

Tala - Distributed Computing

Jeff McGough

Department of Computer Science and Engineering
South Dakota School of Mines and Technology
Rapid City, SD 57701, USA

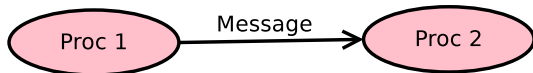
September 29, 2020

Tala

Distributed Communications

Problem

What if you want to send a message between two processes:



There are a lot of approaches to this problem. File system (files and pipes) are easy, but not efficient. For processes on different address spaces, using shared memory does not work. And neither files or SHM will work on when using completely separate systems.

Using sockets is the popular answer.

There are several issues using sockets

- ▶ You need to know the location of the process. IP address and Port must be known.
- ▶ There must be some synchronization, buffering or queuing between the processes.
- ▶ Your communication pattern needs to be client server where the sending process is the client and receiving process is the server. It is asymmetrical in nature.
- ▶ Socket libraries support low level communications and not the communication patterns we would like such as Publish/Subscribe.
- ▶ In all but the most simple applications, it requires a threaded design.

Communication Patterns

There are two patterns we will start with:

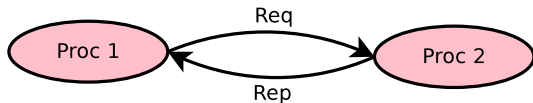
- ▶ Request - Reply : ReqRep
- ▶ Publish - Subscribe : PubSub

There are other patterns but for getting started on this project these two will suffice.

REQREP

Request-Reply is a more traditional client server communication pattern.

It is one to one direct communication.



A client will send a message to the server requesting some information. That is returned in the reply. It is a client initiated exchange.

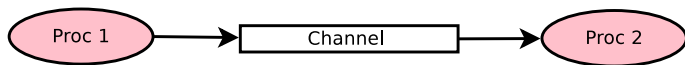
This will be useful for one-off types of communication.

PUBSUB

Publish-Subscribe is our main communication pattern for this project.

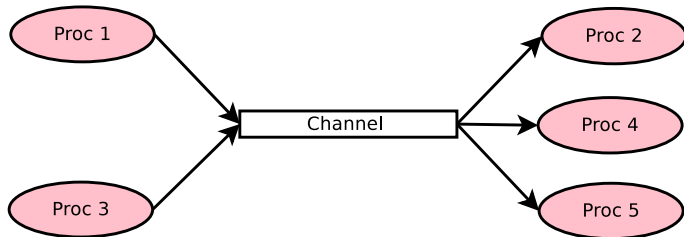
Formally, it is a buffered, indirect, many to many communication pattern.

In PubSub, a process will *publish* a message on a named channel or *topic*. This means that the message is placed in a message queue which is normally a FIFO.



A process that is subscribed can receive the message off of the other end of the topic (FIFO). Where this FIFO is actually located is to be determined.

Any number of processes can publish on the topic. These messages are queued up in the FIFO.



Any number of processes may subscribe to the topic.

More details

I should not need to know where my other processes are running.

I should only need to publish to a topic or subscribe to a topic;.

Topics can carry multiple messages and have a message type associated with them.

Tala

A fast distributed message system for remote interprocess communication.

The following graphical user story describes a common use case.

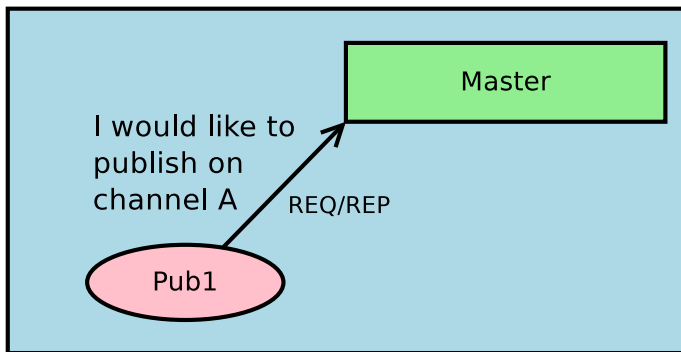
Although not all aspects of Tala are represented here, this use case is essential to understand before moving on.

In this example, we look at a design that has the FIFO or topic as a standalone or separate entity. The reasoning is that if the FIFO is bound to either publisher or subscriber, it might limit the number of publishers or subscribers (this will get us started for now).

This is a design decision. The following “User Stories” are to give a flavor of how the tool is used. How the directory service is constructed or the details on communication topics are not defined.

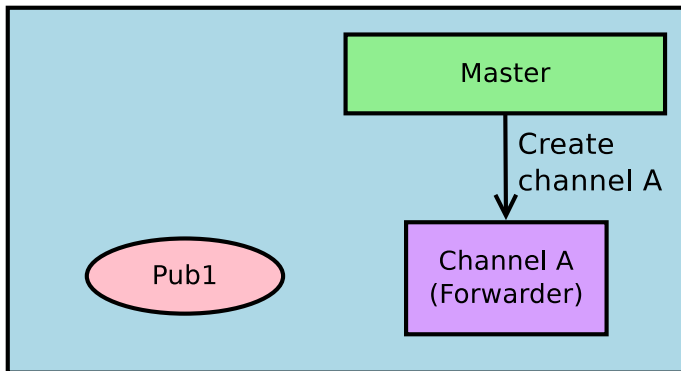
User Story

Request to connect to a specific topic. This can be built into a “topic open” type function call. Master process described later.



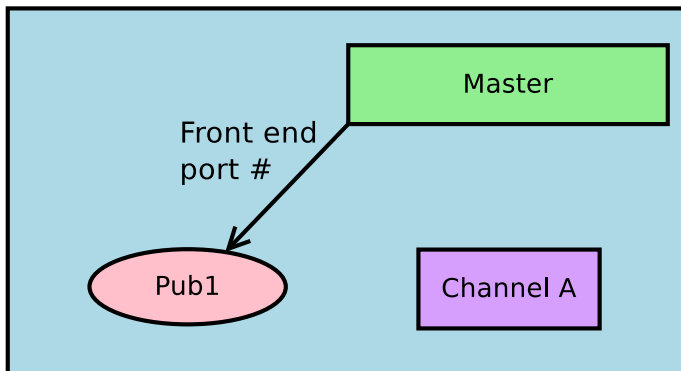
User Story

First time, the topic does not exist and needs to be created.



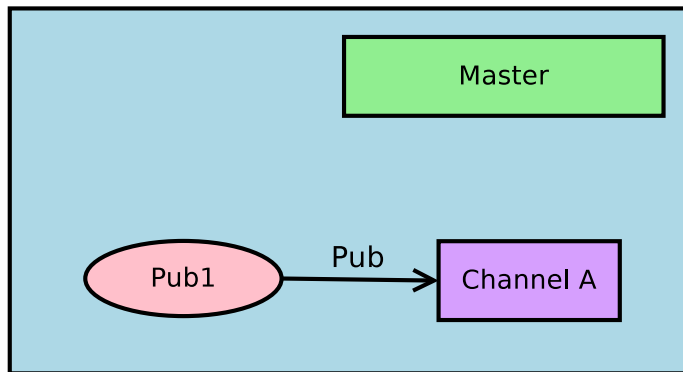
User Story

The relevant connection credentials are provided by the master to the prospective publisher (IP, port, etc for publisher side).



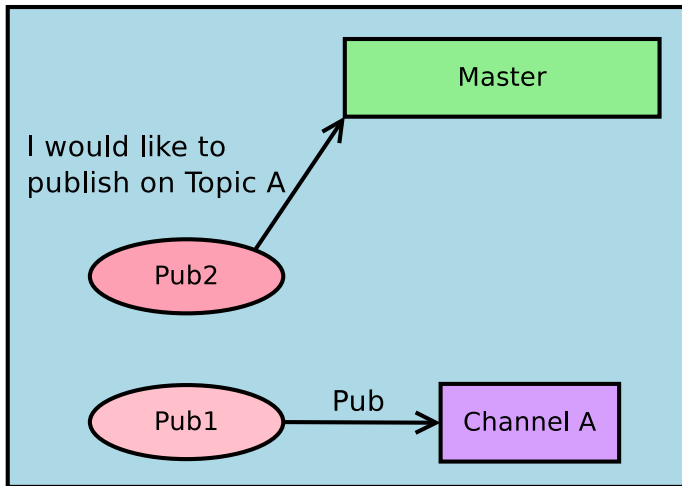
User Story

Using the credentials, the process (publisher) can connect with the topic “publish” messages.



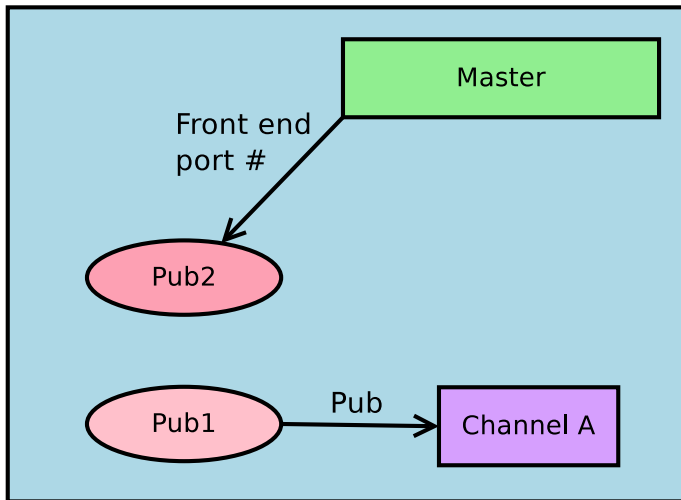
User Story

Another process may enter the process cluster and request to publish on topic A.



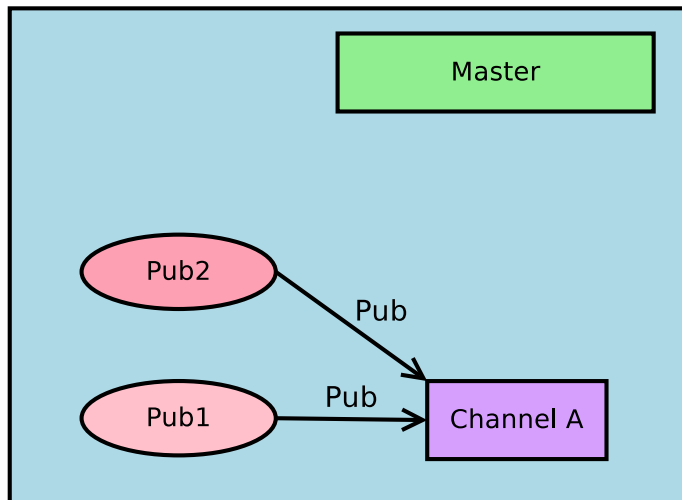
User Story

Topic A exists and so the credentials are sent from the Master's database.



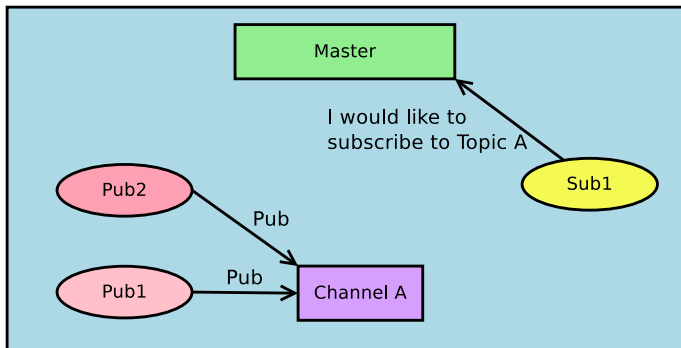
User Story

So both processes can publish on the topic.



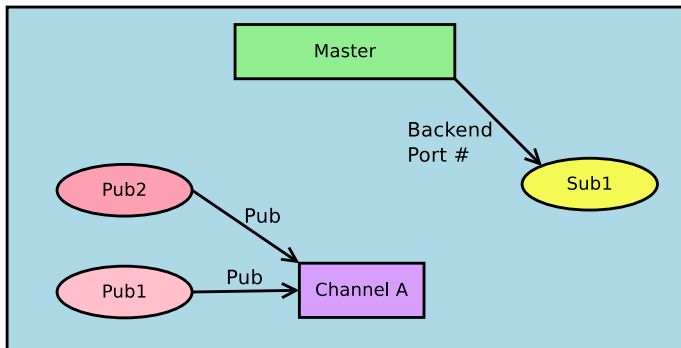
User Story

So far no process is receiving. Now, another process enters the cluster and requests to receive messages from topic A.



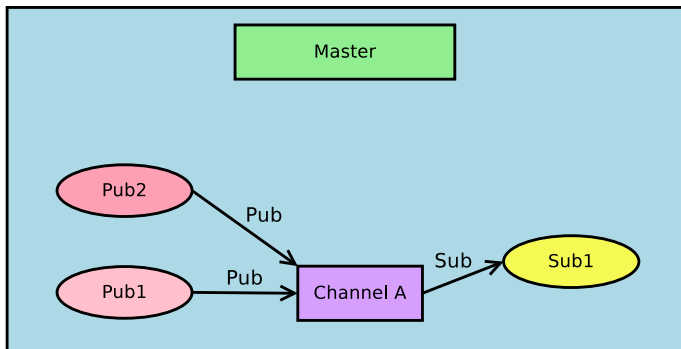
User Story

The master can send the credentials for the receiver side.



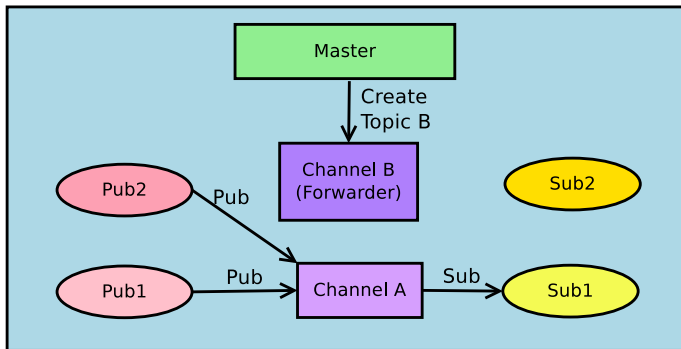
User Story

And this process can connect to the topic and receive the messages.



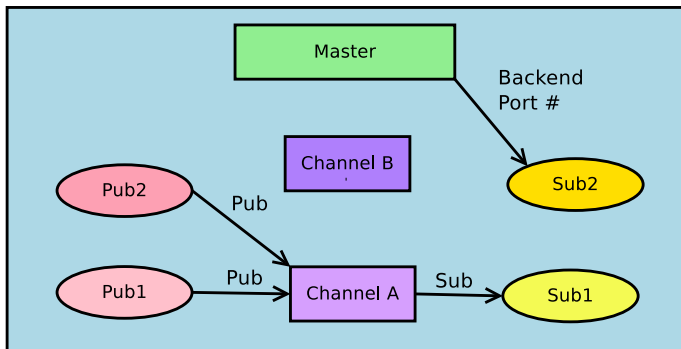
User Story

A fourth process can enter and request to receive messages on topic B. Since this is new, the master will create, store the data in its database.



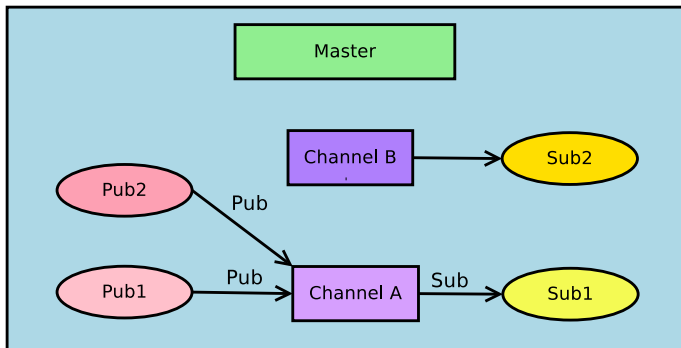
User Story

The master will then provide the relevant credentials to the new subscriber.



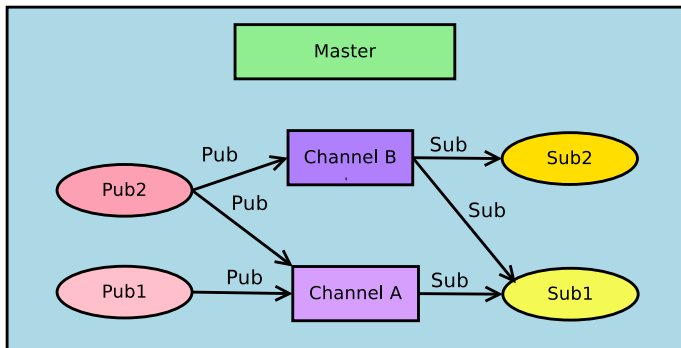
User Story

The second subscriber can connect and begin listening. The point here is that it should not matter the order that pubs and subs connect.



User Story

A process can publish to multiple topics and can subscribe to multiple topics.

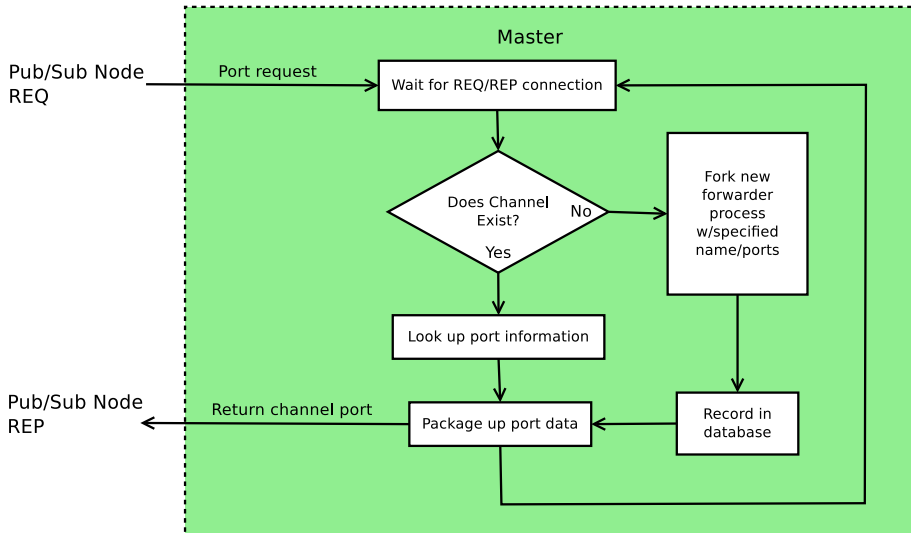


The master process is a directory service but not the router or message queue. Messages go directly from process to topic or topic to process.

This avoids bottlenecks in scaling.

Master keeps database of ports and IPs all pubs and sub. Keeps ports and IPs for processes managing the topics.

Master



Python Code Concept - Publisher

```
import tala as tl

master = "default"
name = "talkernode"
tl.join(name, master)

topic = "realityTV"

pub = tl.publisher(topic)

key = "0"
message = tl.pack("It is all scripted.")
tl.send(pub, message, key)
```

Python Code Concept - Subscriber

```
import tala as tl

master = "default"
name = "listennode"
tl.join(name, master)

topic = "realityTV"

sub = tl.subscribe(topic)

flag = 0 # timeout codes
message, key = tl.receive(sub, flag)
print(tl.unpack(message))
```

There are configuration aspects (in a config file and dynamic).

- ▶ Message queue size.
 - ▶ What is the basic message size limit?
 - ▶ How many messages?
- ▶ Message removal behavior and if multiple buffers required.
 - ▶ Does one subscriber get the message or all?
 - ▶ How do you manage this with multiple pubs and subs?
- ▶ Full FIFO behavior.
- ▶ Message data types.
 - ▶ Messages are tex/binary, strings/floats, arrays?
- ▶ Message data marshalling.
 - ▶ Pack up into a documented format?

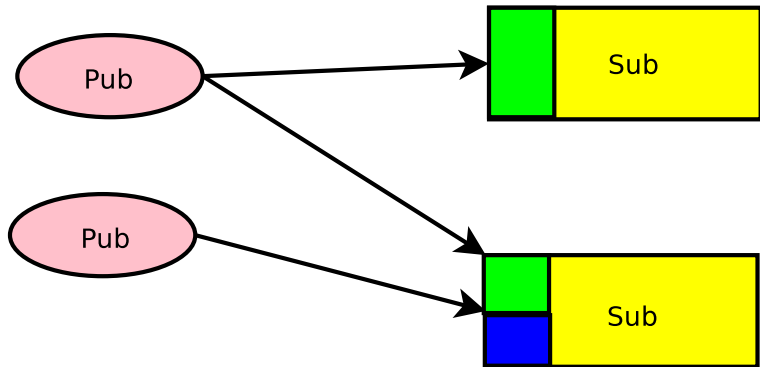
Alternate Design

The topic can be bound to the subscriber.

The master will hand the subscriber's IP/port to the publisher and point to point communication would happen.

For multiple subscribers, the master would need to contact each publisher and add a subscriber to their broadcast list.

Much of the user stories are the same, just the change in the location of the FIFO.



- ▶ Binding languages: Julia, Python, Rust
- ▶ Language for master (Rust or C++)
- ▶ Threading architecture
- ▶ Broker vs. Brokerless
- ▶ Communication lib (OMQ vs Sockets)
- ▶ Communication architecture
- ▶ FIFO architecture
- ▶ Performance (speed, latency, ...)
- ▶ Shared memory hook

- ▶ GUI management tool
- ▶ Node/topic Graph plotting tool
- ▶ Listener and Talker examples in various languages.

Open source aspects (host, license, etc)

Requirements

- ▶ Directory service
 - ▶ Separate process
 - ▶ Channel setup
 - ▶ Tracks nodes in the Tala cluster
 - ▶ Nodes and topics use unix filesystem naming (inherit channel properties)
 - ▶ Tracks channels in the Tala cluster
 - ▶ Answers requests for node communication data
 - ▶ Pub/Sub
 - ▶ Req/Rep
 - ▶ Send/Recv
 - ▶ Signal / Sync
 - ▶ Internal database for network graph
 - ▶ API for GUI management
 - ▶ Launch file - can launch a Tala process cluster

Requirements cont.

- ▶ GUI Management
 - ▶ Launch master
 - ▶ Runtime config
 - ▶ Port ranges
- ▶ Language APIs
 - ▶ Languages: Julia, Python, C/C++, Rust, Java
 - ▶ Node registration
 - ▶ Node Communication setup functions
 - ▶ Pub/sub/req/rep functions
 - ▶ Push/pop/send/receive functions
 - ▶ Signal functions
 - ▶ Node exit
 - ▶ Admin functions
 - ▶ Process sync functions

- ▶ Graph tool
 - ▶ Present the node-topic graph
 - ▶ Nodes are ellipses
 - ▶ Channels are rectangles
 - ▶ Connected by lines with arrows
 - ▶ Auto layout
 - ▶ Drag nodes and channels
 - ▶ Show nodes and select for node data
 - ▶ Show channels and select for channel data
 - ▶ Tap a channel and watch traffic
 - ▶ Place triggers on channels for callbacks

Software Expectations

- ▶ A prototype of the directory service - aka master.
- ▶ An API module for calling Tala in the target languages.
- ▶ A prototype of the GUI management tool.
- ▶ A prototype of the network graph tool.
- ▶ Some sample programs to illustrate Tala.

Project Expectations

- ▶ Design - because this is Senior Design.
- ▶ Communication - keep a dialog.¹
- ▶ Teamwork - a sense of service to others.
- ▶ A great attitude - because every day is a gift.
- ▶ Open source project - with software described above.
- ▶ Documentation - of the code and the experience.
- ▶ Performance analysis (comparison to sockets, MPI, Kafka)

¹I can handle bad news - I don't want any surprises