

# **A multi-protocol communication layer for FastFlow's distributed runtime**

**Federico Finocchio** - 516818

MSCS: Data and Knowledge. [f.finocchio@studenti.unipi.it](mailto:f.finocchio@studenti.unipi.it)

## **Abstract**

Handling the explosion in data production rates both for scientific and commercial applications require amount of storage and computational power which can't be handled in single machines, thus a shift to distributed computing is being recorded. This shift in programming environment requires new tools and frameworks, as well as new abstractions, to allow application developers to provide efficient applications which deal with real-time decisions in an accurate and actionable way.

FastFlow framework provides a set of abstraction that targets both the application and RTS developer, allowing via structured parallel programming the development of efficient parallel applications via a compositional model based on LEGO-style building blocks. In this work we introduce a new abstraction at the communication layer, allowing to connect remote nodes via a multitude of transports in different parts of the FastFlow application, allowing heterogeneity in computing nodes.

# Contents

Introduction . . . . .	4
<b>1 Background</b>	<b>7</b>
1.1 Towards Exascale computing . . . . .	8
<b>2 FastFlow</b>	<b>10</b>
2.1 Building Blocks . . . . .	10
2.2 FastFlow distributed runtime . . . . .	11
<b>3 Mercury</b>	<b>13</b>
3.1 RPC: a Mercury's perspective . . . . .	14
3.2 Architecture . . . . .	14
3.3 Resilience and Fault Tolerance . . . . .	20
3.4 Dependencies . . . . .	20
3.5 Mercury analysis . . . . .	21
<b>4 Argobots</b>	<b>24</b>
4.1 Building blocks . . . . .	26
4.2 Operations . . . . .	27
4.3 Memory Management . . . . .	28
<b>5 Margo</b>	<b>29</b>
5.1 RPC registration process . . . . .	30
5.2 Margo analysis . . . . .	35
<b>6 Implementation</b>	<b>37</b>
6.1 Communication classes . . . . .	37
6.2 RPC based service . . . . .	40
6.3 Splitting taxonomy . . . . .	41

<b>7</b>	<b>Testing</b>	<b>44</b>
7.1	A sample application . . . . .	44

# Introduction

In the past years we are experiencing an explosion in the quantities of produced data, coming from scientific simulations, experimental facilities, connected devices (IoT) and commercial applications. These huge amount of data challenge, the storage and computing capabilities of single machines as well the network being used to transfer data between multiple machines due to I/O bottleneck [1]. Moreover, the amount of computation needed to handle the necessary data for modern applications, like (near) real-time applications that generate near-continuous data, is much bigger than the one offered by current monolithic systems. Thus, as we witnessed the shift from sequential to parallel programming, we are now observing a new shift from parallel to distributed computing. Moreover, as the difference between Big Data and High Performance Computing becomes more and more blurred, and due to this new evolution in the way application are developed and deployed, new frameworks, programming models and tools have surfaced to help the programmer to deal with the complexity of designing parallel applications above the underlying distributed infrastructures [1, 2]. In the transition to distributed computing, new extensive parallel processing and new analytics algorithms are needed in order to provide timely and actionable information [3]. However, dealing with scalable and parallel applications that can be deployed across the network is not an easy task, and abstractions are needed in order to allow the application developers to focus on the optimization of the uppermost layer, leaving, among all, communication and synchronization details to experts on the respective fields.

The shift towards distributed computing introduces new challenges, like resource management, data distribution, coordination of participating processes and monitoring of the application. Given that each of these challenges acts at a different abstraction layer, application developers often prefer to rely over high-abstraction frameworks to abstract from the distributed challenges and focus only on the application development.

The thesis aims at extending the existing distributed-memory runtime system

for FastFlow, a C++ structured parallel programming library originally targeting shared-memory platforms. The thesis introduces a novel communication layer, which offers a standardized API in order to abstract from the underlying transport used for communications between remote nodes. The provided abstraction allows the *application programmer* to design new applications, as well as extending existing ones, by leveraging the classes introduced in this work in order to connect different nodes on the network without having to deal with communication specific implementation. The implemented layer provides an automatized way of plugging different protocols inside existing application without breaking the code and with little to zero effort. The *RTS programmer* is provided with a set of extensible classes in order to allow the integration of existing transports or even future ones. Furthermore, in order to follow and maintain the LEGO style *philosophy* of the FastFlow library, also the implemented communication layer is composed of a set of building blocks which can be composed in order to address the needs of the specific application. Allowing composability of components is very important in heterogeneous environments such as the cloud, in order to allow easy porting of applications and adaptation to environment changes. To fulfill this need, the communication layer presented in this thesis allows each application subset to communicate using the protocol that is most suited for the task at hand or for the architecture specific limitations. Vendor specific transports can represent a limitation in horizontally scaling applications, in particular if the provided transports are completely *obscure* to the application programmer, which should not care of the underlying communication aspects when programming a distributed application. The implemented building blocks allow for this reason to be run above a multitude of different transports by simply plugging the protocol-specific building block.

A natural preliminary phase of this work was strictly tied to the analysis of communication frameworks which are available to the application programmer in order to leverage efficiently multiple protocols during communication between distributed groups. We aim at building a communication layer which is agnostic of the underlying protocol, thus, we first implemented the communication classes by means of a multi-protocol RPC framework in order to enable communication between nodes in a protocol-agnostic fashion. In the first part of this work we analysed the Margo [4] library, its underlying RPC framework called Mercury [4, 5] and the threading library which is used by Margo as a runtime framework, called Argobots [6]. The three analyzed frameworks are part of a much broad set of frameworks which are referred to as *Mochi core*, as specified in [4]. The Mochi core provides a set of frameworks to enable communication, concurrency management and storage with a composition model. As appeared by the testing we performed, but as also described by its de-

velopers, the Mercury RPC framework is purposefully built to leverage HPC fabrics vendor protocols, sometimes resulting in a lack of support and faulty implementation for more classical transports like MPI and TCP. This led us to adapt the existing MPI and TCP drivers in the distributed FastFlow RTS into the novel API.

The thesis uses a bottom-up approach to illustrate the overall picture of the used frameworks and to describe internally each of those frameworks. The thesis is organized as follows:

- **Chapter 1:** illustrates the state-of-the-art regarding communication frameworks in the HPC and BDA environments;
- **Chapter 2:** presents the FastFlow library and its current stage of development, illustrates the distributed runtime which is further extended by this work;
- **Chapter 3:** presents the core framework which provides the building blocks needed to implement RPC calls leveraging multiple transports. We describe the general structure of the framework, which will serve to describe the higher-level framework used to implement the communication layer in FastFlow's distributed runtime;
- **Chapter 4:** describes the runtime system which is used by the higher-level RPC framework in order to manage concurrent execution of different RPCs, to allow asynchronous execution of RPC-related callbacks and to simplify Mercury's progress loop;
- **Chapter 5:** describes the high level framework used for the implementation of the multi-protocol interface for the communication layer in FastFlow's distributed runtime. We show how this framework leverages Mercury and Argobots functionalities to provide a parallel and efficient multi-protocol RPC framework;
- **Chapter 6:** gives insights on the integration of Margo's functionalities into the FastFlow's distributed communication layer. We show pseudo-code of the implemented classes and common use-cases for both the application and runtime programmer;
- **Chapter 7:** focuses on the presentation of sample applications and testing we performed to evaluate the effectiveness of the implemented communication layer.

# Chapter 1

## Background

Challenges brought by the new-generation of data-intensive applications require co-operation in the two fields of High Performance Computing and Big Data Analytics [1]. As a result, we have witnessed an evolution towards distributed computing, that allows to leverage a larger amount of storage and computing power by means of multiple interconnected machines [2]. In the pursuing of the “parallel revolution aggressive goal” of “writing parallel programs that scale with the number of cores and are as easy to write and efficient as sequential programs” [9], we are experiencing an increase in the amount of abstractions, tools, frameworks and libraries specifically targeted at helping the application developers writing parallel programs in an efficient and structured fashion. In this chapter we particularly focus on all those tools that help the programmer in structuring their programs without being required to handle the underlying communication infrastructure challenges.

We consider the *HPC communication layer* as “*any framework, library, tool, protocol or pattern that eases the communication between distributed processes*”, and its only purpose is to abstract the user from the underlying infrastructure and protocols by providing a set of APIs to send and receive data. Following this line, and by considering the distinctions made by [2], we can consider two main leading paradigms:

- **Remote Invocation:** allows the user to call functions from one node to the other across the network. We consider in this category both *Remote Method Invocation* (RMI) and *Remote Procedure Call* (RPC), with the only distinction between the two being they are used to invoke object’s methods or program functions, respectively.
- **Message Oriented:** also in this category we can distinguish two paradigms,



*Message Passing Interface* (MPI) and *Message Queuing* (MQ). The former being the de-facto standard in HPC, with minimal overhead, low level minimal abstraction and no fault tolerance; the latter, instead, offers queues to store intermediate messages still not received and it offers fault tolerance.

By following this distinction, we list some of the most common communication framework and transports used in HPC and BDA environments, as also reported by [2]. However, we only give a brief description of each one of them, since the main aim of this chapter is to show the plethora of transports which can be found in the wild and which need to be considered when developing a communication related framework or library.

**Java RMI** - Java RMI system [10] is the de-facto standard in remote communication between Java Virtual Machines, since it provides an efficient mechanism for method invocations on Java objects residing in different machines. Its communication protocol is based on TCP/IP connections and is strictly integrated with the Java environment [11].

**MPI** - The Message Passing Interface (MPI) [12] is a *message-passing library interface specification*, it describes requirements targeting the message passing parallel programming model by providing a set of function specification that must be implemented to provide communication functionalities between processes. The main MPI's aim is standardization of message-passing functionalities to improve scalability, portability and ease-of-use. It is the de-facto standard for communications among processes in distributed systems, and various implementation of the specification exists, like OpenMPI [13] and MPICH [14].

**ZeroMQ** - ZeroMQ [15] is a communication library built on top of sockets which supports in-process, inter-process, TCP and multicast transports. It allows different communication patterns, like pub-sub, request-reply fan-out and offers an asynchronous I/O model for scalable multicore applications. It was the framework used in the first versions of FastFlow's distributed runtime.

## 1.1 Towards Exascale computing

As illustrated in [1], the path towards exascale computing requires leveraging the high-performance fabric networking interface offered by the various HPC systems. Thus, frameworks have emerged in order to completely leverage the functionality

and performance requirements of customized High Performance Computing system's network hardware. We present here some of the frameworks which addressed this issue in recent years.

**Unified Communication X (UCX)** - UCX [16, 17] is a network API framework for high throughput computing targeting modern interconnects with massive parallelism. It is designed to provide a set of interfaces for implementing multiple programming model libraries that are portable, scalable and efficient.

**OpenFabrics Interface (OFI)** - OFI [18, 19] is a collection of libraries, specifically designed to target scalability requirements of HPC systems, focused on exporting communication services to applications. The main aim of the set of libraries in OFI is to expose an interface to application programmers of the underlying network fabrics, in order to fulfill HPC users' needs.

**Mercury** - Mercury [4, 5] is a framework implementing RPC and specifically designed for use in HPC environments. It offers an abstracted network API whose implementation is provided by a set of *plug-ins*. It allows asynchronous execution of remote procedures and also offers Remote Memory Access (RMA) whenever the underlying transport fabric supports it.

We have found Mercury's simple interface and broad support for various protocols particularly interesting, so we decided to integrate it as a communication framework for the distributed runtime of FastFlow. Moreover, Mercury's compatibility with multiple frameworks implementing high-performance HPC-related transports, allows maintaining portability and scalability over the heterogeneous environment of HPC network ecosystem.

# Chapter 2

## FastFlow

FastFlow is the result of a joint effort of the Parallel Programming Model Group of the University of Pisa and the Parallel Computing Research Group of the University of Turin. FastFlow is a C++ structured parallel programming framework, that leverages a streaming data-flow approach and originally targeted cache-coherent shared-memory architectures [7]. It was recently extended with a distributed-memory runtime in order to allow developers to build and deploy FastFlow applications over distributed environments. The FastFlow library was built with a layered design by keeping efficiency in the base mechanisms to maintain efficiency across the whole framework [8]. Each of the implemented layers provides abstractions at different levels in order to help both the application and the RTS developer. The top level abstractions provide classical parallel patterns targeting the application developer, the lowest layers instead expose abstractions for the RTS programmer, by providing a reduced set of **Building Blocks** (BBs from now on), that can be used to efficiently implement most of the existing parallel applications [8].

### 2.1 Building Blocks

**Building Blocks** are the core elements of FastFlow programming, since they are used as the basic abstraction layer to properly build FastFlow streaming parallel patterns. BBs come in two flavor, *sequential* and *parallel* and they both rely over Single-Producer Single-Consumer (SPSC) lock-free queues to share memory and allow efficient communication.

The implemented building-blocks are the *pipeline*, *farm* and the *all-to-all*. The *pipeline* building block has a double usage, in fact it can be used to connect build-

ing blocks and to build pipeline parallelism. The *farm* is used to express functional replication which is controlled both in input and output by **Emitter** and **Collector** entities which can be specialized in order to implement different dispatching/gathering policies. Finally, the *all-to-all* building block offers the possibility to implement the shuffle communication pattern and to remove bottlenecks coming from centralized entities in the *farm* building block.

The main characteristic of **building blocks** is their composition-oriented nature, which allows the programmer to assemble them in a LEGO-style fashion in order to build complex parallel applications. This reflects the style of the structured parallel programming methodology and enables a more organized way to build parallel applications.

## 2.2 FastFlow distributed runtime

Since the last version of FastFlow, a distributed runtime has been added which allows for FastFlow applications to be organized and deployed over distributed computing nodes via a set of wrappers class and renewed **BB** components. The introduced mechanisms for the distributed runtime allows the programmer to easily and efficiently divide and deploy an existing application into *distributed groups*, as well as creating a new application directly for the distributed environment retaining the same LEGO-style approach for parallel application on shared-memory systems.

The distributed extension for the FastFlow library is provided by modification at the **BB** level in order to provide the easy porting of existing shared-memory FastFlow application into the distributed environment without dealing with the burden of handling distributed communications; provide a set of mechanisms to implement higher-level ready-to-use distributed patterns.

### 2.2.1 Transition to distributed-memory

The transition from shared-memory to distributed-memory required the creation of new concepts, as the one of *distributed groups*, as well as new abstractions, as the classes for serialization and communication of the task to be sent across the network. Existing applications can be deployed in the distributed environment by retaining the semantic of the original FastFlow application. The only modifications required for an application programmer to transform a shared-memory FastFlow application into a mixed shared/distributed-memory application are related to the identification

of disjoint *distributed groups* via the provided API calls and the mapping of each group to an endpoint by means of a JSON configuration file.

The distributed groups (*dgroups* henceforth) represent logical partitions of the original FastFlow application which are executed on remote nodes. Each one of the *dgroups* is realized as a shared-memory application to fully exploit the potential of the single nodes. Moreover, the resulting *dgroup* is internally plugged with specific sequential BBs and node wrappers to realize serialization/deserialization and communication between the deployed *dgroups* by retaining the original application semantic.

# Chapter 3

## Mercury

Mercury [4, 5] is an asynchronous RPC framework purposefully built to efficiently provide communication services to HPC systems with high-performance fabrics. It provides an abstracted network implementation to enable transparent support to future systems and protocols, efficient use of existing native transport mechanisms, and support of large data arguments via the RDMA enabled interface. Moreover, it provides an asynchronous RPC interface, based on a callback system, that allows transfer of parameters and procedure calls in the context of both local and remote execution in order to completely remove the differentiation of communications between on and off-node.

The library allows the execution of arbitrary functions via a system of function tagging and callbacks, coupled to a queue system that saves procedure calls which are pending and still not executed. The receiving process drives the progress loop in which arguments are retrieved, function calls are executed, and the results are sent back to the origin node.

Additionally, the library offers the possibility of being ported to various systems since the network layer is abstracted and the application interface is based on a simple set of network primitives for both point-to-point messaging and one-sided RDMA. The network abstraction layer functionalities are implemented on top of different plugins. The plugin system can be considered as a *compatibility* layer, which implements the communication functionalities offered by different protocols. Examples of plugins are Libfabric [20, 21], MPI [22] and UCX [23], and multiple protocols are offered by each of the plugin, ranging from vendor-specific protocols to more classic ones like TCP and UDP.

## 3.1 RPC: a Mercury’s perspective

General RPC frameworks provide the possibility to serialize function parameters and ship them to a remote node that will execute the respective function call. As stated in [5], Mercury tries to address two main problems:

- inability to take advantage of HPC high performance communication protocols: standard frameworks are usually designed on top of TCP/IP protocols, which represent a limitation over the performances often required by HPC systems;
- inability to transfer large amount of data: standard RPC frameworks doesn’t allow (or discourage) transfer of large amount of data through the implemented mechanisms.

Mercury addresses these limitations by offering an asynchronous and flexible RPC interface specifically tailored for HPC systems. To do this, Mercury exposes an abstracted API to perform asynchronous RPC as well as large data transfer, and completely relies on the underlying network implementation provided by the plugins, in order to be independent from the used transport mechanism. As stated in [5], Mercury’s main purpose is to serve as a basis for higher-level frameworks that need to remotely exchange data in a distributed environment, by offering a flexible interface completely decoupled from the underlying protocol and system specifications.

## 3.2 Architecture

Mercury is built on a two-layer architecture, as shown in **Figure 3.2.1**, where each layer provides an abstraction for specific functionalities. The lowermost layer offers abstractions needed to provide networking functionalities such as point-to-point messaging, address lookup, remote memory access, progress and cancellation. The uppermost layer is further divided in two service-level components, referred to as *RPC interface* and *bulk data interface*. The RPC interface allows the programmer to remotely execute function calls, shipping function arguments and receiving results from the remote node; the bulk interface, instead, complements the RPC interface and allows large data transfer via the creation of memory descriptors, which enables the possibility to initiate raw memory transfers using remote memory access, whenever the underlying system provides it.

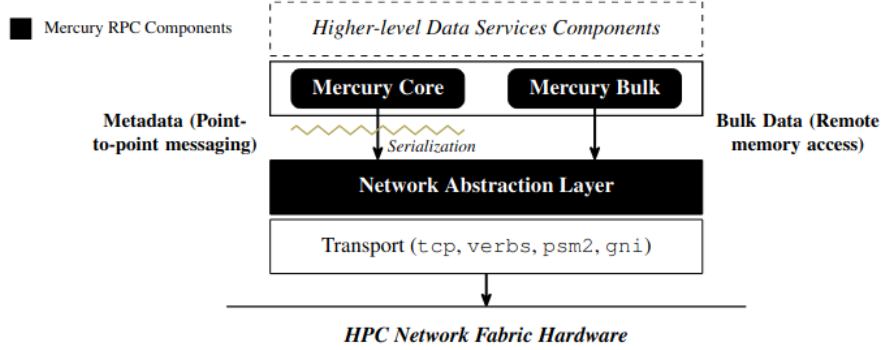


Figure 3.2.1: Mercury RPC components in the software stack.

Mercury extends the functionalities of another project, called *I/O Forwarding Scalability Layer* (IOFSL)[24], which allows RPC calls specifically related to file-system-specific I/O operations. By extending this layer, Mercury allows to generate RPC calls to generic functions that can be dynamically defined and registered in the application using Mercury.

### 3.2.1 Plugins

The plugin functionalities are referred to as a “*support for various network protocols that can be easily added and selected at runtime*” [25]. Given that Mercury’s main aim is to leverage the high performance solutions provided by HPC network fabrics, which requires specific low-level vendor APIs, Mercury relies on plugins as an intermediate layer for network functionalities. In order to overcome the burden of implementing the network abstraction layer directly on top of those APIs, Mercury relies over various plugins for the implementation of functionalities like RDMA and point-to-point messaging. A multitude of plugins are provided by the framework, however the vast majority of them is under testing or being deprecated, leaving as “stable” only the ones related to Libfabric, UCX and shared-memory [26] for local nodes communication. We show in §3.5.2 the main limitations we have encountered during testing of these plugins.

Switching between various plugins is very simple and can be done by specifying a predefined string containing the desired plugin paired with the needed protocol to use during communications. Each plugin defines its own format, but common fields are shared among the configuration strings of various plugins, and they mostly refer



to the type of plugin and the protocol to be used. Format of strings is provided in **Table 3.1**.

Plugin	Protocol	Initialization format
ofi	tcp	ofi+tcp[://<hostname,IP,interface name>:<port>]
	verbs	ofi+verbs[://[domain/]<hostname,IP,interface name>:<port>]
	psm2	ofi+psm2
	gni	ofi+gni[://<hostname,IP,interface name>]
ucx	all	ucx+all[://[net_device/]<hostname,IP,interface name>:<port>]
	tcp	ucx+tcp[://[net_device/]<hostname,IP,interface name>:<port>]
	rc,ud	ucx+<rc,ud>[://[net_device/]<hostname,IP,interface name>:<port>]
na	sm	na+sm[://<shm_prefix>]
mpi	dynamic, static	mpi+<dynamic, static>

Table 3.1: Mercury plugin’s initialization format.

Note, however, that plugins may behave differently regarding how the format string is provided, in fact plugins do not manage incomplete or incorrect configuration strings uniformly. Some of these issues are reported in §3.5.2.

### 3.2.2 Network Abstraction Layer

The Network Abstraction Layer (NAL) provides an abstraction of the network infrastructure above which the communications are executed. It is a simple abstraction which only provides limited functionalities, like address lookup, point-to-point messaging, remote memory access, progress, and cancellation.

The abstractions provided by this layer allows the uppermost layers to be completely agnostic of the underlying communication protocol implemented via the plugin system. Moreover, The API is non-blocking and uses a callback mechanism to provide asynchronous execution. Progress is driven by API calls which allows user callbacks to be placed in completion queue and retrieved for execution.

Mercury refers to communicating nodes as **origin** and **target**, indicating respectively the node issuing the request and the node receiving it. Both origin and

target nodes must specify the desired plugin/protocol pair at initialization phase, by providing a string as described in §3.2.1. Since a node can both provide and ask for services, the only time a server-specific behaviour is defined only refers to initialization, where the user can specify if the current node will be listening for incoming RPCs. Besides this, no more server/client concepts are used in the following.

The functionalities offered by the NAL refers to three main mechanisms:

- expected messages: requires a *receive* operation to be pre-posted by the target. Therefore, this requires the origin node to be known in advance, before the receive operation is posted. If the receive operation is not posted before the message is sent, it can be dropped;
- unexpected messages: does not require the target to post a receive operation for the message, and they can arrive from any source. The target can retrieve received messages in an asynchronous way. These messages are allowed to be dropped, but the plugin can decide to queue them anyway;
- remote memory access: allows registration of memory chunks which can be later accessed by target nodes. Abstractions are provided through API which contains operations generally provided by most RDMA protocols.

The network abstraction is designed to allow emulation of one-sided operations, such as RDMA, on top of two-sided operations. In this way, Mercury can easily adapt to protocols which only supports fixed operations, like TCP/IP ones, where one-sided communications are not possible.

### 3.2.3 RPC Layer

The RPC layer allows nodes to issue remote calls, and it is based on the messaging model described in §3.2.2. RPC requests are based on the common knowledge of origin and target nodes on how to encode and decode function parameters and return values.

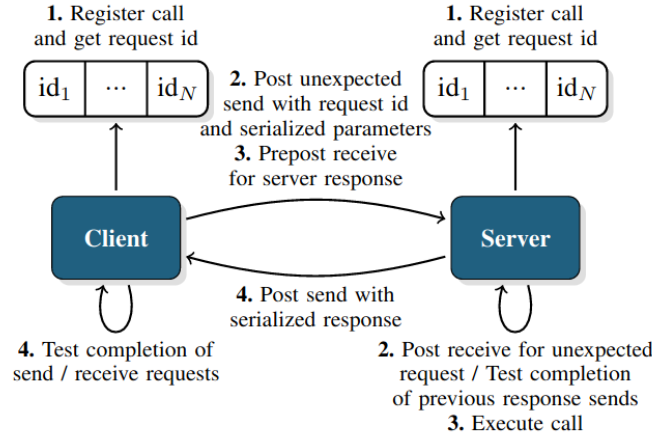


Figure 3.2.2: Execution flow of RPC call.

Mercury provides mechanisms to support a set of generic function calls avoiding hard-coded routines. To allow this, origin and target nodes must register a unique function name along with encoding and decoding routines by using shared input/output types. We postpone the precise description of the registration process to §5.1, since Margo introduces few additional steps, which are the ones that are actually used to implement FastFlow’s communication classes. We show in **Figure 3.2.2** a usual flow of execution needed by both origin and target nodes to execute a RPC request. The registered function is mapped to an ID which will be used in all the communications between the two nodes. A further step is needed by the target, which must register a callback that will be executed every time that ID is received. Once the functions are registered, the origin sends an *unexpected* message to the target.

Two situation may occur during an RPC request, and different mechanisms are used to guarantee full asynchrony:

- the RPC call expects a response from the target: the origin node prepares its memory buffer to receive the response and uses the buffer to pre-posts an *expected receive*. At the reception of the response, the origin node can retrieve the response from the callback queue and proceed;
- the RPC call does not expect a response from the target: the origin node, at the moment of RPC registration, declares that this RPC does not expect a response. An RPC of this kind allows the origin node to proceed without

posting a receive operation, progress can be made as soon as the request is sent to the target.

From the target side, receiving an *unexpected* message with a specified ID translates in the execution of the callback registered at startup by decoding the parameters sent by the origin, and sending the outputs by encoding them at the end of execution, if the registered RPC expects a response.

### 3.2.4 Bulk Layer

This layer allows to send large data by avoiding intermediate memory copies. It is performed by creating a local memory handle which points to previously registered areas of memory (not necessarily contiguous) which the target node can access via RDMA operations. The bulk layer is directly built on top of the RDMA interface defined in the network abstraction layer.

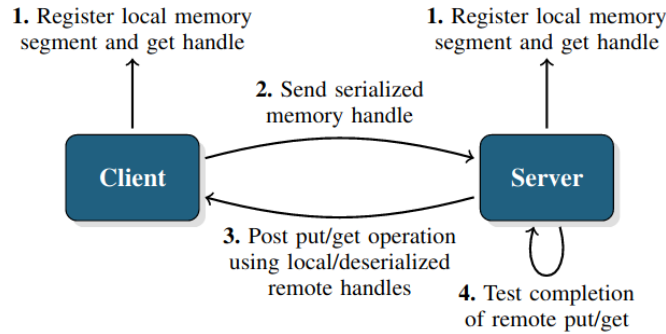


Figure 3.2.3: Execution flow of bulk request.

A typical execution flow for a bulk request is depicted in **Figure 3.2.3**. The target node manages all the bulk transfers in order to be able to control the data flow and protect its memory from concurrent accesses. This operation is one-sided, and it is started by the origin node which creates a bulk data descriptor, serializes it, and send the serialized descriptor to the target as an argument via a specific RPC request. Once the descriptor has been deserialized by the target, two situations can occur. The request can be related to *consumption* or *production of data*. In both cases the target node creates a memory handle to manage the request, allocating the necessary memory. However, in the first case the target initiates a remote **read** operation before performing the function call, in the second case the target executes

the call and initiates a remote **write** operation with the produced results.

Memory handles are an important building block for memory transfer, in particular in such cases where non-contiguous memory is to be transferred. The memory of the communicating nodes is abstracted by the memory handle and allows to access memory in a transparent way without modifications to the process described above.

### 3.3 Resilience and Fault Tolerance

The fault tolerance mechanism is mainly provided by the possibility to interrupt calls and reclaim resources of pending operations after these have been signaled as *cancelled*. Cancellation is an asynchronous and local operation. From a user perspective, completion of offloaded operations is known only when the associated callback is placed in the local queue of pending operations. When a callback is triggered and the operation was locally canceled by the user, that operation is reported as canceled and aborted.

A basic flow of events related to the cancellation of an operation is shown in **Figure 3.3.1**.

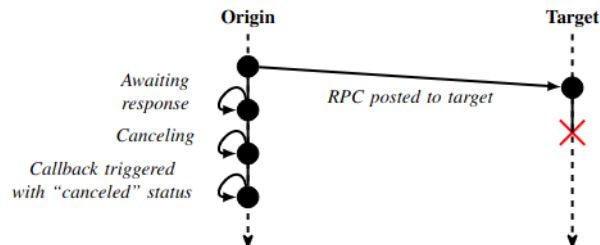


Figure 3.3.1: Cancellation of an RPC operation.

### 3.4 Dependencies

As reported by [27], the dependencies are mostly related to the intended plugin to integrate inside a specific application. No further requirements are specified by the Mercury’s developers as needed for the installation and use of their framework.

## 3.5 Mercury analysis

Mercury comes with lots of functionalities and drawbacks. We list in this section all the advantages and the main limitations of the Mercury framework encountered during the testing phase, which are therefore shared by all the higher-level libraries which are based on it.

### 3.5.1 Advantages

Mercury’s interface allows to easily enable communication using a system of plugin which offers compatibility with a multitude of network and vendor specific protocols. The Mercury library, once again, shows the importance of having an abstraction layer to provide efficient use of underlying protocols without having a complete expertise on how they work. Being able to address different needs and to leverage the functionalities of different systems, by the means of a general API, makes implementing communication on such systems painless and less error prone. Moreover, Mercury’s capability of handling different protocols with no code changes at all allows portability of applications to new systems which may provide only vendor specific protocols or which does not offer support to common transports such as TCP or MPI.

### 3.5.2 Limitations

Mercury, however, is not exempt from problems and limitations. In this section we gathered all the main limitations that are known and which we discovered during the testing phase of the presented framework. As specified in §3.2.1, MPI plugin is deprecated and not maintained [28]. In fact, Mercury’s developers suggest the use of libfabric plugin to leverage efficient HPC communication mechanisms, and UCX plugin as a general purpose plugin. As specified in [25], MPI is considered to be present in almost all systems, and the functionality offered by the NAL are only intended for prototyping and testing, this suggesting that in case an MPI implementation is needed, one shouldn’t rely on Mercury’s interface.

Further limitations found during the testing of this library are mostly related to the limited support of the Libfabric plugin paired with TCP transport [29], which is intended to use for testing and debugging applications on machines that do not provide high performance fabric protocols [30]. Problems about this specific transport were already reported [30, 31, 32], and during testing the main that we noticed are:

- During reply phase of RPC request using `ofi+tcp`, information about public

IP of origin node are not used. This resulting in a lost packet in case the origin node was behind NAT, probably due to the *emulated* source addressing, as stated in [21]. The same does not happen by using `ofi+sockets`, which use is discouraged as per [29];

- Various problems related to termination of nodes which issued a request that couldn't be fulfilled, as also reported in [31];
- As per [33], the sockets provider has been deprecated in favor of the tcp one.

Other plugins, such as UCX [34], showed some problem in connecting nodes situated in different networks, mostly due to the fact that the plugin provided by Mercury struggles to bind the socket address to the specified one. However, it is able to establish a connection between endpoints situated in the same subnetwork.

Given the requirements for the FastFlow distributed version, the internal limitations of the Mercury framework could be too tight to allow its extensive utilization as a final communication library. However, the provided plugins and functionalities could anyway serve as a prototyping library in order to develop FastFlow communication layer APIs, in order to painlessly switch to a different communication framework in the future.

In fact, as we show in §5, the composition of the Mercury framework with higher-level abstractions allows to easily develop communication nodes to handle multiple endpoints at once without struggling with the explicit progress mechanism provided by Mercury. For this reason, in the next sections we introduce the libraries used by the higher-level abstractions which use Mercury framework as a core building block to provide RPC functionalities in an automatized way.

We show, for each of the analyzed libraries and frameworks, which are their advantages and drawbacks, with special focus on lowering the amount of dependencies and minimizing the knowledge of the communication layer about the data that are exchanged between nodes, in order to create communicating nodes completely agnostic about the types used during computation, allowing also for parallelization of the serialization process and reducing the computation needed at the “edge” of each group.

Besides the limitations we have analyzed, which are mostly related to technical aspects and not on the usability of the presented frameworks, the simplicity of creating communication nodes and connecting them via a high-level communication library

such as Margo, without requiring too many modification in the original application code, can be a good driver to experiment and develop an initial set of API calls which will allow to extend the communication layer functionalities and plug in a broader set of protocols which are still not supported by the used frameworks. Having an independent API is a natural and important step to remove strict dependencies from each of the technologies used to implement the communication functionalities, allowing extensibility with little to no changes in existing code. Besides the limitations we have found during the experimentation with the analyzed libraries, our analysis on those frameworks was anyway important to understand the power and utility of having a high level abstraction on top of very specific APIs which requires a deep understanding of how the underlying infrastructure works. Hence, it is conceptually interesting and it could be very helpful in the implementation of the final communication layer to retain some of the concepts introduced by the *Mochi core* building blocks.



# Chapter 4

## Argobots

Argobots [6, 35] is a lightweight low-level threading framework, which offers a portable library interface and allows specialized runtime management to the user. The main aim of Argobots is to provide a mapping between high-level abstraction to low level implementations [6], as well as offering a lightweight layer of execution. To this aim, Argobots implements lightweight parallel work units, such as user-level threads (ULT) and tasklets, which are non-preemptable and offers, respectively, different models of execution. Work units are executed by OS-level threads, which in this context are referred to as Execution Streams (ESs). Each ES can be associated to a set of *pools*, which are containers of work units, and execute tasks in the order provided by *scheduler* entities. Argobots allows to define various scheduler (included custom ones), which determine the order of execution of each work unit inside the pools. Scheduler entities can be “plugged” at runtime to change the strategy of execution, based on requirements adapting to the computation at hand. Moreover, work units can be dynamically moved to a different pool in order to allow computation on a different ES.

A generic Argobots application is depicted in **Figure 4.0.1**, which is built by following the execution flow shown in **Figure 4.0.2**. We can see how different pools can be associated to the same ES and how the scheduler provides the work units to the execution flow, which is sequential and guarantees progress. We describe in the next sections the characteristics of each building block of an Argobots application.

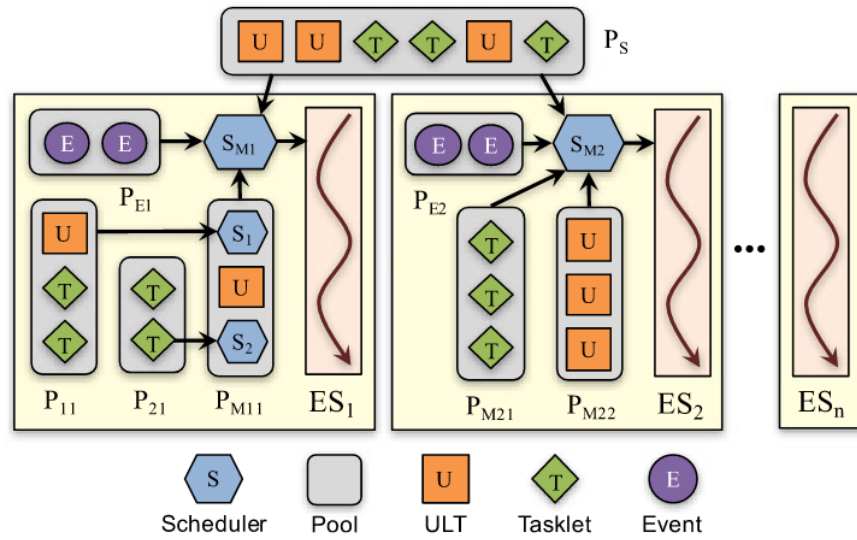


Figure 4.0.1: Argobots execution model.

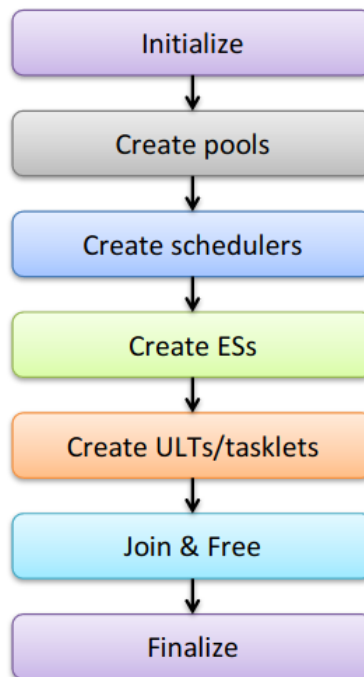


Figure 4.0.2: Argobots application flow of execution.

## 4.1 Building blocks

### 4.1.1 Execution Streams

Execution streams represent a sequential and independent instruction stream, which can consist of one or multiple non-preemptive work units (namely, ULTs and tasklets), it is mapped to a Pthread and can be bound to a hardware processing element such as CPU core or hardware thread. The work units to be executed are contained in one (or more) pools, and they are retrieved from one (or more) schedulers. Scheduling policies can be tweaked during runtime, and they determine the order in which work units flow to the execution stream and are thereby executed sequentially.

Given the sequential nature of an ES, work units which are executed by the same ES do not require expensive synchronization mechanisms, since they do not run concurrently. However, synchronization is needed between work units executed in different ES running in parallel.

### 4.1.2 Schedulers

Schedulers are responsible to deliver work units to the ES they are associated with, and they can implement different scheduling policies, which determines how work units are retrieved from pools and handled to the ES. Schedulers can be composed and scheduled just like work units, in order to compose different scheduling approaches dynamically, related to the computation at hand. Since a scheduler is considered as a work unit, it can be inserted into a pool and executed, which translates in a change of scheduling policy for the ES in charge, which will shift back to its original scheduler when the computation associated with the current scheduler is over.

### 4.1.3 ULTs and Tasklets

User-level threads and tasklets, also called *work units*, represent two different types of workflow. User-level threads are independent execution units, which are executed in user space and can yield control to the scheduler, they can be dynamically migrated on a different execution stream during execution, since they have a private stack and context-saving capabilities. Tasklets, on the other hand, are considered to be “lighter” than ULTs, since they do not incur in costs related to context saving and stack management. They should be considered as atomic units of execution, which can’t yield to a different execution and run to completion without context switching or suspension.

#### 4.1.4 Pools

Pools are containers of work units, and their sole purpose is to act as a uniform way of associating work to ES, indirectly, through schedulers. The only property associated with pools is the one regarding access, which can be set as private or shared. Shared pools may be used to implement a “work-stealing” strategy between two or more ESs. Pools are also used internally by the ES to receive asynchronous events.

A pool associated with a running or stacked scheduler defines a set of work units ready to execute, which have to be controlled by the application programmer, since Argobots does not implicitly define dependencies between work units. Hence, synchronization mechanisms offered by the library must be used to control the control flow of different work units. Various operations can be performed over pools, like migration of work units to or from a specified pool, creation, destruction, and of course pop and push operations.

## 4.2 Operations

Argobots defines a set of operations which are common to all work units, made exception for some of them which are not available to Tasklets, given their execution model limitations. Implemented operations refer to:

- Creation: allows the creation of work units, which are then inserted in a specific pool in a ready state. The type of pool and the associated scheduler(s) will determine the time and context of execution;
- Join: work units can be joined by other ULTs, which wait for their termination;
- Yield: a work unit can cooperatively yield control to the scheduler which is executing in the current ES at the time of yielding. The scheduler will then gather the next work unit to be executed following the implemented policy;
- Yield.to: a ULT can decide to yield control to a specific ULT in the same ES. This operation avoids the overhead of one context switch by bypassing the scheduler;
- Migration: allows work units to be migrated between pools;
- Synchronization: ULT can rely over usual synchronization mechanisms implemented by Argobots, such as mutexes, condition variables etc.

## 4.3 Memory Management

Argobots is intended to be used in fine-grained dynamic environments, where creation and destruction of work units, as well as context switch between them, take place at high frequencies.

Since, as shown in [6], memory allocation and deallocation contribute to most of the time required to create and destroy a work unit, Argobots developed its own memory allocator. The memory allocator creates, for each ES, a memory pool which is subsequently used to manage all creation and destruction calls of work units. The memory pool is held private per ES, in order to avoid heap access synchronization upon creation of work units by means of the same ES. The size of the pool is tied to the number of work units which are spawned and allows to return memory to the system upon their destruction, if the pool reached a given threshold.

Context switching, on the other hand, has a lower cost in terms of memory, compared to creation and destruction of work units, and it is only related to ULTs, since tasklets are atomic operations and can't yield to other work units. Context switching costs are mostly associated with:

- ULT suspension: context of the currently running ULT must be saved and context of the next ULT must be resumed. However, the first step can be avoided if the yielding ULT is terminating and will not be resumed later;
- ULT join: the join operation suspends the caller and yield control to the scheduler, until the joined ULT terminates its execution. However, if two ULTs are in the same ES, the joining ULT can directly yield control to the joined ULT, bypassing the context switch to the scheduler entity.

# Chapter 5

## Margo

Margo [4] is a Mercury binding that uses Argobots as a runtime library. Just like Mercury, Margo requires both origin and target nodes to register for RPC calls associated with a specific identifier (a string) and an input/output type, which must necessarily be registered using the relative macros, as we describe in the following. By using Argobots as a runtime library, Margo hides the handling of Mercury's progress loop by delegating it to a specific Argobots ES (namely, a physical thread). Moreover, Margo allows to freely manage how the various RPC associated callbacks will be dispatched among internal pools and external user defined pools. When initializing Margo, the user can specify its own pools and ESs that will be responsible for both progress loop and RPC calls, otherwise the programmer can decide to completely let Margo handle the calls with its own pools and ESs.

Being a binding between Mercury and Argobots, Margo retains all the concepts we defined in the previous section, however, Margo greatly simplifies the development of RPC-based services by introducing extensions to the Mercury-based model:

- more intuitive communication: Margo defines wrappers based on Argobots runtime in order to present asynchronous communication mechanisms as a user-level thread communication model. Communication facilities are handled by ULTs, which can be suspended and resumed as communication proceeds;
- progress loop abstraction: polling and communication events are internally managed by a ULT. Policies for handling RPC-related execution can be manually defined and tweaked to better suit the application needs. In this way, multiple providers can be easily handled and multiplexed by the same core process, without the need of handling multiple progress loops at once;

- renewed polling strategy: Margo allows both busy and idle mode, which allows to tweak Margo for both performance and resource consumption needs.

## 5.1 RPC registration process

In this section we describe the RPC registration process, which follows a standardized procedure and requires RPC input and output types, as well as routines to pack the types into network buffers, to be defined by both origin and target nodes. Additionally, the callback function to be executed must be defined by the target node with a predefined signature in order for Margo to encapsulate it in the internal service functions which triggers RPC execution upon receiving the request from the origin node.

The registration process must proceed by:

- defining input and output types: as shown in **Listing 5.1.1**, RPC arguments types must be defined by the means of a C-style `struct` datatype, which can internally contain all the necessary types for the RPC call;
- defining the packing routine: **Listing 5.1.2** defines the routine which is internally used by Margo to correctly write to and read from the network buffers which are then used to send and receive RPC data arguments. Margo also offers a “simplified” way of defining such routines, which we omit for sake of presentation and to better explain how the internal network buffers are built and filled with data. Each packing routine is strictly tied to a specific type. Having a type `X`, the routine must be defined as `hg_proc_X`, and it must encode and decode each of the field contained in the `X` struct type. Additionally, Margo can be built with XDR capabilities which will in turn change the way data are represented in the buffer exchanged between nodes, however, since FastFlow distributed version shares already serialized streams, we don’t rely upon this functionality.
- defining the actual RPC callback: this step is only required by the node which will act as a “server” for the specific RPC. In **Listing 5.1.3** we provide a sample RPC callback declaration and definition, which allows the server node to define the function to be executed upon an RPC request from another node. RPC callback declaration and definition require also a call to specific macros defined by Margo, which will wrap the RPC callback defined by the user with

Argobots-aware code, necessary for Margo to dispatch RPC callbacks among different ULTs.

- associating an ID to the RPC: after having registered all the types and routines, both origin and target nodes have to associate a common ID for the RPC that they intend to use. This can be done as shown in Listings **5.1.4** and **5.1.5**. The registering process is complete, and the origin node can now use the returned RPC identifier to issue requests to the listening target node by using a `margo_forward` call. Note that, after the registration process is complete, the listening node can receive RPCs calls with the specified ID from all the origin nodes which registered an RPC with the same ID.



```

1 typedef struct {
2     hg_int64_t    hash_val;
3     hg_uint64_t   size;
4     hg_bulk_t     bulk_handle;
5 } ff_rpc_in_t;

```

Listing 5.1.1: RPC type definition.

```

1 hg_return_t
2 hg_proc_ff_rpc_in_bulk_t(hg_proc_t proc
3     ,
4     void* data) {
5     ...
6     // Retrieve input data structure
7     ff_rpc_in_t* in = (ff_rpc_in_t*)
8     data;
9
10    // Write/read data to/from Margo's
11    // send/receive buffer carried by '
12    proc'
13    hg_proc_hg_int64_t(proc,
14        &in->hash_val);
15
16    // We call a specific routine for
17    // each
18    // field in the RPC input struct
19    hg_proc_hg_uint64_t(proc,
20        &in->size);
21
22    hg_proc_hg_bulk_t(proc,
23        &in->bulk_handle);
24
25    ...
26 }

```

Listing 5.1.2: Packing routine definition. Allows the Margo framework to manage data from/to the network buffer used internally to ship data during RPC calls. Each of the type-specific routines allows data copies which are aware of the size of data to be packed/unpacked.

```

1 void ff_rpc(hg_handle_t handle);
2 DECLARE_MARGO_RPC_HANDLER(ff_rpc);
3
4 void ff_rpc(hg_handle_t handle) {
5     ff_rpc_in_t      in;
6     const struct hg_info* hgi;
7     margo_instance_id mid;
8
9     // Get input data
10    margo_get_input(handle, &in);
11
12    // Retrieve objects to identify current
13    RPC
14    // and get back registered data
15    hgi = margo_get_info(handle);
16    mid = margo_hg_info_get_instance(hgi);
17
18    // Here data may be retrieved and used
19    // internally in the RPC call
20
21    margo_free_input(handle, &in);
22    return;
23 }
24 DEFINE_MARGO_RPC_HANDLER(ff_rpc)

```

Listing 5.1.3: RPC declaration and definition

```

1 int main(int argc, char** argv) {
2     // Margo initialization code
3     ...
4
5     // Initialize the Margo instance used to perform
6     // Margo-related calls
7     margo_instance_id mid;
8     mid = margo_init(listening_addr,
9     MARGO_CLIENT_MODE, 1, 1);
10    ...
11    my_rpc_id = MARGO_REGISTER(mid, "ff_rpc",
12                                ff_rpc_in_t, ff_rpc_out_t, NULL);
13
14    // Looks up for a server
15    hg_addr_t svr_addr;
16    margo_addr_lookup(mid,
17                        svr_addr_str, &svr_addr);
18
19    // Creates the handle for the RPC call
20    hg_handle_t handle;
21    margo_create(mid, svr_addr,
22                  my_rpc_shutdown_id, &handle);
23    margo_forward(handle, NULL);
24    ...
25 }

```

Listing 5.1.4: Finalization of the registration process in the origin node. The register macro specifies the input/output expected types as well as the packing routines as described above.

```

1 int main(int argc, char** argv) {
2     // Margo initialization code
3     ...
4
5     // Initialize the Margo instance used to perform
6     // Margo-related calls
7     margo_instance_id mid;
8     mid = margo_init(protocol, MARGO_SERVER_MODE, 1,
9                     1);
10    my_rpc_id = MARGO_REGISTER(mid, "ff_rpc",
11                               ff_rpc_in_t, ff_rpc_out_t, ff_rpc
12                               );
13    margo_wait_for_finalize(mid);
14 }

```

Listing 5.1.5: Finalization of the registration process in the target node. The listening node only needs to register the RPC types and routines, as the origin node, and additionally it needs to specify the RPC callback for the registered ID.

## 5.2 Margo analysis

As we pointed out in §3.5, most of the limitations are tied to the support of the various plugins and their functionalities. Margo, besides sharing all these limitations, introduce an incredibly easy interface to implement multi-homed RPC services, which greatly simplifies the development of RPC-based FastFlow nodes.

The main limitation of the Margo library are mostly related to how RPC calls are handled, in particular with reference to data copies happening between the input data and the send/receive buffers used to ship RPC data through the network. However, further investigation is needed at this point in time to analyze how those limitations can be avoided, without relying on the bulk functionalities, which however remain a suitable fallback method if data copies are unavoidable. However, a further problem is introduced in the utilization of bulk transfers via RDMA where no responses are expected from the sender, in fact in this particular case one must be extremely careful about freeing memory that has still not been read from the remote end.

Margo, on its hand, introduces additional dependencies which are related to the libraries it uses in the underlying layers and how it handles configuration strings. The introduced dependencies are the following:

- Mercury: as described in §3, in particular requires for this framework to be built with `-DMERCURY_USE_SYSTEM_BOOST:BOOL=OFF -DMERCURY_USE_BOOST_PP:BOOL=ON`, which enable pre-processors macro in Mercury in case BOOST library is not present in the system. Since we do not want to be dependend on BOOST library, this is necessary since Margo uses these macros internally;
- Argobots: as described in §4;
- json-c: A JSON implementation for C language[36], needed internally by Margo in order to generate configurations based on json-formatted strings provided at initialization of Margo instances.

# Chapter 6

## Implementation

### 6.1 Communication classes

We developed FastFlow’s communication classes by extending the `ff_node_t` class, in such a way the developed classes can be used in any context where a `ff_node_t` can be used. However, we show that they naturally sit at the extremes of a pipeline building block, since they offer functionalities to receive and forward data from and through the network.

The classes we implemented are strictly tied to the concepts of receiver and sender node, and they make use of Margo’s capabilities of developing multi-endpoint services without requiring to manually handling progress loops. At the current stage of development, there is no personalized behaviour based on the endpoint which received a call to the registered RPCs. This represents the next main step in the evolution of the implemented classes. (De)Multiplexing, however, should be fairly easy to handle thanks to the internal mechanisms that Margo offers in order to identify the origin node which issued a request. Having personalized behaviour, depending on the specific endpoint, is very important to deal with the grouping concept introduced since the distributed version of FastFlow has been developed. Particularly pathological cases, where a personalized behaviour must be implemented given the specific endpoint on which the request is received, benefit from the automatized mechanisms offered by Margo. Moreover, the simplicity in handling multiple endpoints via the implemented classes allow to test a multitude of situations with little to no code changes.

### 6.1.1 Receiver node

Receiver node is a FastFlow `ff_node_t` that relies on a Margo instance initialized with `MARGO_SERVER_MODE`. This allows the receiver node to wait for incoming RPC requests at the specified addresses. A list of addresses can be provided at initialization, translating in a receiver node listening for the same set of RPC functionalities on all of the provided endpoints. **Listing 6.1.1** shows pseudo-code of the implemented receiver node.

```
1 struct receiverStage: ff_node_t<Task> {
2     std::vector<margo_instance_id> mids // Margo instances to be used during
    servicing
3
4     // Parametrized constructor which registers all specified endpoints with the
    // provided configuration strings.
5     receiverStage(std::array<char*> endpoints, std::array<char*> configs) {
6         ...
7         for(i < num_endpoints) {
8             // init_endpoint initializes a Margo server instance on the provided
            address
9             // and config string. The endpoint allocates a specific ES and pool
            to handle
10            // all the requests coming through the specified address.
11            // Returns an handle to the created Margo instance
12            mids[i] = init_endpoint(endpoints[i], configs[i]);
13
14            // Code to get the end address (to be forwarded and used by the
            sender node)
15            // and various debugging can be put here, by using the mids generate
            by the init
16            // phase
17            ...
18
19            // Use the initialized mids to register the set of RPCs that will be
            offered by
20            // this receiver node.
21            // Responses for the registered RPCs are disabled.
22            register_service(mids[i])
23        }
24    }
25
26    // The receiver node has nothing to do in the service method, it can simply
    wait
27    // for termination of all the listening endpoints and forward the EOS at the
    end.
28
29
```

```

30 // 'task' parameter is ignored in this stage, since it acts as a sort of '
    endo-stream'
31 // generator for the other stages in the pipe.
32 Task* svc(Task* task) {
33     wait_for_finalize(mids);
34     return EOS;
35 }
36 }

```

Listing 6.1.1: Receiver node pseudo-code.

### 6.1.2 Sender node

Sender node, on the other hand, relies over Margo instance initialized with `MARGO_CLIENT_MODE`, since we do not expect this node to be listening on any RPC request for now. The current version of the sender node allows to contact only a single endpoint at a given address. The address to contact for this specific node must be provided at initialization, and at the moment it is not possible to change the remote service address. In the next versions the possibility to contact multiple addresses, in order to implement the pathological grouping case described before, will be implemented following the same approach used for the receiver stage. Pseudo-code for this class is provided in **Listing 6.1.2**.

```

1 struct senderStage: ff_node_t<Task> {
2     margo_instance_id    mid;        // Margo object to use for communications
3     hg_addr_t            svr_addr;    // Server address listening for incoming
    RPCs
4
5     senderStage(char* addr) {
6         // Initialize the Margo instance as a client node and creates the
        necessary
7         // objects to register RPC calls and addresses to issue the requests to
        the
8         // specified address.
9         mid = init_endpoint(addr);
10
11         // This is the exact same function that is called by the receiver, but in
        this
12         // case we don't need to provide implementation of the actual RPC
        function.
13         register_service(mid);
14     }
15
16     Task* svc(Task* task) {

```



```

17     // The task received by the previous stage is packed into the RPC
    registered
18     // type and forwarded to the sender via an RPC call. The current node can
19     // proceed since no response is expected.
20     ff_rpc_in in = pack_task(task);
21     forward_task(mid, in);
22     return GO_ON;
23 }
24
25 void svc_end() {
26     // Upon termination forwards a shutdown RPC to the connected receiver
    node.
27     // The sender has nothing to do more than cleaning Margo resources.
28     forward_shutdown(mid, shutdown_id);
29     finalize(mid);
30 }
31 };

```

Listing 6.1.2: Sender node pseudo-code.

## 6.2 RPC based service

The developed classes allow distributed groups to communicate through a fixed set of RPC calls, which are shared among receiver and sender nodes and used for all the communications between the various groups. Two RPCs are offered by the communication service, and they refer to:

- **ff\_rpc**: this RPC is used during the whole lifetime of the FastFlow application, and it is used by the sender node, upon reception of a stream element, to forward the current stream element to the receiver node it is connected to. The sender node packs the data in the RPC input type, as described in §5.1, and ships them issuing a forward call to the connected endpoint. Since no response is expected from the RPC call, the sender can proceed immediately.
- **ff\_rpc\_shutdown**: this RPC is only used upon reception of an EOS object from the stream. In this case no data needs to be sent through the network, but only a signal that the stream has ended, in order to allow the receiving node to gracefully terminate its execution and forward EOS accordingly to the local group's nodes. Signaling an EOS propagates through the network, since the receiver node will generate an actual EOS object when the **ff\_rpc\_shutdown** callback is executed that will allow local nodes to terminate as per FastFlow's

execution flow. At that point, the next sender node will issue an end-of-stream RPC and the whole FastFlow application will terminate. Note, however, that a receiver node will not forward an EOS object unless all the listening endpoints have received a shutdown RPC.

It is important to point out that this set of RPC callbacks need to be registered only once per Margo instance. This means that, once for each listening endpoint, the receiver node must register the same set of RPCs, but it is completely agnostic on the number of nodes that will issue RPC requests on the same endpoint.

## 6.3 Splitting taxonomy

In this section we analyze the characteristics of different connections that take place between remotely connected groups, we show that very specific types can be determined whether the splitting happens in a *horizontal* or *vertical* fashion. We describe the specific categories depending on different splitting strategies. To better explain the identified categories, we introduce the concepts of *level*, *horizontal* and *vertical* splits. Considering a generic FastFlow application as a pipeline of stages, we naturally associate a *level* to each stage incrementally. In addition, we define as *horizontal* a split which creates two groups sitting at the same level of the original pipeline, instead we define as *vertical* a split which creates groups operating at two different levels of the pipeline. Note, however, that a *vertical* split can internally contain further *vertical* or *horizontal* splits. Given these definitions, we can identify two categories that describe each of the connections in remotely-connected groups, specifically:

- *internal*: elements to be sent/received are related to a connection that takes place between groups that are split horizontally. Connections of this type can be linked to splittings that divides an original FastFlow building block in two groups belonging to the same pipeline level, for example an *all-to-all* building block divided by distributing left and right workers in two groups;
- *external*: elements to be sent/received are related to connections that take place between two groups at different levels of the original pipeline. These connections are typically related to vertical splits of the original FastFlow application.

To better illustrate this taxonomy, we present a simple example which contains every concept we introduced up to now. In **Figure 6.3.1** there is a sample FastFlow

application composed of an *all-to-all* building block with two left workers and two right workers. By performing an *horizontal* “cut”, we want to split the *all-to-all* in two different groups, each of them with one left worker and one right worker. If we consider this *all-to-all* block to be in the middle of a pipeline, the groups will finally be as depicted in **Figure 6.3.2**. We show the structure of the main software entities that handle the splitting. As we can see, the amount of nodes participating in the application increases to facilitate communication between remote groups, both via *internal* connections (LS1 to RS2) and via *external* ones (ff\_receiver to LS2).

We will proceed now by describing the bigger picture by delving into particulars and explaining each of the involved parts and messy communication lines. The reasoning is very simple once the way communications are handled is clear, but it can seem pretty complicated at the start. We depicted in different colours the data flow which involve a stream element as the same “origin” object. When the colour of a communication line changes, also the provenance attached to the considered object will change. We describe in the following the different concepts depicted in the picture below:

- left/right box: they are abstracted representations of *internal* nodes in remotely connected groups. Left and right boxes of horizontally split groups, belonging to the same original building block, are logically connected by mean of their own local sender and receiver nodes. Since they respectively emulate a left and right worker, a stream element received by a left box (and hence originally shipped by a right box in a remote group) will be forwarded to a right worker in the current local group.
- red channel: carries stream elements which are part of the normal flow of execution, for example streaming elements flowing through the stage of a pipeline one after the other.
- purple/green channel: represents *internal* connection channels. They allow communication between boxes of different groups in order to maintain the *all-to-all* semantic in a distributed environment.

The connections, internally, are determined by each remote group by means of the configuration file provided upon initialization and by the information that are exchanged during the startup phase. In this way, each receiver/sender node can associate correctly the ID of each channel to the various entities that are created after the splitting is performed.

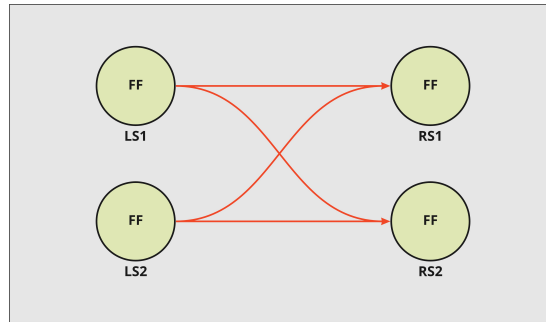


Figure 6.3.1: Sample A2A building block.

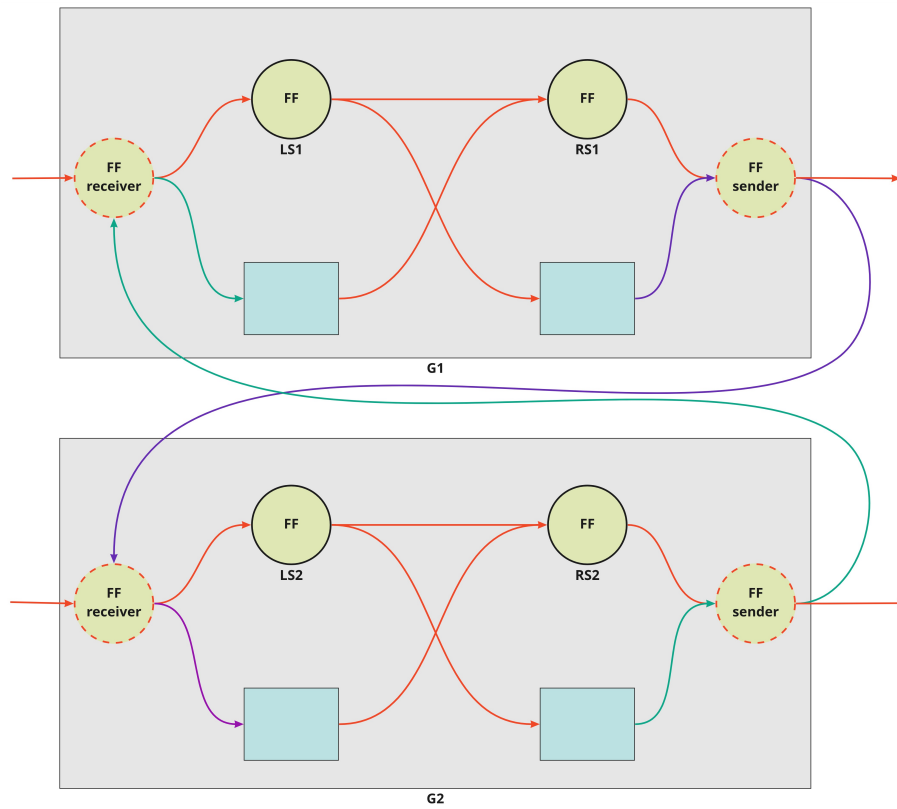


Figure 6.3.2: Pathological case representing the splitting of an A2A node into two groups which needs to communicate internally in order to retain the original semantic of the A2A building block.

# Chapter 7

## Testing

In this section we show the testing application we have developed to assess the functionalities of the implemented classes described in §6.

### 7.1 A sample application

In **Figure 7.1.2** we show a sample application developed to test the implemented classes for communication. As it can be seen, the structure of the application is very simple, we have three groups communicating in a pipeline fashion, by means of implemented `sender` and `receiver` nodes. Communication between groups can happen with the different plugins provided by the Mercury library, and switching from one plugin (transport) to another is only a matter of passing as configuration string the desired plugin (transport) to use for communication.

**Figure 7.1.1** shows the legend for the nodes and communication channels used in the sample application. We have two types of nodes, sequential ones which are standard `ff_node_t`, and Margo-injected nodes which are used for communication between distributed groups. Simple plain arrows indicates usual FastFlow shared memory channel, instead bullet-tailed arrows indicate a distributed connection between two nodes. Note, however, that “distributed” here only means that two groups are connected through Margo calls, but the nodes can reside in the same machine and communicate with the shared memory plugin provided by Margo, as it happened in part of the testing phase.

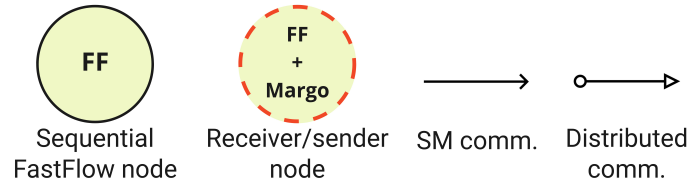


Figure 7.1.1: Node legends

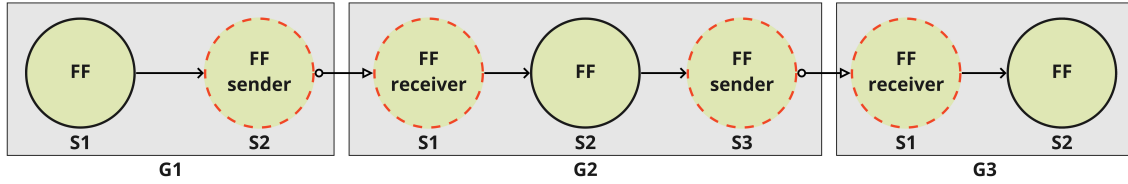


Figure 7.1.2: A sample testing application composed of three main groups which are connected by using the implemented classes using Margo library.

Moreover, the standard FastFlow nodes depicted in **Figure 7.1.2** can be any composition of standard FastFlow building blocks. We showed here, for simplicity, a single sequential `ff_node_t`. The important thing to notice is that **sender** and **receiver** nodes are mandatory, respectively as last and first nodes, in order to compose distributed groups correctly. The depicted groups represents, potentially, the three situations that may occur when splitting a standard FastFlow application in a distributed one, that are:

- (G1): group with only out remote connections. In this case only a **senderStage** is needed as a last node in the application pipeline. The pipeline is most likely to generate elements from an endo/eso-stream;
- (G2): group with both in/out remote connections. Both **receiverStage** and **senderStage** are needed, respectively as first and last nodes in the pipeline. All internal nodes receive and send stream elements using Margo-injected nodes;
- (G3): group with only in remote connections. Only a first **receiverStage** pipeline node needed.

In **Listings 7.1.1**, we show a skeleton with required library calls necessary to initialize Margo and Argobots environment, which are necessary steps before composing FastFlow building blocks with the implemented communication nodes.

```

1 int main(int argc, char** argv)
2 {
3     ...
4
5     // Setting up main Argobots instance
6     margo_set_environment(NULL);
7     ABT_init(0, NULL);
8
9     {
10         // Here we place code to build each of the groups composing
11         // the application, namely G1, G2, G3.
12     }
13
14     // Finalizing Argobots
15     ABT_finalize();
16 }

```

Listing 7.1.1: Skeleton of a sample FastFlow application using Margo as communication layer.

After having initialized the environment with the required library calls, we proceed by building the groups that have to be executed as follows:

- group G1: implemented in **Listing 7.1.2**, creates two nodes, the first one is a simple FastFlow `ff_node_t` generating a stream of elements to forward to other nodes, the second one is a `senderStage` initialized with the address of the `receiverStage` to which elements will be forwarded using the Margo communication mechanisms;
- group G2: implemented in **Listing 7.1.3**, builds three stages, where the first and the last one are communicator nodes, respectively a `receiverStage` and a `senderStage`. The middle node is a simple `ff_node_t` which simply lets tasks flow between the two nodes communicating with groups G1 and G3;
- group G3: implemented in **Listing 7.1.4**, builds a `receiverStage` and a standard `ff_node_t` node which simply prints the tasks received.

```

1 firstStage first(stream_len);
2 senderStage sender(receiver_addr);
3 ff_Pipe<float> pipe(first, sender);
4 if (pipe.run_and_wait_end()<0) {
5     error("running pipe");
6     return -1;
7 }

```

Listing 7.1.2: G1 node composition

```

1 // Build addresses vector
2 ...
3
4 receiverStage receiver(addresses);
5 forwardStage first;
6 senderStage sender(receiver_addr);
7 ff_Pipe<float> pipe(receiver, first, sender);
8 if (pipe.run_and_wait_end()<0) {
9     error("running pipe");
10    return -1;
11 }

```

Listing 7.1.3: G2 node composition

```

1 // Build the addresses vector
2 ...
3
4 receiverStage receiver(addresses);
5 forwardStage first;
6 ff_Pipe<float> pipe(receiver, first);
7 if (pipe.run_and_wait_end()<0) {
8     error("running pipe");
9     return -1;
10 }

```

Listing 7.1.4: G3 node composition

With the implemented classes, building various groups communicating between each other is very simple and straightforward. Moreover, the way communications are abstracted by the underlying frameworks, makes it easy and painless to extend functionalities of the communication nodes. Adding new protocols requires zero effort and no code modifications from the user point-of-view, since the only step to



follow is to define a Mercury-accepted string with the preferred protocol to use during communication and use it at initialization of the receiver node.

Note that, since the receiver stages can be initialized with a vector of strings which represents all the endpoints to which the receiver will listen on, we can (and we must, to allow correct termination) compose multiple instances of groups (G1) and (G2), based on the amount of endpoints created, respectively by (G2) and (G3). For example, if we create a group (G2) with two listening endpoints, in order to allow the whole application to terminate correctly, we must run two (G1) connecting, in turn, to both the endpoints.

# Bibliography

- [1] Daniel A. Reed and Jack Dongarra. “Exascale Computing and Big Data”. In: *Commun. ACM* 58.7 (June 2015), pp. 56–68. ISSN: 0001-0782. DOI: [10.1145/2699414](https://doi.org/10.1145/2699414). URL: <https://doi.org/10.1145/2699414>.
- [2] Cristian Ramon-Cortes et al. “A survey on the Distributed Computing stack”. In: *Computer Science Review* 42 (2021), p. 100422. ISSN: 1574-0137. DOI: <https://doi.org/10.1016/j.cosrev.2021.100422>. URL: <https://www.sciencedirect.com/science/article/pii/S1574013721000629>.
- [3] Stephen Kaisler et al. “Big Data: Issues and Challenges Moving Forward”. In: *2013 46th Hawaii International Conference on System Sciences*. 2013, pp. 995–1004. DOI: [10.1109/HICSS.2013.645](https://doi.org/10.1109/HICSS.2013.645).
- [4] Robert B. Ross et al. “Mochi: Composing Data Services for High-Performance Computing Environments”. In: *J. Comput. Sci. Technol.* 35.1 (Jan. 2020), pp. 121–144. ISSN: 1000-9000. DOI: [10.1007/s11390-020-9802-0](https://doi.org/10.1007/s11390-020-9802-0). URL: <https://doi.org/10.1007/s11390-020-9802-0>.
- [5] Jerome Soumagne et al. “Mercury: Enabling remote procedure call for high-performance computing”. In: *2013 IEEE International Conference on Cluster Computing (CLUSTER)*. 2013, pp. 1–8. DOI: [10.1109/CLUSTER.2013.6702617](https://doi.org/10.1109/CLUSTER.2013.6702617).
- [6] Sangmin Seo et al. “Argobots: A Lightweight Low-Level Threading and Tasking Framework”. In: *IEEE Transactions on Parallel and Distributed Systems* 29.3 (2018), pp. 512–526. DOI: [10.1109/TPDS.2017.2766062](https://doi.org/10.1109/TPDS.2017.2766062).
- [7] Marco Aldinucci et al. “Fastflow: High-Level and Efficient Streaming on Multicore”. In: Mar. 2014. ISBN: 9780470936900. DOI: [10.1002/9781119332015.ch13](https://doi.org/10.1002/9781119332015.ch13).
- [8] Marco Aldinucci et al. “FastFlow: High-level and Efficient Streaming on Multi-core”. In: 2017.

- [9] Krste Asanovic et al. “A View of the Parallel Computing Landscape”. In: *Commun. ACM* 52.10 (Oct. 2009), pp. 56–67. ISSN: 0001-0782. DOI: [10.1145/1562764.1562783](https://doi.org/10.1145/1562764.1562783). URL: <https://doi.org/10.1145/1562764.1562783>.
- [10] Oracle, *Java remote method invocation specification*. 2017. URL: <https://docs.oracle.com/javase/9/docs/specs/rmi/>.
- [11] Fabian Breg et al. “Java RMI Performance and Object Model Interoperability: Experiments with Java/HPC++ Distributed Components”. In: *zz* (Dec. 2012).
- [12] Message Passing Interface Forum. *MPI: A Message-Passing Interface Standard Version 4.0*. June 2021. URL: <https://www.mpi-forum.org/docs/mpi-4.0/mpi40-report.pdf>.
- [13] Edgar Gabriel et al. “Open MPI: Goals, Concept, and Design of a Next Generation MPI Implementation”. In: *Recent Advances in Parallel Virtual Machine and Message Passing Interface*. Ed. by Dieter Kranzlmüller, Péter Kacsuk, and Jack Dongarra. Berlin, Heidelberg: Springer Berlin Heidelberg, 2004, pp. 97–104. ISBN: 978-3-540-30218-6.
- [14] *MPICH Overview — MPICH*. en-US. URL: <https://www.mpich.org/about/overview/> (visited on 05/24/2022).
- [15] Pieter Hintjens. *ZeroMQ*. O’Reilly Media, Mar. 2013. ISBN: 9781449334062. URL: <https://www.oreilly.com/library/view/zeromq/9781449334437/>.
- [16] Pavel Shamis et al. “UCX: An Open Source Framework for HPC Network APIs and Beyond”. In: *2015 IEEE 23rd Annual Symposium on High-Performance Interconnects*. 2015, pp. 40–43. DOI: [10.1109/HOTI.2015.13](https://doi.org/10.1109/HOTI.2015.13).
- [17] *Unified Communication X*. en. URL: <https://openucx.org/introduction/> (visited on 05/26/2022).
- [18] Paul Grun et al. “A Brief Introduction to the OpenFabrics Interfaces - A New Network API for Maximizing High Performance Application Efficiency”. In: *2015 IEEE 23rd Annual Symposium on High-Performance Interconnects*. 2015, pp. 34–39. DOI: [10.1109/HOTI.2015.19](https://doi.org/10.1109/HOTI.2015.19).
- [19] *Libfabric*. URL: <https://ofiwg.github.io/libfabric/> (visited on 05/26/2022).
- [20] Paul Grun et al. “A Brief Introduction to the OpenFabrics Interfaces - A New Network API for Maximizing High Performance Application Efficiency”. In: *2015 IEEE 23rd Annual Symposium on High-Performance Interconnects*. 2015, pp. 34–39. DOI: [10.1109/HOTI.2015.19](https://doi.org/10.1109/HOTI.2015.19).
- [21] *Mercury - Network Abstraction Layer. OFI plugin*. URL: <https://mercury-hpc.github.io/user/ofii/>.

- [22] *Mercury - Network Abstraction Layer. MPI plugin.* URL: <https://mercury-hpc.github.io/user/na/#deprecated-plugins>.
- [23] Pavel Shamis et al. “UCX: an open source framework for HPC network APIs and beyond”. In: *2015 IEEE 23rd Annual Symposium on High-Performance Interconnects*. IEEE. 2015, pp. 40–43.
- [24] Nawab Ali et al. “Scalable I/O forwarding framework for high-performance computing systems”. In: *2009 IEEE International Conference on Cluster Computing and Workshops*. 2009, pp. 1–10. DOI: [10.1109/CLUSTER.2009.5289188](https://doi.org/10.1109/CLUSTER.2009.5289188).
- [25] *Mercury - Network Abstraction Layer.* URL: <https://mercury-hpc.github.io/user/na/>.
- [26] *Mercury - Network Abstraction Layer. Shared Memory plugin.* URL: <https://mercury-hpc.github.io/user/sm/>.
- [27] *mercury-hpc.* URL: <https://github.com/mercury-hpc/mercury>.
- [28] *Mercury - Issue #489.* URL: <https://github.com/mercury-hpc/mercury/issues/489>.
- [29] *Mercury - Network Abstraction Layer. Available plugins.* URL: <https://mercury-hpc.github.io/user/na/#available-plugins>.
- [30] *OFI/Libfabric - tcp sockets fabric provider.* URL: [https://ofiwg.github.io/libfabric/v1.11.1/man/fi\\_tcp.7.html](https://ofiwg.github.io/libfabric/v1.11.1/man/fi_tcp.7.html).
- [31] *Mercury - Issue #418.* URL: <https://github.com/mercury-hpc/mercury/issues/418>.
- [32] *Mercury - Issue #332.* URL: <https://github.com/mercury-hpc/mercury/issues/332>.
- [33] *Open Fabric Interfaces.* URL: <https://github.com/ofiwg/libfabric>.
- [34] *Mercury - Network Abstraction Layer. UCX plugin.* URL: <https://mercury-hpc.github.io/user/na/#ucx>.
- [35] *Argobots: A lightweight low-level threading framework.* URL: <https://www.argobots.org/>.
- [36] *JSON-C - A JSON implementation in C.* URL: <https://github.com/json-c/json-c/wiki>.