS* topics, through a *graphical user interface (GUI)*.

The */camera/rgb* topic provides the frames of the robot's camera as coloured images, while the */camera/depth* provides frames that provide depth information. The */tf* topic publishes parent-child relations of the internal topics of the robot's configuration, and finally the */scan/* topic provides information from the robot's laser that gives horizontal depth information in polar coordinates.

The GUI displays the laser data embedded on the camera stream, while allowing for real-time face detection. Additionally, it displays the depth frames and the *tf* relations as a tree. Finally, a mock-up battery bar is displayed to show-case the framework's ability for simulation. Figure 26 illustrates the dataflow solution to the above problem on the left and the corresponding *RHEA* code on the right.

The implementation details (i.e. the visualization class and methods *faceDetect*(line 7), *embedLaser*(line 8) and *toGray*(line 26)) are not shown for brevity's sake. It is evident that this model of programming encourages a clean separation of concerns between the




---

```

1 Stream.configure(new HazelcastDistributionStrategy(
2   RxjavaEvaluationStrategy::new,
3   RosEvaluationStrategy::new,
4 ));
5
6 Stream<Mat> image =
7   Stream.from(new RosTopic<>("/camera"))
8     .map(CvImage::toCvCopy);
9
10 Stream<Mat> initial = image.first().repeat();
11
12 Stream.zip(image, initial, Pair::new)
13   .sample(100, TimeUnit.MILLISECONDS)
14   .timeout(1, TimeUnit.MINUTES)
15   .filter(Surveillance::motionDetected)
16   .map(Pair::snd)
17   .subscribe(window::showImage);
18
19 boolean motionDetected(Pair<Mat,Mat> pair) {
20   Mat m1 = pair.getLeft(), m2 = pair.getRight(), m = new Mat();
21   Core.absdiff(m1, m2, m);
22   Imgproc.threshold(m, m, 80, 255, Imgproc.THRESH_BINARY);
23   Imgproc.erode(m, m, Imgproc.MORPH_RECT(3,3));
24   for (int i = 0; i < m.rows(); i++)
25     for (int j = 0; j < m.cols(); j++) {
26       double[] pixel = m.get(i, j);
27       double sum = pixel[0] + pixel[1] + pixel[2];
28       if (sum > 50) return true;
29     }
30   return false;
31 }

```

---

Figure 25: Camera surveillance

individual components, namely between the sensor data manipulation and the actual visualization on the GUI.

#### 5.4 Robot hospital guide

As a final example, we will examine a more IoT-based application. Consider a robot that guides patients to different parts of a hospital, such as the gym or cafeteria. Assuming map localization, path finding and obstacle avoidance are already implemented, there still remains a problem with calibrating the robot's speed according to the patient's status.

To keep track of the patient's distance from the robot, each patient carries a smart-phone that acts as a *bluetooth low-energy (BLE) beacon*. The robot uses its bluetooth receiver to publish the distance from the signal source to an *MQTT* topic, which is then transformed by our stream application to velocity commands for the robot, in the form of slowing down or speeding up.

The first module constitutes the main program logic, where a declared dataflow graph acts

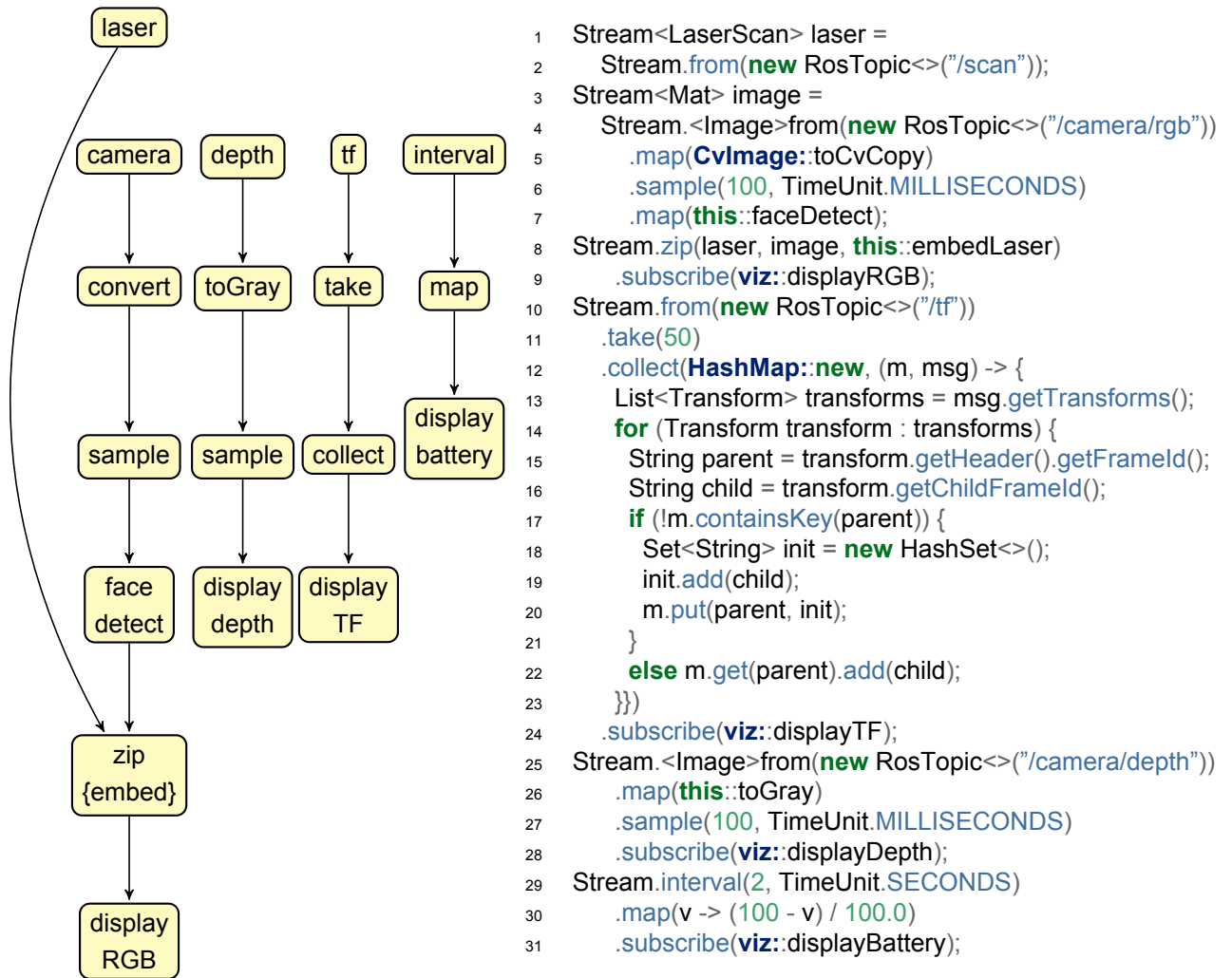


Figure 26: Robot control panel

as a stream transformation from beacon information to velocity commands to the robot. Figure 27 shows the dataflow graph on the left and the corresponding *RHEA* code on the right.

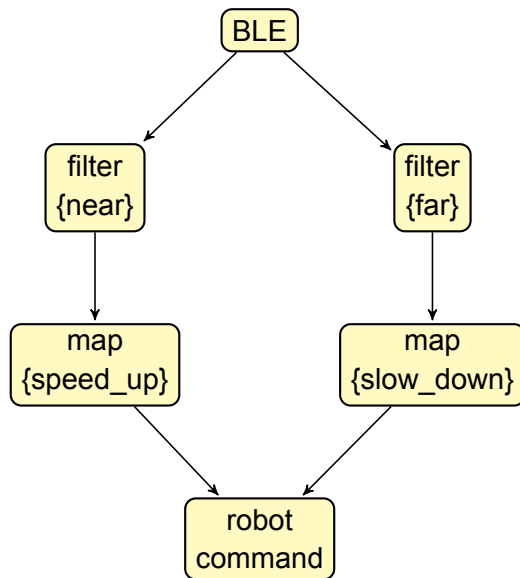
The second module just uses the *ReactiveBeacons* library<sup>14</sup> to get a stream of beacon data via *rxjava*, and then publishes it to a *MQTT* topic, which is the input of the first module. The corresponding *RHEA* code follows:

```

1 Stream.configure(new HazelcastDistributionStrategy(
2   RxjavaEvaluationStrategy::new,
3   MqttEvaluationStrategy::new
4 ));
5
6 Stream.from(ReactiveBeacons.observe())
7   .map(Beacon::getProximity)

```

<sup>14</sup><https://github.com/pwittchen/ReactiveBeacons>



```

1 Stream.configure(new HazelcastDistributionStrategy(
2     RxjavaEvaluationStrategy::new,
3     RosEvaluationStrategy::new,
4     MqttEvaluationStrategy::new
5 ));
6
7 Topic<RobotCommand> vel = new RosTopic<>("/robot/cmd");
8
9 Stream<Proximity> ble =
10     Stream.from(new MqttTopic<>("/ble"));
11
12 ble.filter(Proximity::isNear)
13     .map(d -> Commands.SPEED_UP)
14     .subscribe(vel);
15
16 ble.filter(Proximity::isFar)
17     .map(d -> Commands.SLOW_DOWN)
18     .subscribe(vel);
  
```

Figure 27: Robot hospital guide

```

8     .subscribe(new MqttTopic<>("/ble"));
  
```

This example clearly show-cases the framework's ability to combine different technologies and act as a high-level, declarative coordination language.



## 6. RELATED WORK

This section discusses related work in the fields of *Big Data*, *Robotics* and *IoT*.

### 6.1 Big Data

The necessity for implicit parallelism and distribution of more and more applications, dealing with huge and/or complex data, has brought increasingly more attention to the dataflow programming model. Its easy to understand and maintain structure and declarative approach to programming, while not losing expressibility, has attracted many frameworks to utilize it.

#### 6.1.1 GoogleDataflow

Google recently released the *GoogleDataflow* framework<sup>15</sup>, which is an evolution of *FlumeJava*[14], which in turn is a successor of the famous *MapReduce* framework[15].

*MapReduce* was a very simple model that allowed automatic concurrency/distribution on a cluster by allowing only a very minimal program structure. First, the user specifies a *map* function that processes a key/value pair to generate a set of intermediate key/value pairs, and a *reduce* function that merges all intermediate values associated with the same intermediate key. Although it was widely adopted at first, quickly many problems that could not be expressed with the above formalism were found and therefore a more expressive model was required.

This gave birth to *FlumeJava*, which is a generalization of the *MapReduce* framework that allows more expressive pipelines of more primitive operations like *MapReduce*.

Although *FlumeJava* was more attractive due to its expressibility, still the pipeline constructed could not formulate all problems that are needed in some big-data applications. For instance, the constructed dataflow could not contain cycles, which is an integral part of *incremental computation*, used extensively nowadays for machine learning and data analysis.

And this is how *GoogleDataflow* came to exist, offering a fully generic dataflow framework integrated with many other closely-related technologies from Google<sup>16</sup>, like Cloud Storage, Cloud PubSub, Cloud Datastore, Cloud Bigtable and BigQuery.

It is open-source, offers fully automatic resource management that auto-scales for optimal throughput and provides increased reliability and data consistency. Moreover, it provides a unified programming model through its API, while allowing data monitoring and demand-driven execution.

---

<sup>15</sup><https://cloud.google.com/dataflow/>

<sup>16</sup><https://cloud.google.com>

### 6.1.2 TensorFlow

Another dataflow framework from Google is *TensorFlow*<sup>17</sup>, which is an open-source polyglot library for machine learning and especially construction of neural networks.

The interesting fact is that, although it started out as a rigid neural network library, it quickly generalized to a dataflow construction library, much similar to my own project, which started out as a robotics library.

Its main features are its portability to multiple computational architectures (e.g. CPU, GPU, etc...) and multiple language APIs (e.g. C++, Python), although its main advantage are its domain-specific operators for neural nets (i.e. common subgraphs, auto-differentiation).

Though the edges/streams connecting the nodes, only a single but flexible data type is allowed, namely the *Tensor* type, which essentially is a multi-dimensional array that usually represents features or weights. Figure 28 illustrates a neural network as a dataflow graph.

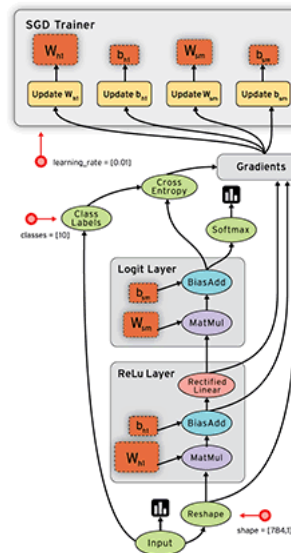


Figure 28: TensorFlow graph

### 6.1.3 Akka

Definitely one of the most mature frameworks for distribution targeting the JVM, *Akka*<sup>18</sup> is a toolkit and runtime for highly concurrent, distributed and resilient message-driven applications. It is also one of the founders of the *Reactive Streams*[13] initiative.

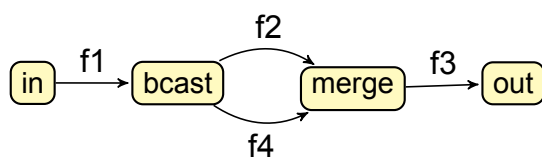
Its approach follows the *Actor* model[16], where one perceives abstract computational agents, called actors, that are distributed in space and communicate with point-to-point messages that are buffered in a queue. In reaction to a message, an actor can create more actors, make local decisions, send more messages and determine how to respond

<sup>17</sup><https://www.tensorflow.org/>

<sup>18</sup><http://akka.io/>

to the next message received.

Similar to the problem of *ROS* that my framework solved, which is the inappropriate nature of callbacks for complex scenarios, *Akka* developers also felt the necessity for a more flexible and composable programming model, so they developed the *AkkaStreams* library<sup>19</sup> which provides a convenient API for stream processing and also dataflow graph construction with an interesting DSL. Figure 29 demonstrates a dataflow graph on the left, generated by the DSL code on the right.



```

1  val g = FlowGraph { implicit b =>
2    import FlowGraphImplicits._
3    val in = Source(1 to 10)
4    val out = Sink.ignore
5    val bcast = Broadcast[Int]
6    val merge = Merge[Int]
7    val f1, f2, f3, f4 = Flow[Int].map(_ + 10)
8    in ~> f1 ~> bcast ~> f2 ~> merge ~> f3 ~> out
9      bcast ~> f4 ~> merge
10 }

```

Figure 29: Akka DSL

#### 6.1.4 Spark

A very well-known and well-adapted framework for scalable large-data processing is Apache's *Spark*<sup>20</sup>. It was developed to overcome the shortcomings of the *MapReduce* framework mentioned above, providing a much more efficient and flexible runtime.

It follows the same general approach as *RHEA*, meaning that it is completely generic and encourages domain-specific libraries to be built upon it. For instance, *MLib*<sup>21</sup> is a library for machine learning and *GraphX*<sup>22</sup> is a library for iterative graph algorithms, both stacked upon *Spark*.

It offers a rich set of data-parallel operators ( $\simeq 80$ ) that can be used interactively from Scala, Python, Java or R. The code below shows the classic word-counting example in Spark's Scala API.

```

1  Spark.textFileStream("hdfs://...") /* Get file stream */
2    .flatMap(_ .split(" ")) /* Split into words */
3    .map(x => (x, 1)).reduceByKey(_ + _) /* Count words */

```

<sup>19</sup><http://doc.akka.io/docs/akka-stream-and-http-experimental/1.0-M2/scala.html>

<sup>20</sup><http://spark.apache.org/>

<sup>21</sup><http://spark.apache.org/mlib/>

<sup>22</sup><http://spark.apache.org/graphx/>

### 6.1.5 Naiad

Offering the high throughput of batch processors, the low latency of stream processors and the ability to perform iterative and incremental computations at the same time is extremely challenging and none of the aforementioned frameworks manage to provide it. Applications that need all these features need to rely on multiple platforms, at the expense of efficiency, maintainability and simplicity.

Naiad[17] combines all of these features in a unifying framework, that provides a generic low-level platform, that a wide variety of high-level programming models can be built upon, enabling such diverse tasks as streaming data analysis, iterative machine learning, and interactive graph mining.

Its main contribution is the definition of a new computational model, namely the *Timely Dataflow* model, which is an extension to the dataflow model I introduced in the first chapter, by allowing a more efficient and lightweight coordination mechanism for capturing opportunities for parallelism. This is achieved by enriching the dataflow model with timestamps that represent logical points in the computation.

Figure 30 shows a Naiad application that supports real-time queries on continually updated data.

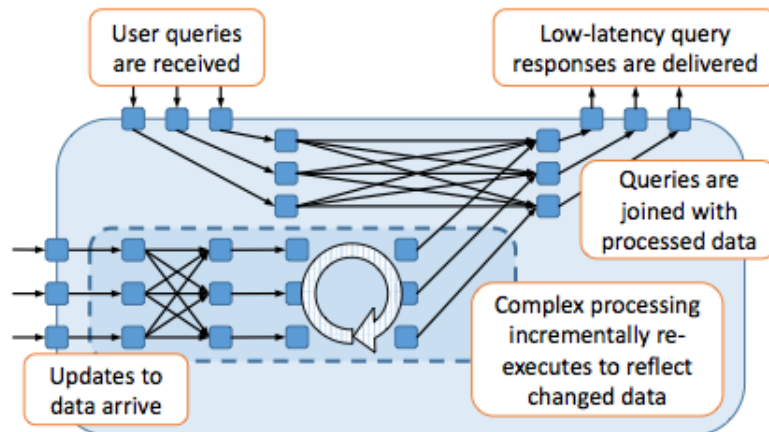


Figure 30: Naiad application

### 6.1.6 dispel4py

A less-known framework for Python is *dispel4py*<sup>23</sup>. It provides the ability to describe abstract workflows for distributed data-intensive applications.

Similar to my *EvaluationStrategy* concept, it allows different mappings to enactment systems, such as MPI<sup>24</sup> and Apache Storm<sup>25</sup>.

Its main disadvantages are that it has only an API for Python and only allows low-level

<sup>23</sup><https://github.com/dispel4py>

<sup>24</sup><http://www.mcs.anl.gov/project/mpich-high-performance-portable-implementation-mpi>

<sup>25</sup><http://storm.apache.org/>

specification of the graph's nodes, through the definition of *Processing Elements*. Therefore, it is inconvenient to compose larger graphs from simpler ones and the source code becomes chaotic and difficult to maintain.

## 6.2 Robotics

It is only natural that the dataflow model would make its way through the field of robotics, as many behaviours in control theory are expressed as dataflow diagrams.

### 6.2.1 roshask

*Roshask*[18] is a binding from the Haskell programming language to the basic *ROS* interfaces. Like *RHEA*, the approach is to overcome the shortcomings of *ROS* callbacks by viewing topics as streams. This allows for, and encourages, a higher level of abstraction in robot programming, while making the fusing, transforming and filtering of streams fully generic and compositional.

Below is the classic Talker-Listener *ROS* example, where one node publishes messages to a topic and another one listens for them.

---

```

1 sayHello :: Topic IO S .String
2 sayHello = repeatM (fmap mkMsg getTime)
3 where mkMsg = S .String . ("Hi "++) . show
4
5 tn :: Node ()
6 tn = advertise "chatter" (topicRate 1 sayHello)
7
8 main :: IO ()
9 main = runNode "talker" tn
10
11 showMsg :: S .String -> IO ()
12 showMsg = putStrLn . ("Msg: "++) . S.data
13
14 main = runNode "listener" $
15 runHandler showMsg = << subscribe "chatter"

```

---

### 6.2.2 Yampa

*RHEA* and *roshask* were heavily influenced by the work of Hudak's group (Yale Haskell Group) on robot DSLs and FRP in general[1, 5, 6, 7, 3, 4].

*Yampa*<sup>26</sup> is a DSL embedded in Haskell that realizes the FRP model, using arrows to minimize time-/space- leaks. Figure 31 shows the primitive stream operators that are supported.

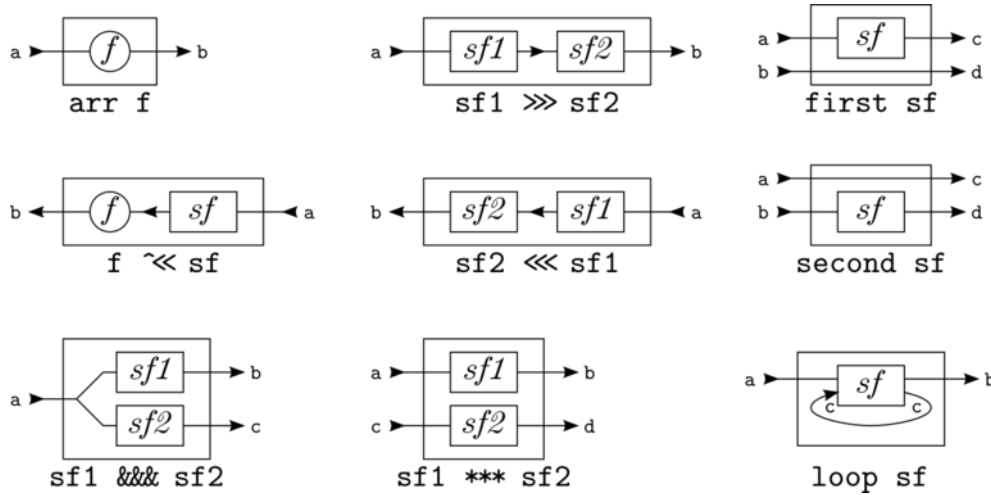


Figure 31: Yampa operators

### 6.2.3 Flowstone

*Flowstone*<sup>27</sup> is a programming environment that mixes graphical and text based programming in Ruby that can be used for robotics, image/signal processing and interconnecting heterogeneous sources. It follows a variant of the dataflow model, where applications are built by linking together functional blocks called components. Figure 32 shows a screenshot of the graphical environment, where we discern both graphical and textual elements.

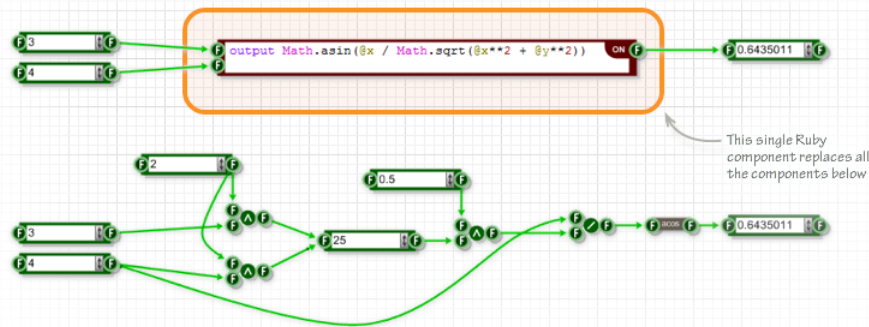


Figure 32: Flowstone screenshot

Its main advantage is that it is stand-alone, so no compiling is necessary, which allows for rapid prototyping.

<sup>26</sup><https://wiki.haskell.org/Yampa>

<sup>27</sup><http://www.dsrobotics.com/flowstone.html>

## 6.3 Internet of Things

IoT applications often deal with much heterogeneity, due to the variety of sources that different devices introduce. Therefore, a component-based approach suits well to solve this problem and there are some dataflow frameworks that follow that approach.

### 6.3.1 NoFlo

*NoFlo*<sup>28</sup> is a JavaScript implementation of *Flow-based Programming*[19], which is a particular form of dataflow programming, based on bounded buffers, information packets with defined lifetimes, named ports and separate definition of connections.

*NoFlo* applications consist of components that are connected together in a graph. This allows for clear separation of control flow from the actual software logic, helping you organize large applications easier than traditional OOP paradigms.

You can design *NoFlo* applications using a web-based graph editor<sup>29</sup>, which is depicted in figure 33.

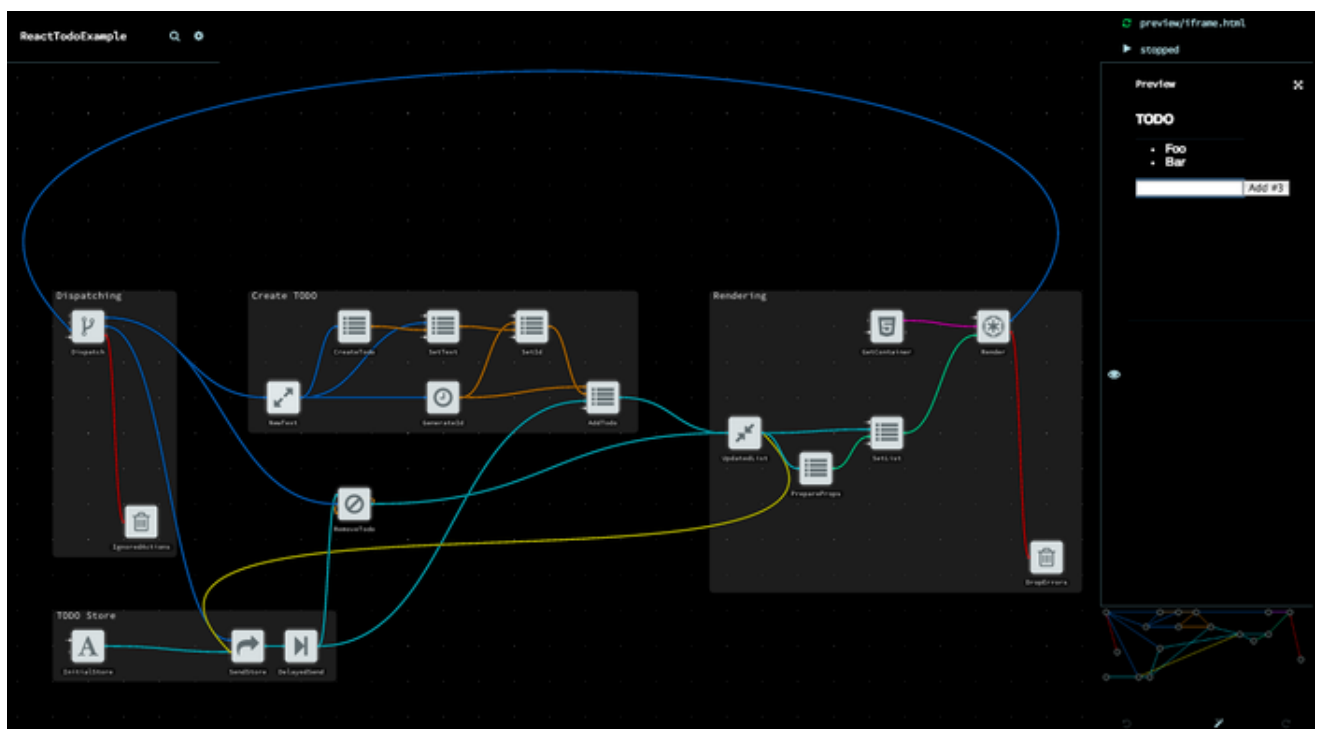


Figure 33: NoFlo graph editor

### 6.3.2 Node-RED

Another interesting IoT tool for JavaScript following a dataflow approach is *Node-RED*[20], which is a visual tool for wiring together hardware devices, APIs and online services in new and interesting ways.

<sup>28</sup><http://noflojs.org/>

<sup>29</sup><https://flowhub.io/>

Applications called flows, are built immediately on a browser, and can be deployed on the Cloud with just a single click. The main advantage of this tool is that it encourages social development, due to the fact that flows are stored in JSON format, which can be easily imported and exported for sharing with others. The online flow library<sup>30</sup> has had a huge number of contributions so far. Figure 34 shows a screenshot of the editor.

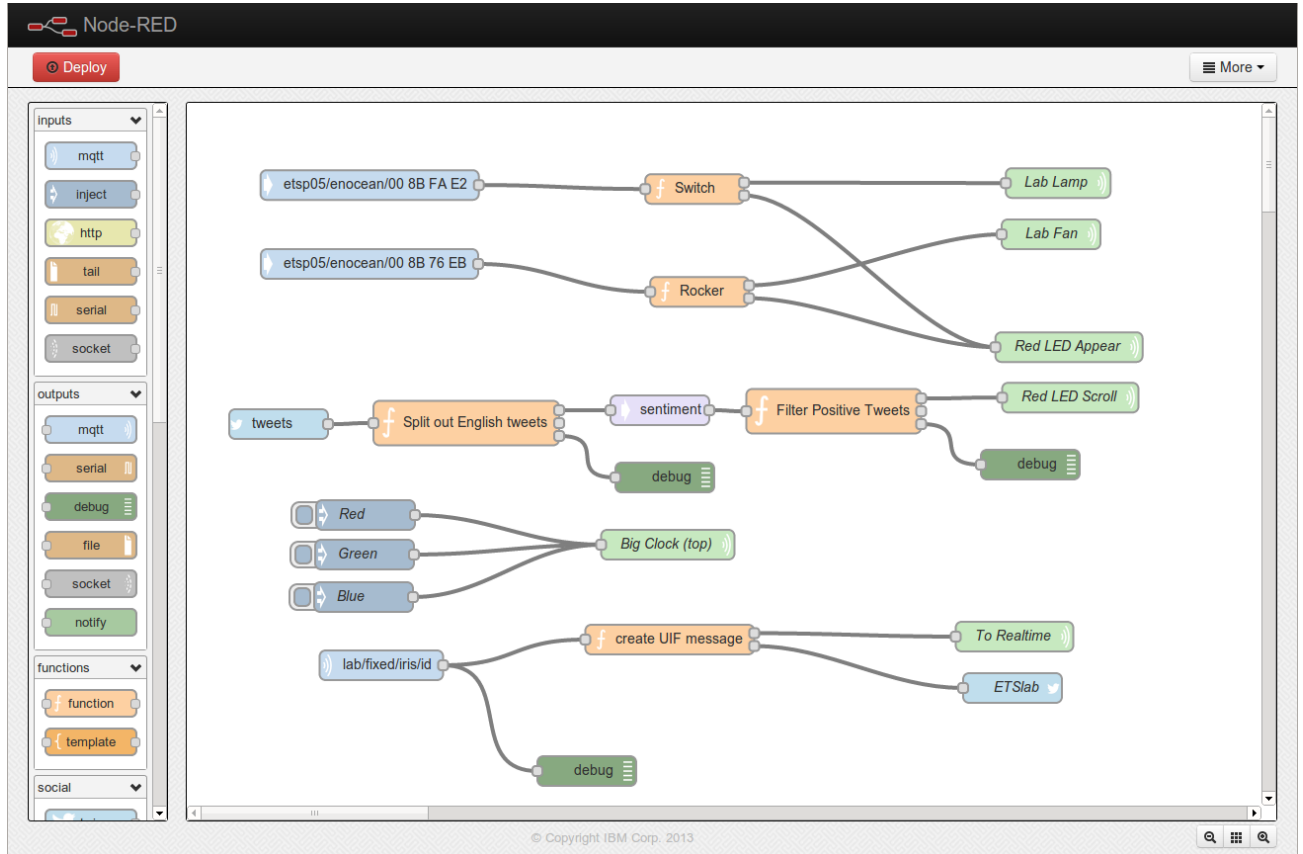


Figure 34: Node-RED graph editor

<sup>30</sup><http://flows.nodered.org/>



## 7. FUTURE WORK

This section discusses interesting topics and ideas for future contribution and extension of the framework.

### 7.1 More strategies

Every component that was initially intended to be replaceable for the sake of flexibility, was implemented using the strategy design pattern. So the definite starting point for contribution would be to implement more concrete strategies.

For evaluation, a good idea would be a low-level implementation in C/C++ to allow compatibility with older/smaller systems, not capable of running the JVM. Additionally, more JVM-based strategies are possible, by using the standard *java.util.Stream* library<sup>31</sup>, *Akka Streams* or *Scala Iteratees*[21].

For distribution, again a more low-level approach written in C/C++ (e.g. using *MPI/OpenMP*) would be beneficial to the variety of systems that are allowed to cooperate freely. Moreover, many frameworks similar to Hazelcast could replace it to allow for easy integration with applications that are already committed to another framework.

### 7.2 Dynamic reconfiguration

A definite shortcoming of *RHEA* is that all configuration takes place initially and remains static throughout execution. This poses a problem for many scenarios, where environment is constantly changing and available resources may be introduced or become obsolete.

For instance, in a distributed application that controls the behaviour of a *robot swarm* and all communication is done through a *wireless local area network (WLAN)*, a robot may lose signal and become unavailable at runtime. Later on it may rejoin the network, so it is vital that its sensor input and computational power become quickly available to the system. A nice DSL that is used to specify adaptive behaviour in robot navigation can be found on [?].

Apart from robotics, the need for adaptive behaviour of software systems has been evident for a long time[?]. There have also been attempts for runtime adaptation of, specifically, stream processing systems, such as the *Flexstream* framework[?].

A JVM-based technology that is relatively easy to integrate into *RHEA* is *HotWave*[?], which is an *aspect-oriented programming (AOP)* framework, based on the famous *AspectJ*<sup>32</sup> open-source extension of Java, that supports dynamic (re)weaving of previously loaded classes, and ensures that all classes loaded in a JVM can be (re)woven. The contribution idea here is to integrate *HotWave* into the *org.rhea\_core* package and then use it to specify the desired adaptive behaviour for reconfiguring where nodes are executed,

<sup>31</sup><https://docs.oracle.com/javase/8/docs/api/java/util/stream/package-summary.html>

<sup>32</sup><http://www.eclipse.org/aspectj/>

what operation they perform, and so forth.

### 7.3 Advanced network profiling

The current strategy for network profiling is based on *round-trip time (RTT)* which is measured explicitly by pinging from each machine to every other one. As the number of machines increase, calculating the *RTT* becomes extremely expensive and may outweigh the benefits of exploiting network proximity.

To overcome the aforementioned problem, one could estimate, instead of measuring exactly, the *RTT* with less computationally expensive methods. There has been extensive research on *RTT* estimation for *peer-to-peer (P2P)* networks[22]. Most of the methods proposed suffer from the constraint that estimation is calculated on a single machine, therefore introducing a central point of failure. This is highly unsuited for distributed system, where a decentralized approach should be used. One such attractive approach is the *Vivaldi* coordinate system[23], which is a simple and lightweight algorithm that assigns synthetic coordinates to hosts such that the distance between the coordinates of two hosts accurately predicts the communication latency between them. The algorithm is completely decentralized (i.e. the same piece of code runs on every host) and experiments show that it achieves a median relative error of 11%, where the error is characterised by the squared-error function:

$$E = \sum_i \sum_j (L_{ij} - ||x_i - x_j||)^2$$

where  $L_{ij}$  is the actual RTT between nodes  $i$  and  $j$  and  $||x_i - x_j||$  is the euclidean distance between their coordinates.

### 7.4 Advanced fault-tolerance

An aspect that *RHEA* is far behind from most of its competing frameworks is fault-tolerance. Although there is a minimal control on *back-pressure*, available through the small provided set of such operators, still there are no advanced methods for specifying behaviour for graceful error-recovery.

The above is essential for large machine clusters, in which systems it is certain that host failures and other faults will be a common occurrence. The functional nature of the dataflow model enables fault-tolerance, in addition to parallelism, due to the fact that a node can be moved to another machine for execution, while preserving the original semantics.

The extension proposed is to first provide low-level functionality for moving nodes across the network at runtime, and then use that to design high-level error recovery mechanisms. The issue confronted here is that, at the current state of the system, everything is con-

figured statically before evaluation. To enable runtime configuration, it is mandatory to refactor major internal components, as previously noted.

This contribution path can draw heavy influence from recent research on fault-tolerance for stream processing engines[?, ?, ?]. These provide efficient models for availability and data recovery/consistency, by using data replication and even parallel recovery of lost state across the cluster.

## 7.5 Integration with other technologies

An issue that should not be neglected is interoperability with other dataflow frameworks mentioned throughout this thesis. Although the usage of *RSS* partly aids that cause, not all frameworks for the JVM support it, especially older ones, and is practically useless concerning non-JVM frameworks.

A nice idea would be an *ExportStrategy* that provides a one-to-one mapping between different dataflow representations. For instance, the *NodeRedExportStrategy* would export a single JSON file ready to be imported into *Node-RED* and deployed immediately. An issue that quickly arises is that of expressibility, meaning that the target platform should be at least as expressible as *RHEA*. This can be solved by providing a mechanism/DSL to specify restricted views of the *Stream* class, which allow for a subset of the original operation set.

## 7.6 Visual language

The dataflow programming model also provides a very intuitive graphical representation to the structure of the software system being implemented. Thus, a helpful tool would be a graphical environment, where graph construction can be achieved through the user interface and delegate only the stream declaration part to a text editor. Another useful functionality would be clustering together sub-graphs by zooming out to provide clarity in complex systems.

Another useful utility is to provide visual debugging support, so as to be able to monitor values and errors going through the streams real-time or by playing back the recorded actions. Moreover, if dynamic reconfiguration is supported, the visual debugger could display where each node is being executed at each moment in time and other relative information.

## 7.7 Stream reasoning

A very helpful extension to any dataflow framework is the ability to reason about, in formal logic, the ever-changing streams of data. This creates the need for a new formal logic that is able to capture this flowing notion, which has also troubled the new *Semantic Web* field[?]. Fortunately, various sound and efficient logic formalisms have been conceived

recently[?, ?].

Stream reasoning has also been researched in the robotics domain[?], which led to the development of a stream reasoning framework for *ROS* called *DyKnow*[?]. Integrating *DyKnow* with *RHEA*, as they both can support *ROS*, would certainly be a meaningful contribution. Figure 35 shows the general architecture of the *DyKnow* framework.

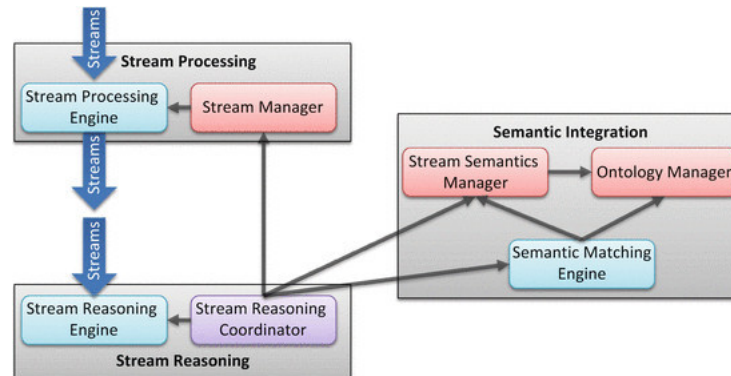


Figure 35: DyKnow architecture

## 8. CONCLUSIONS

The framework described in this thesis was designed with extensibility in mind, aiming to act as a fundamental basis, onto which various domain-specific libraries or DSLs will rely in the future. To that end, a constant effort to generalize and make components as abstract as possible was made.

The set of operators aided expressibility, making it possible to specify any dataflow graph in a concise and readable manner. This disallowed optimizations suitable for less expressive models (e.g. *Map-Reduce*), but recent research has shown that general dataflow topologies have optimization opportunities that are yet to be found[?]. A minimal optimization stage has been implemented, which paves the path to more advanced optimization techniques, such as those used in *Naiad*[17] and *Stratosphere*[?].

Figure 36 illustrates all the pluggable components of the framework around the core, which are normally deployed in separate libraries.

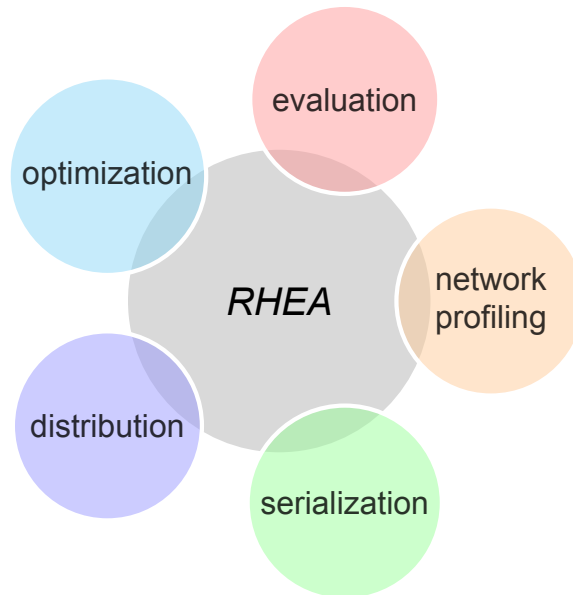


Figure 36: The *RHEA* ecosystem

The applications demonstrated the framework’s ability to provide a higher level of abstraction, where the language only specifies how different components coordinate, without knowledge of the implementation details. This is exactly what *Ziria* accomplishes in the domain of wireless systems programming[?]. The driving force for both frameworks (i.e. *RHEA* and *Ziria*) is that some specific domains have fixated their methods on low-level programming, whereas more satisfactory paradigms can solve many shortcomings.

This is a general notion in computer science, owing its existence to the fact that the problems we are facing are getting increasingly more complex, while resources meet certain realistic bounds, and therefore a higher abstraction layer is mandatory for maintaining readability, efficiency and expressibility.

## ABBREVIATIONS, INITIALS AND ACRONYMS

A table of all abbreviations used throughout the thesis follows.

FRP	Functional Reactive Programming
JVM	Java Virtual Machine
NCSR	National Centre for Scientific Research
ROS	Robot Operating System
IoT	Internet of Things
CPU	Central Processing Unit
TCP	Transmission Control Protocol
PubSub	Publish/Subscribe
OOP	Object-oriented Programming
UML	Unified Modelling Language
GPU	Graphics Processing Unit
DSL	Domain-specific Language
RSS	Reactive Streams Standard
API	Application Programming Interface
MPI	Message Passing Interface
JSON	JavaScript Object Notation
RTT	Round-Trip Time
P2P	Peer-To-Peer
AOP	Aspect-Oriented Programming
GUI	Graphical User Interface
WLAN	Wireless Local Area Network

## REFERENCES

- [1] C. Elliott and P. Hudak, “Functional reactive animation,” in *ACM SIGPLAN Notices*, vol. 32, pp. 263–273, ACM, 1997.
- [2] E. Amsden, “A survey of functional reactive programming,” *Unpublished*, 2011.
- [3] Z. Wan, W. Taha, and P. Hudak, “Real-time frp,” in *ACM SIGPLAN Notices*, vol. 36, pp. 146–156, ACM, 2001.
- [4] C. M. Elliott, “Push-pull functional reactive programming,” in *Proceedings of the 2nd ACM SIGPLAN symposium on Haskell*, pp. 25–36, ACM, 2009.
- [5] P. Hudak, A. Courtney, H. Nilsson, and J. Peterson, “Arrows, robots, and functional reactive programming,” in *Advanced Functional Programming*, pp. 159–187, Springer, 2003.
- [6] J. Peterson, P. Hudak, and C. Elliott, “Lambda in motion: Controlling robots with haskell,” in *Practical Aspects of Declarative Languages*, pp. 91–105, Springer, 1999.
- [7] Z. Wan, W. Taha, and P. Hudak, “Event-driven frp,” in *Practical Aspects of Declarative Languages*, pp. 155–172, Springer, 2002.
- [8] J. Hughes, “Generalising monads to arrows,” *Science of Computer Programming*, vol. 37, no. 1–3, pp. 67 – 111, 2000.
- [9] R. Paterson, “A new notation for arrows,” *SIGPLAN Not.*, vol. 36, pp. 229–240, Oct. 2001.
- [10] M. Quigley, K. Conley, B. Gerkey, J. Faust, T. Foote, J. Leibs, R. Wheeler, and A. Y. Ng, “Ros: an open-source robot operating system,” in *ICRA workshop on open source software*, vol. 3, p. 5, 2009.
- [11] E. Gamma, *Design patterns: elements of reusable object-oriented software*. Pearson Education India, 1995.
- [12] B. Jonas, D. Farley, R. Kuhn, and M. Thompson, “Reactive manifesto.” <http://www.reactivemanifesto.org/>, 2014.
- [13] “Reactive streams standard.” <http://www.reactive-streams.org/>, 2015.
- [14] C. Chambers, A. Raniwala, F. Perry, S. Adams, R. R. Henry, R. Bradshaw, and N. Weizenbaum, “Flumejava: easy, efficient data-parallel pipelines,” in *ACM Sigplan Notices*, vol. 45, pp. 363–375, ACM, 2010.
- [15] J. Dean and S. Ghemawat, “Mapreduce: simplified data processing on large clusters,” *Communications of the ACM*, vol. 51, no. 1, pp. 107–113, 2008.
- [16] C. Hewitt, P. Bishop, and R. Steiger, “A universal modular actor formalism for artificial intelligence,” in *Proceedings of the 3rd international joint conference on Artificial intelligence*, pp. 235–245, Morgan Kaufmann Publishers Inc., 1973.
- [17] D. G. Murray, F. McSherry, R. Isaacs, M. Isard, P. Barham, and M. Abadi, “Naiad: a timely dataflow system,” in *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles*, pp. 439–455, ACM, 2013.
- [18] A. Cowley and C. J. Taylor, “Stream-oriented robotics programming: The design of roshask,” in *Intelligent Robots and Systems (IROS), 2011 IEEE/RSJ International*

- Conference on*, pp. 1048–1054, IEEE, 2011.
- [19] J. P. Morrison, “Flow-based programming,” in *Proc. 1st International Workshop on Software Engineering for Parallel and Distributed Systems*, pp. 25–29, 1994.
  - [20] M. Blackstock and R. Lea, “Toward a distributed data flow platform for the web of things,” in *Web of Things (WoT), 2014 5th International Workshop on the*, 2014.
  - [21] O. Kiselyov, “Iteratees,” in *Functional and Logic Programming*, pp. 166–181, Springer, 2012.
  - [22] S. Lukin, “A comparison of round-trip time estimation algorithms,” *Loyola University Maryland, UCSC SURF-IT Research*, 2010.
  - [23] F. Dabek, R. Cox, F. Kaashoek, and R. Morris, “Vivaldi: A decentralized network coordinate system,” in *ACM SIGCOMM Computer Communication Review*, vol. 34, pp. 15–26, ACM, 2004.