# Flexible data transport for online reconstruction

*M. Al-Turany*
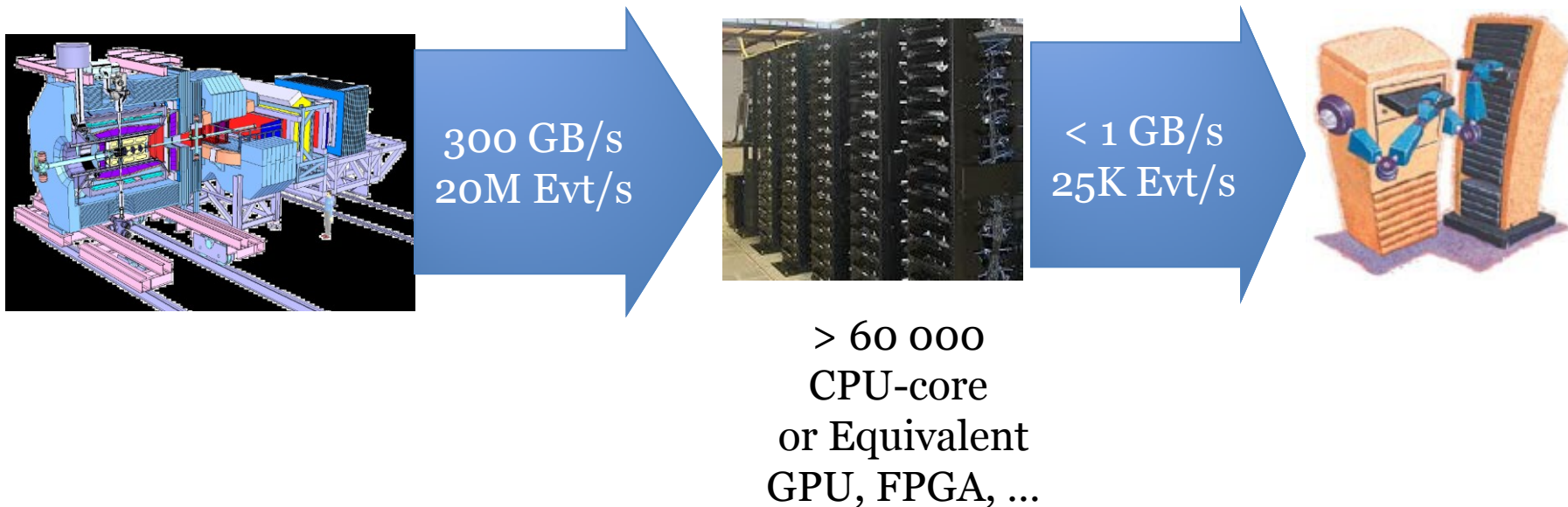
*Dennis Klein*

*A. Rybalchenko*

# This talk:

- Introduction
  - Design requirement
  - Zero MQ
  - Socket Pattern
- Current Status
- Results

# The Online Reconstruction and analysis



300 GB/s
20M Evt/s

> 60 000
CPU-core
or Equivalent
GPU, FPGA, ...

< 1 GB/s
25K Evt/s

How to manage the data flow on such a huge cluster?
How to recover single/multiple processes?
How to monitor it?
......

# Design constrains

- Highly flexible: different data paths should be modeled.

- Adaptive: Sub-system are continuously under development and improvement

- Should work for simulated and real data: developing and debugging the algorithms

- It should support all possible hardware where the algorithms could run (CPU, GPU, FPGA)

- It has to scale to any size! With minimum or ideally no effort.

M. Al-Turany, Panda Collaboration Meeting, Goa

# Before Re-inventing the Wheel

- What is available on the market and in the community?
  - ○ ALICE, ATLAS, CMS, LHCb, …
  - ○ Financial and weather application have also huge data to deal with
- Do we intend to separate online and offline?
- Multithreaded concept or a message queue based one?
  - ○ Message based systems allow us to decouple producers from consumers.
  - ○ We can spread the work to be done over several processes and machines.
  - ○ We can manage/upgrade/move around programs (processes) independently of each other.

# ØMQ (ZeroMQ)  Available since 2011

- A socket library that acts as a concurrency framework.

- Faster than TCP, for clustered products and supercomputing.

- Carries messages across inproc, IPC, TCP, and multicast.

- Connect N-to-N via fanout, pubsub, pipeline, request-reply.

- A synch I/O for scalable multicore message-passing apps.

- 30+ languages including C, C++, Java, .NET, Python.

- Most OSes including Linux, Windows, OS X, PPC405/PPC440.

- Large and active open source community.

- LGPL free software with full commercial support from iMatix.

# Zero in ØMQ

Originally the zero in ØMQ was meant as "zero broker" and (as close to) "zero latency" (as possible). In the meantime it has come to cover different goals:
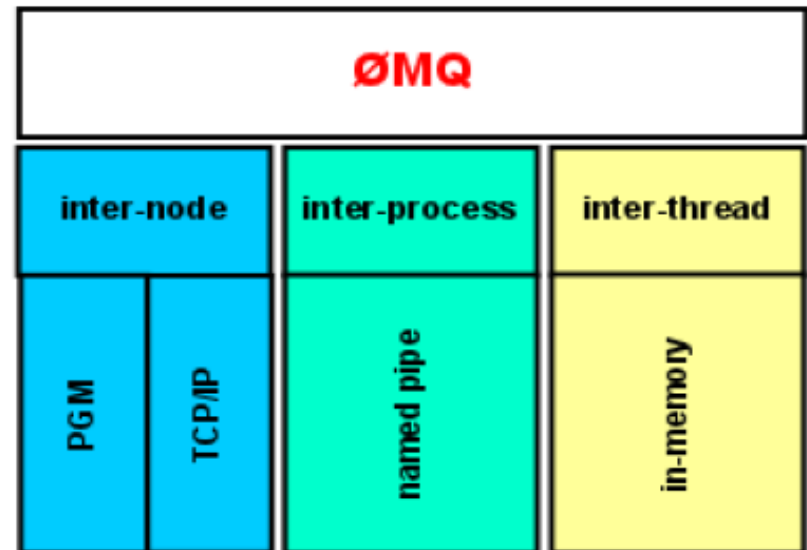
- zero administration,

- zero cost,

- zero waste.

More generally, "zero" refers to the <span style="color:red">culture of minimalism</span> that permeates the project.

<span style="color:red">Adding power by removing complexity rather than exposing new functionality.</span>

# ZeroMQ sockets provide efficient transport options

- Inter-thread

- Inter-process

- Inter-node
  - which is really just inter-process across nodes communication



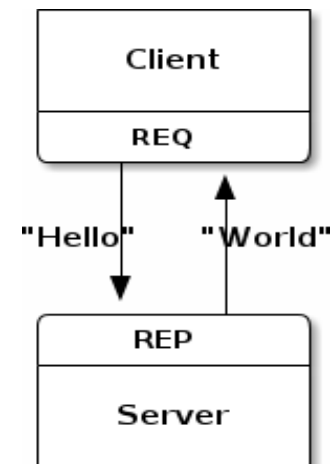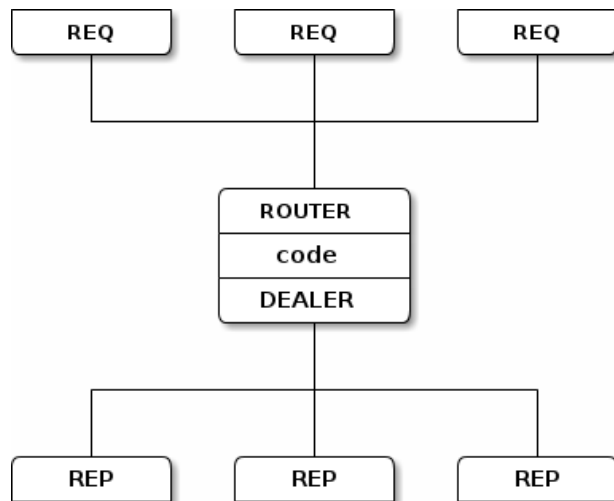PMG : Pragmatic General Multicast (a reliable multicast protocol)
Named Pipe:  Piece of random access memory (RAM) managed by the operating system and exposed to programs through a file descriptor and a named mount point in the file system.  It behaves as a first in first out (FIFO) buffer

# The built-in core ØMQ patterns are:

- **Request-reply**, which connects a set of clients to a set of services. This is a remote procedure call and task distribution pattern.

- **Publish-subscribe**, which connects a set of publishers to a set of subscribers. This is a data distribution pattern.

- **Pipeline**, which connects nodes in a fan-out / fan-in pattern that can have multiple steps, and loops. This is a parallel task distribution and collection pattern.

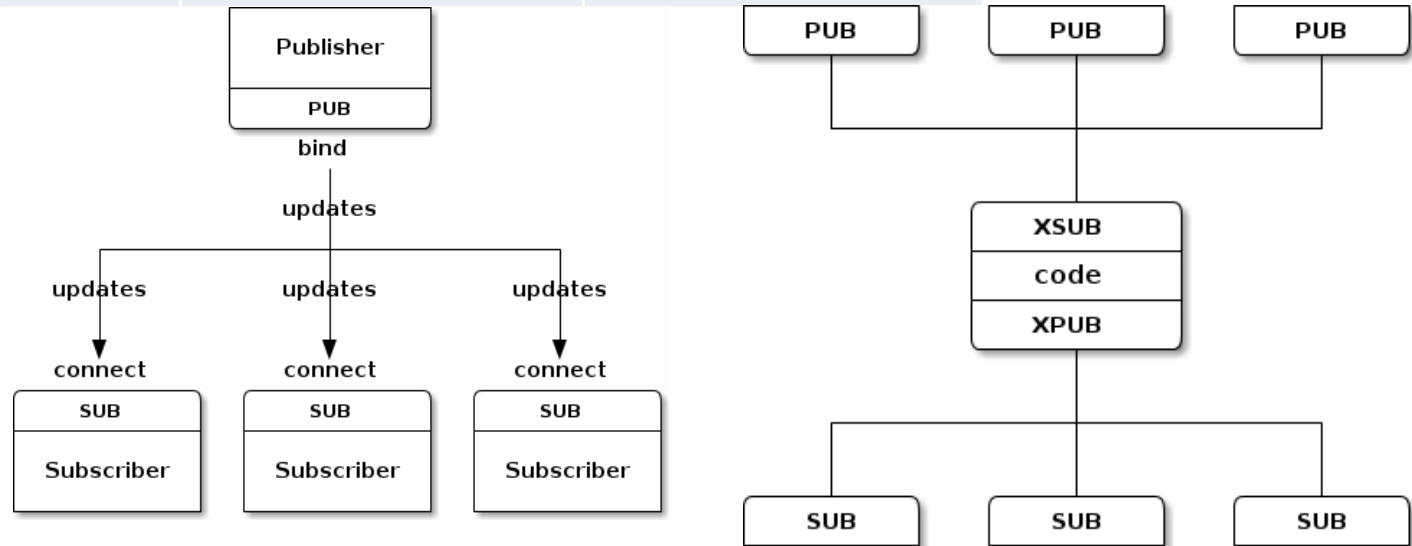- **Exclusive pair**, which connect two sockets exclusively

M. Al-Turany,  Panda Collaboration
Meeting, Goa

# Request-Reply Pattern

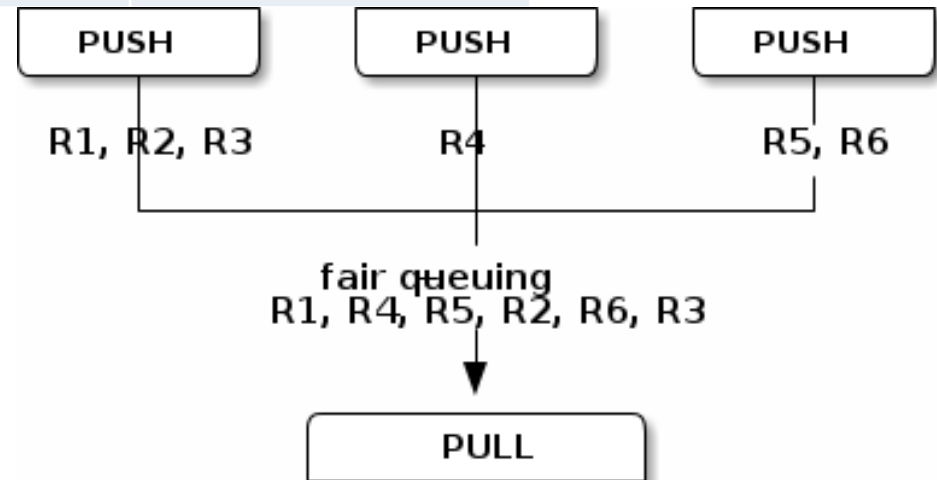| Socket type | **REQ** | **REP** |
|---|---|---|
| Compatible peer sockets | REP, ROUTER | REQ, DEALER |
| Direction | Bidirectional | Bidirectional |
| Send/receive pattern | Send, Receive | Send, Receive |
| Outgoing routing strategy | Round-robin | Last peer |
| Incoming routing strategy | Last peer | Fair-queued |
| Action in mute state | Block | Drop |

M. Al-Turany,  Panda Collaboration Meeting, Goa

# Publish-Subscribe Pattern

| Socket type | PUB | SUB |
| --- | --- | --- |
| Compatible peer sockets | SUB, XSUB | PUB, XPUB |
| Direction | Unidirectional | Unidirectional |
| Send/receive pattern | Send Only | Receive only |
| Outgoing routing strategy | Fan-out | N/A |
| Incoming routing strategy | N/A | Fair-queued |
| Action in mute state | Drop | Drop |

M. Al-Turany,  Panda Collaboration Meeting, Goa

# Pipeline Pattern

| Socket type | **PUSH** | **PULL** |
|---|---|---|
| Compatible peer sockets | PULL | PUSH |
| Direction | Unidirectional | Unidirectional |
| Send/receive pattern | Send Only | Receive only |
| Outgoing routing strategy | Round-Robin | N/A |
| Incoming routing strategy | N/A | Fair-queued |
| Action in mute state | Block | Block |

M. Al-Turany,  Panda Collaboration Meeting, Goa

## Example of sending control commands

- A worker process can manages two sockets (a PULL socket getting tasks, and a SUB socket getting control commands)

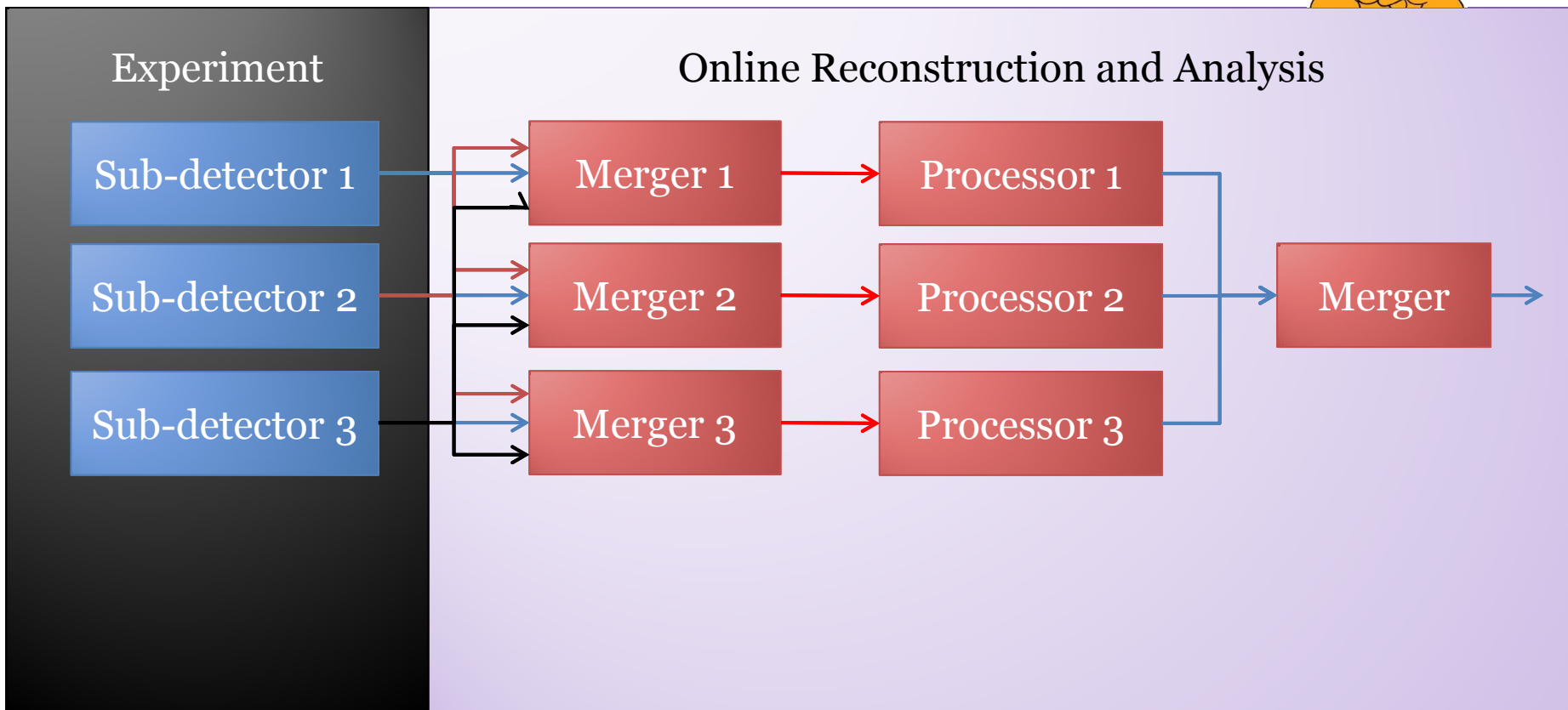- Could be very useful for calibration and alignment parameter

# Data Transfer Framework as Extension to FairRoot! Why?

- Modeling the pipeline processing within the online analysis

- Enable concurrency in FairRoot for offline analysis

- Reliable and efficient data transport through message queuing technology

- The long term plan is to have the same framework for online and offline

# Data flow example



**Experiment**

**Online Reconstruction and Analysis**

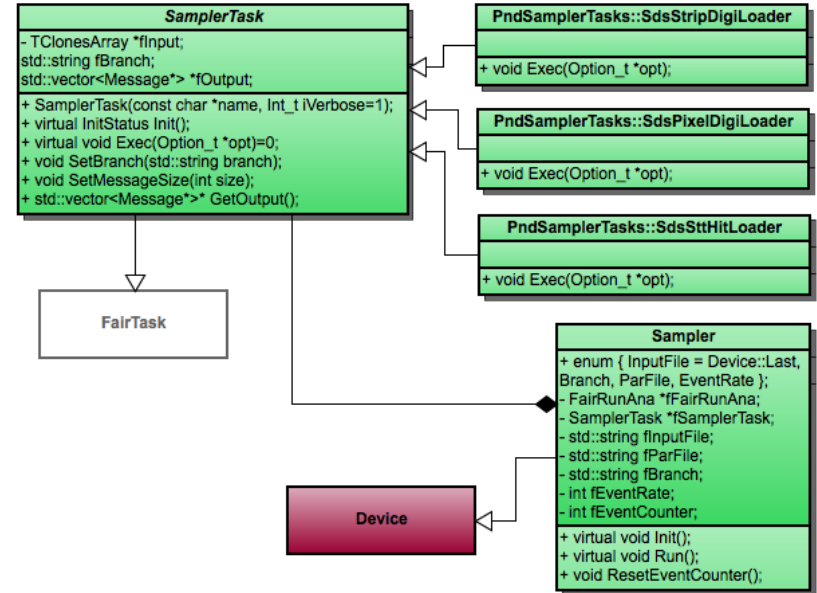| Sub-detector 1 | → | Merger 1 | → | Processor 1 |
| Sub-detector 2 | → | Merger 2 | → | Processor 2 | → | Merger | → |
| Sub-detector 3 | → | Merger 3 | → | Processor 3 |

# Current Status

- The Framework deliver some components which can be connected to each other in order to construct a processing pipeline.

- All component share a common base called Device (ZeroMQ Class).

- All devices are grouped by three categories:
  - Source: Sampler
  - Message-based Processor:
    - Sink, BalancedStandaloneSplitter, StandaloneMerger, Buffer
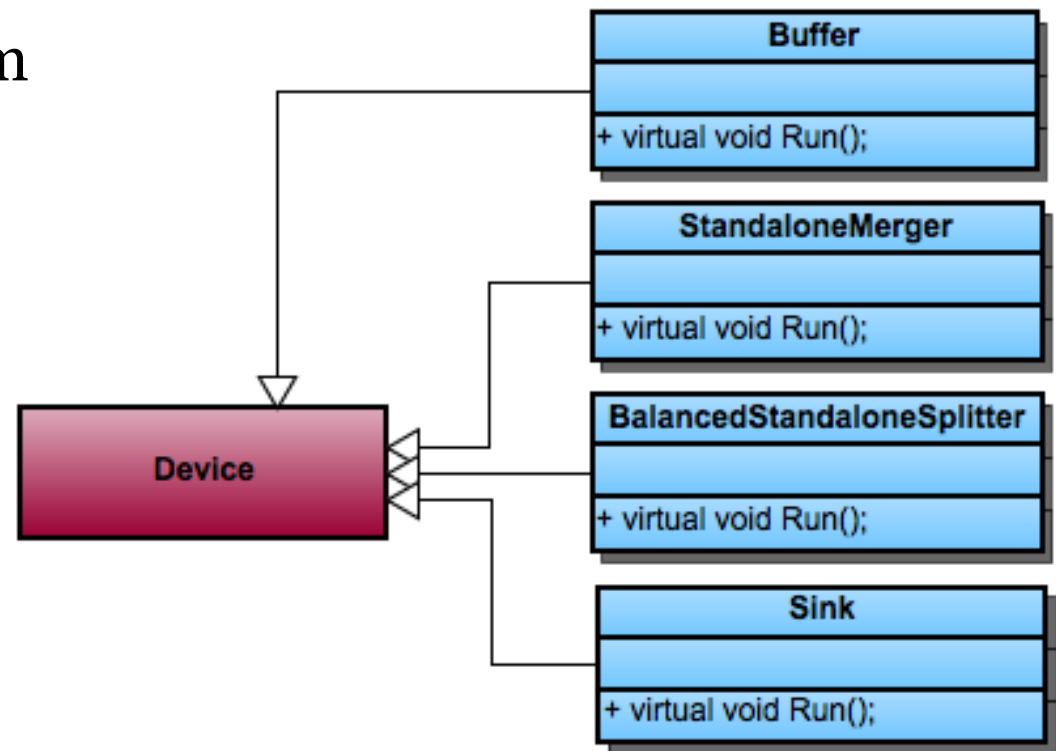  - Content-based Processor: Processor

# Sampler

- Devices with no inputs are categorized as sources

- During RUN state the sampler loops infinitely over the loaded events and send them through the output socket.

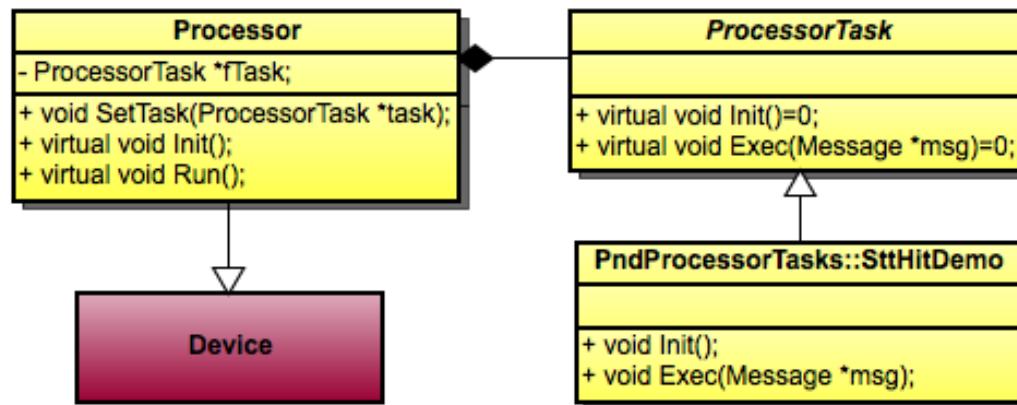- A variable event rate limiter has been implemented to control the sending speed

# Message-based Processor

- All message-based processors inherit from Device and operate on messages without interpreting their content.

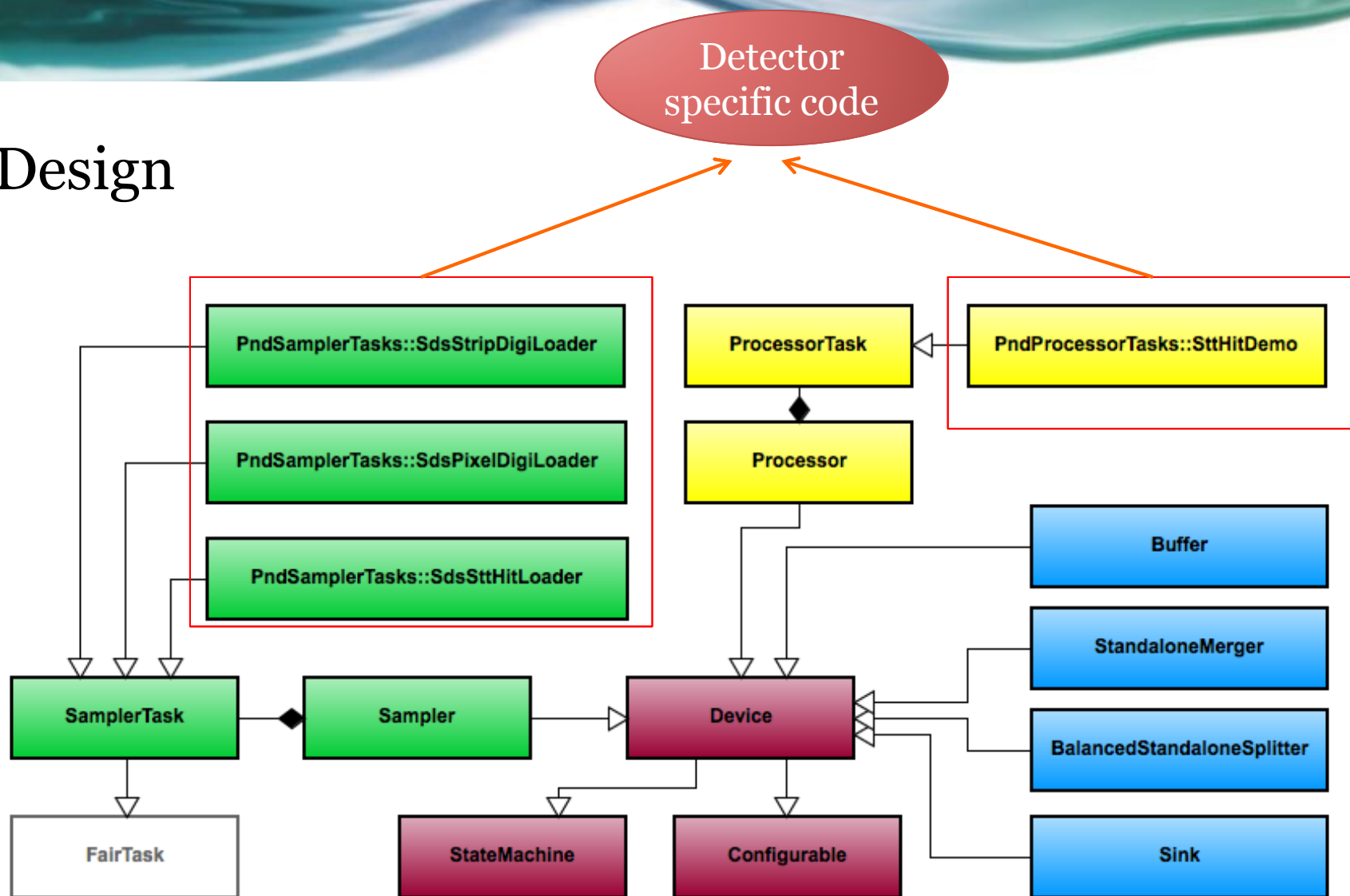- Four message-based processors have been implemented so far
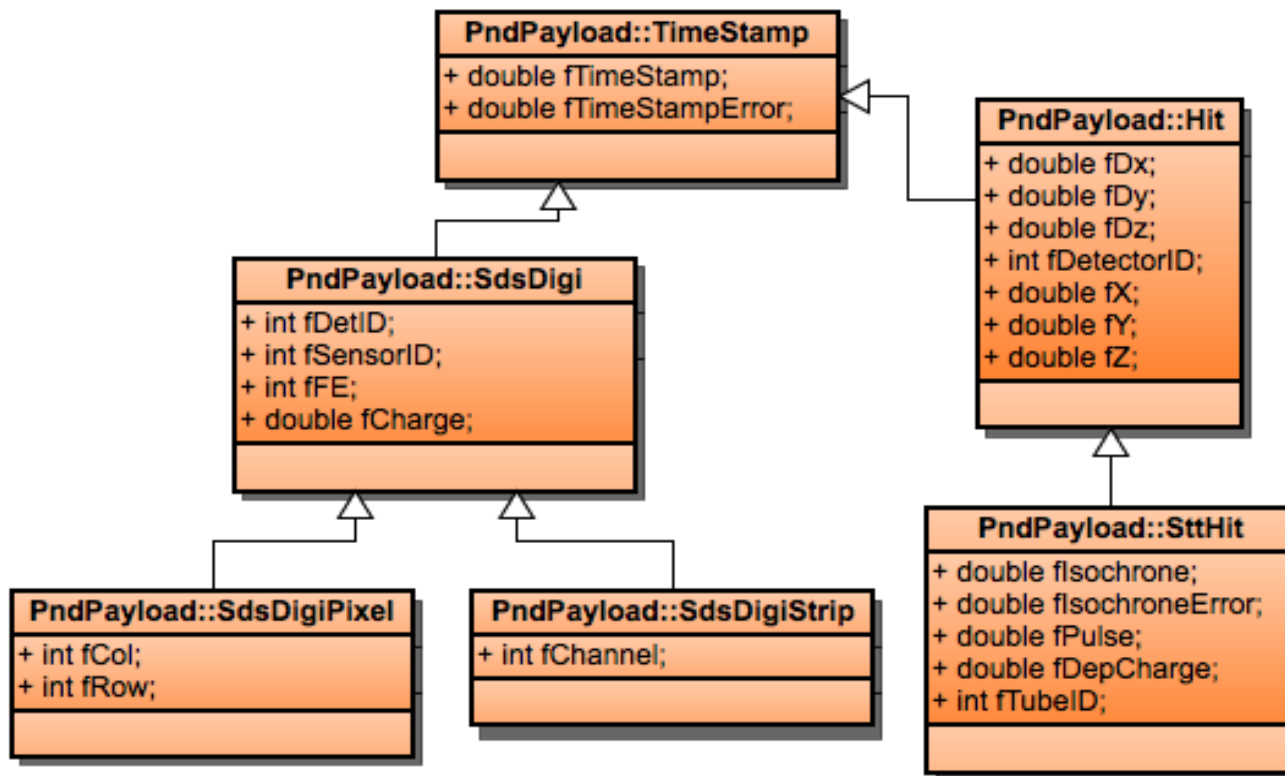
# Content-based Processor

- The Processor device has one input and one output socket.
- A task is meant for accessing and potentially changing the message content.

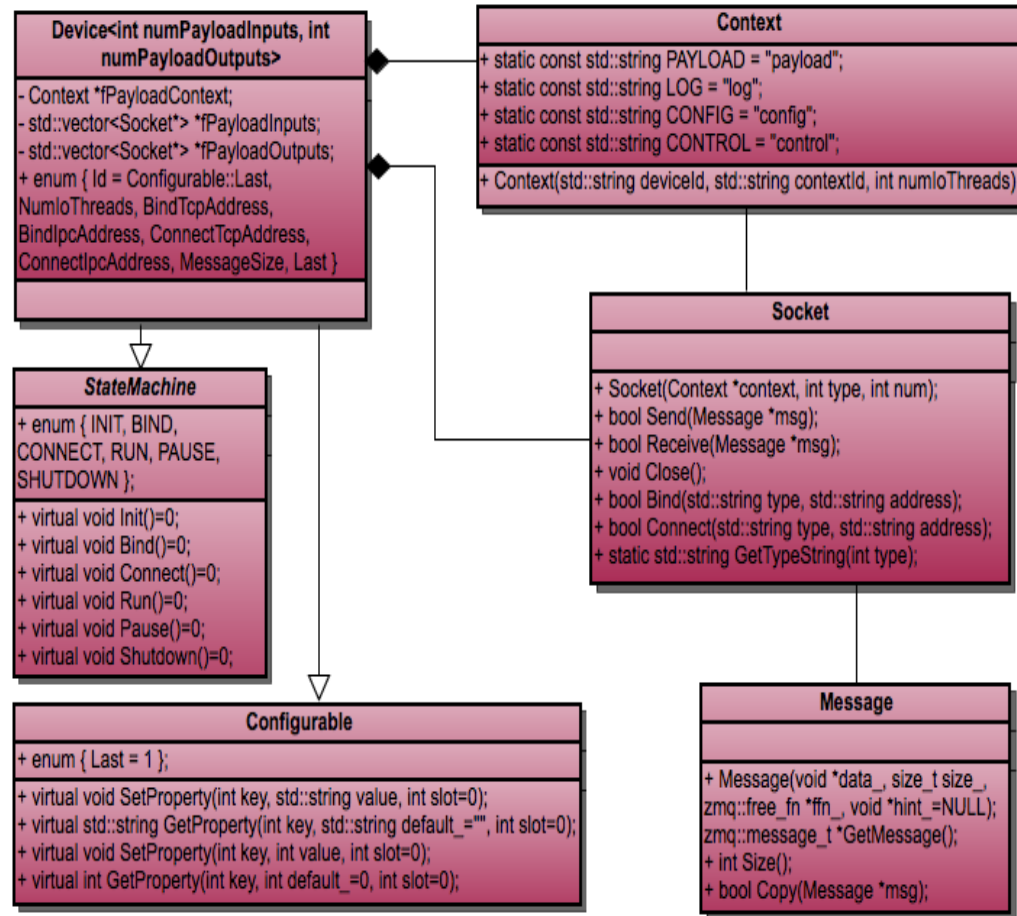M. Al-Turany, Panda Collaboration Meeting, Goa

# Design

M. Al-Turany, Panda Collaboration
Meeting, Goa

New simple classes without ROOT are used in the Sampler (This enable us to use non-ROOT clients) and reduce the messages size.

M. Al-Turany, Panda Collaboration Meeting, Goa

# Device

- Each processing stage of a pipeline is occupied by a process which executes an instance of the Device class

M. Al-Turany, Panda Collaboration Meeting, Goa

# Message format (Protocol)

- Potentially any content-based processor or any source can change the application protocol.

- The framework provides a  generic Message class that  works with any arbitrary and continuous junk of memory.

- One has to pass a pointer to the memory buffer and the size in bytes, and can optionally pass a function pointer to a destructor, which will be called once the message object is discarded.
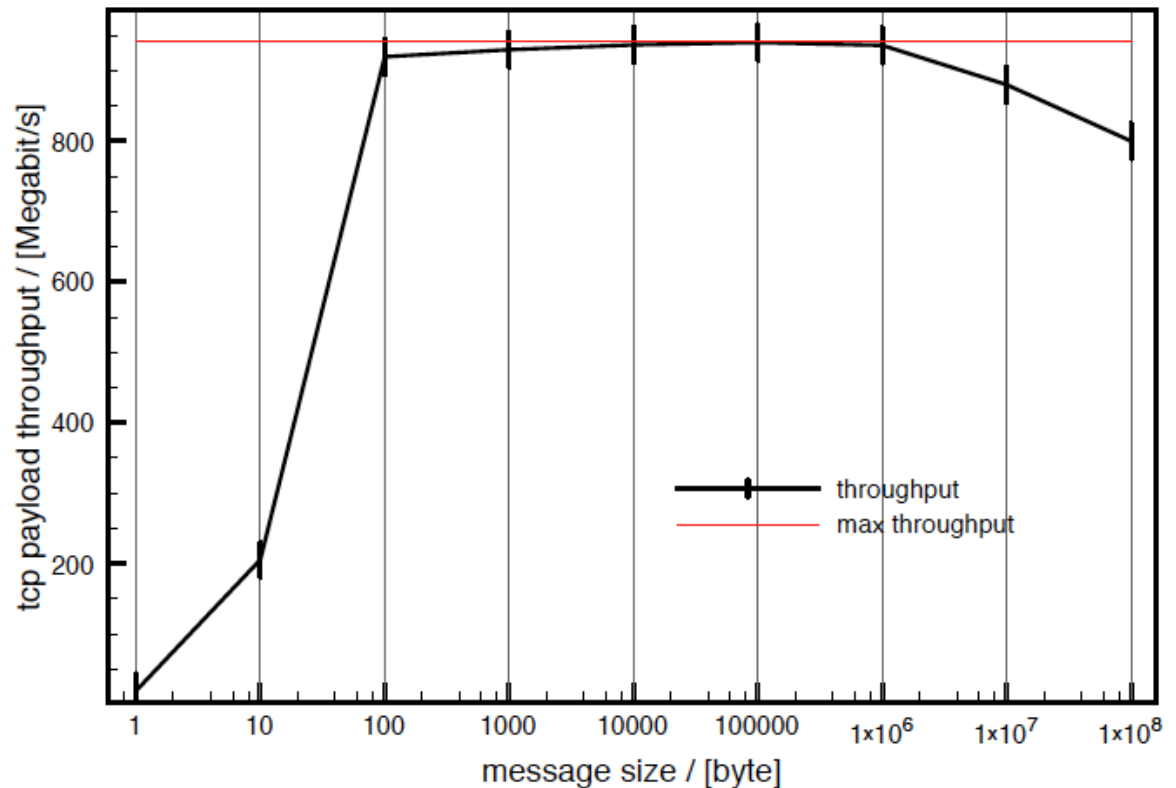
M. Al-Turany,  Panda Collaboration Meeting, Goa

# Test setup and results

M. Al-Turany,  Panda Collaboration Meeting, Goa

# Four identical nodes were connected to a GigabitEthernet switch for testing

- CPU:  Intel Xeon L5506 @ 2.13 GHz

- Memory:  24 GiB, 6  4 GiB

- Network:  Intel 82574L Gigabit Network Connection, speed=1Gbit/s

- Operating system GNU/Linux 3.2.32-1 x86_64, Debian 7.0

- ZeroMQ 3.2.0

- FairRoot PandaRoot oct12 release,

- Fairsoft development version from 18.12.2012

ZeroMQ reaches the upper TCP throughput for message sizes larger than a hundred bytes.



The last two measured values at 10 MB and 100 MB message size appear to be less efficient due to non-fractional output of the benchmark program. The near to maximum throughput for these last two values has been confirmed by monitoring the throughput with the linux tool iftop

# Results

- TCP throughput of <span style="color:red">117.6MB/s</span> was measured which is very close to the theoretical limit of <span style="color:red">117.7 MB/s</span> for the TCP/IPv4/GigabitEthernet stack.
- This was achieved using the Linux default values for the Ethernet MTU (1500 B) and TCP buffer size (85.3 KB).

- The throughput for the <span style="color:red">named pipe</span> transport between two devices <span style="color:red">on one node</span> has been measured around <span style="color:red">1.7 GB/s</span>

# Thanks!

M. Al-Turany, Panda Collaboration Meeting, Goa
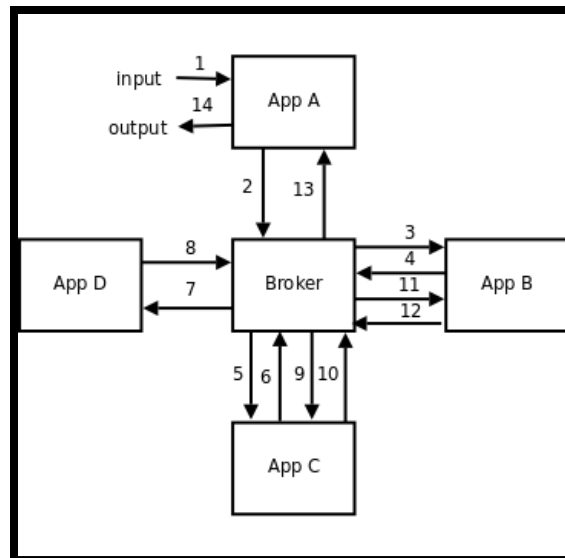
# Backup slides

M. Al-Turany,  Panda Collaboration Meeting, Goa

# Broker

- Architecture of most messaging systems is distinctive by the messaging server ("broker") in the middle.

- Every application is connected to the central broker.

- No application is speaking directly to the other application. All the communication is passed through the broker.
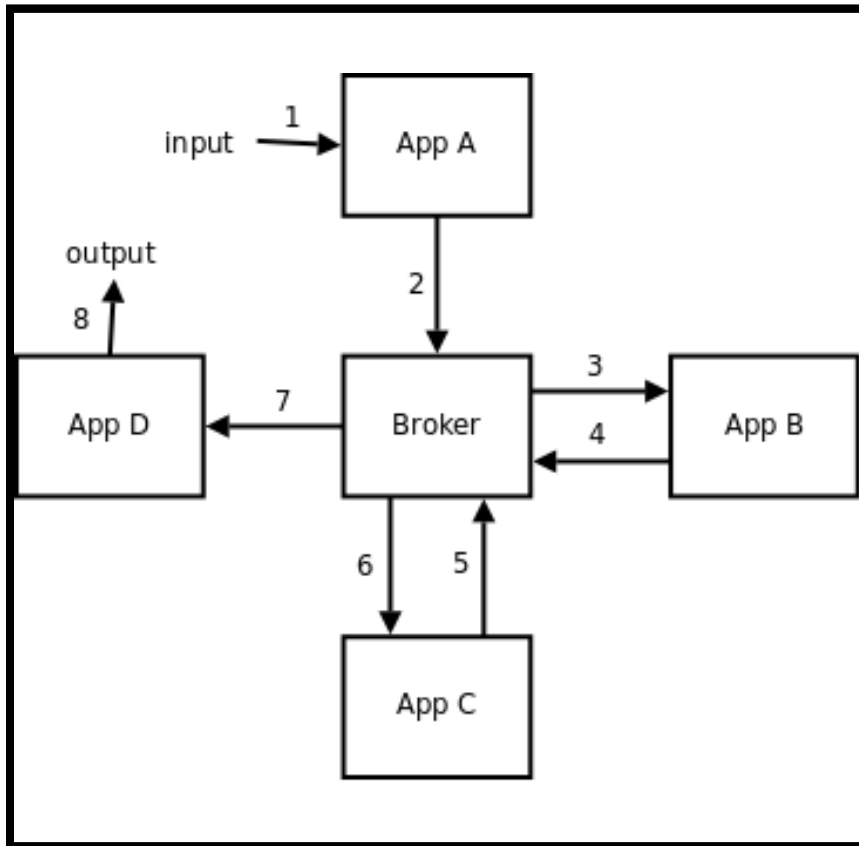
# Advantages

- Applications don't have to have any idea about location of other applications. The only address they need is the network address of the broker.

- Message sender and message receiver lifetimes don't have to overlap. Sender application can push messages to the broker and terminate. The messages will be available for the receiver application any time later.

- Broker model is to some extent resistant to the application failure. So, if the application is buggy and prone to failure, the messages that are already in the broker will be retained even if the application fails.
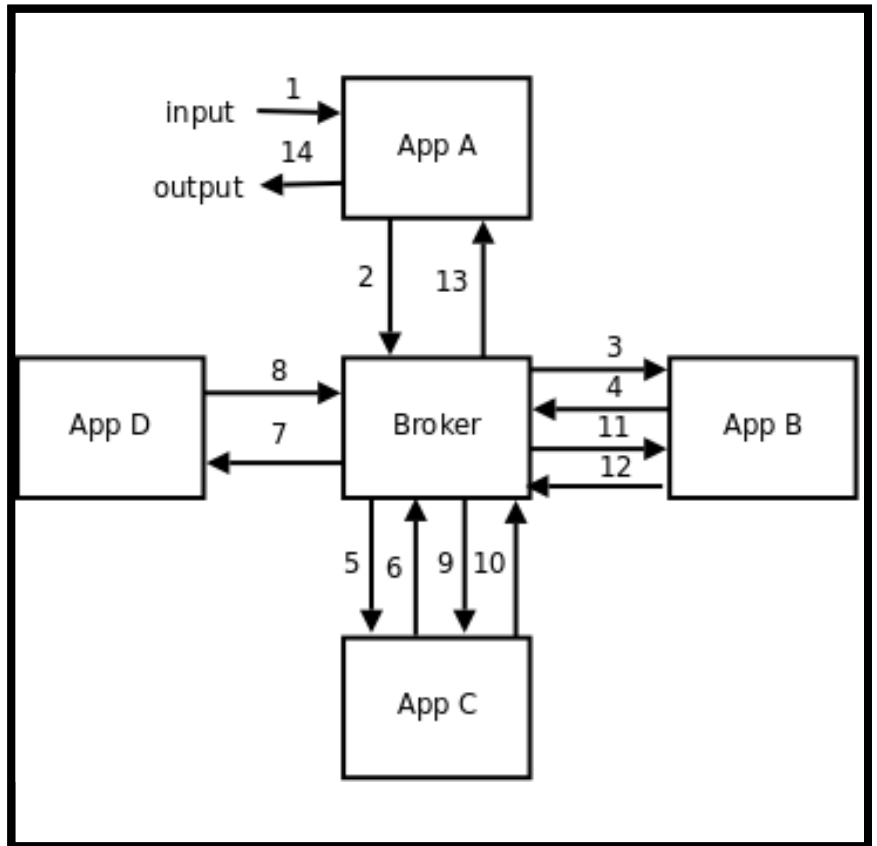
# Drawbacks

- It requires excessive amount of network communication.

- The fact that all the messages have to be passed through the broker can result in broker turning out to be the <span style="color:red">bottleneck</span> of the whole system.
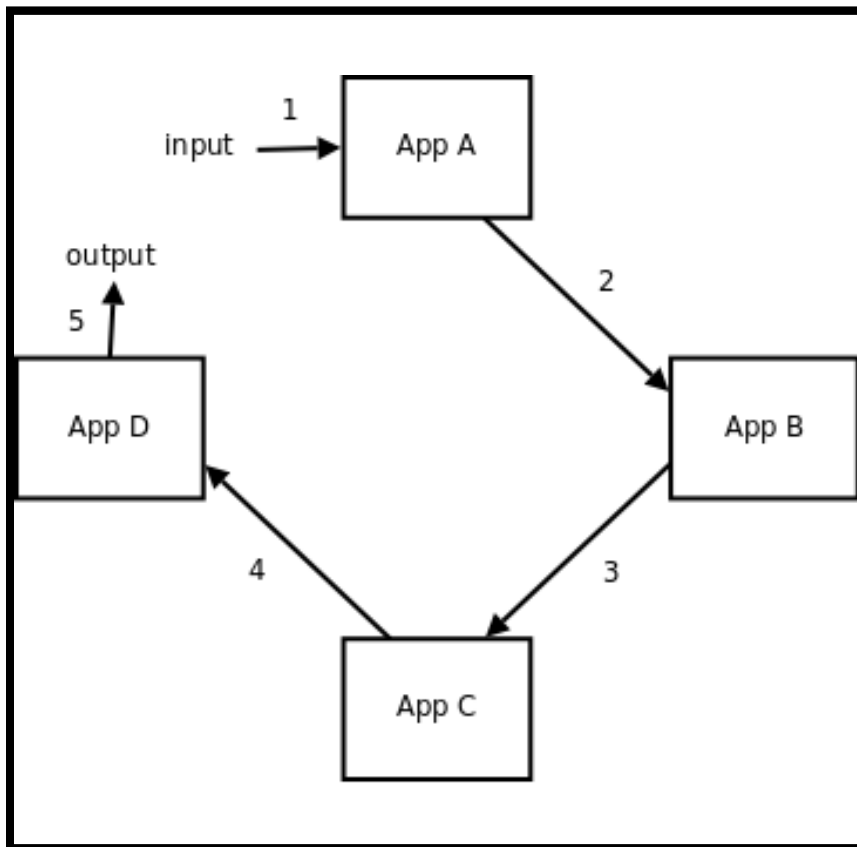
M. Al-Turany, Panda Collaboration
Meeting, Goa

# Broker pattern

pipelined fashion

service-oriented architectures (SOA)
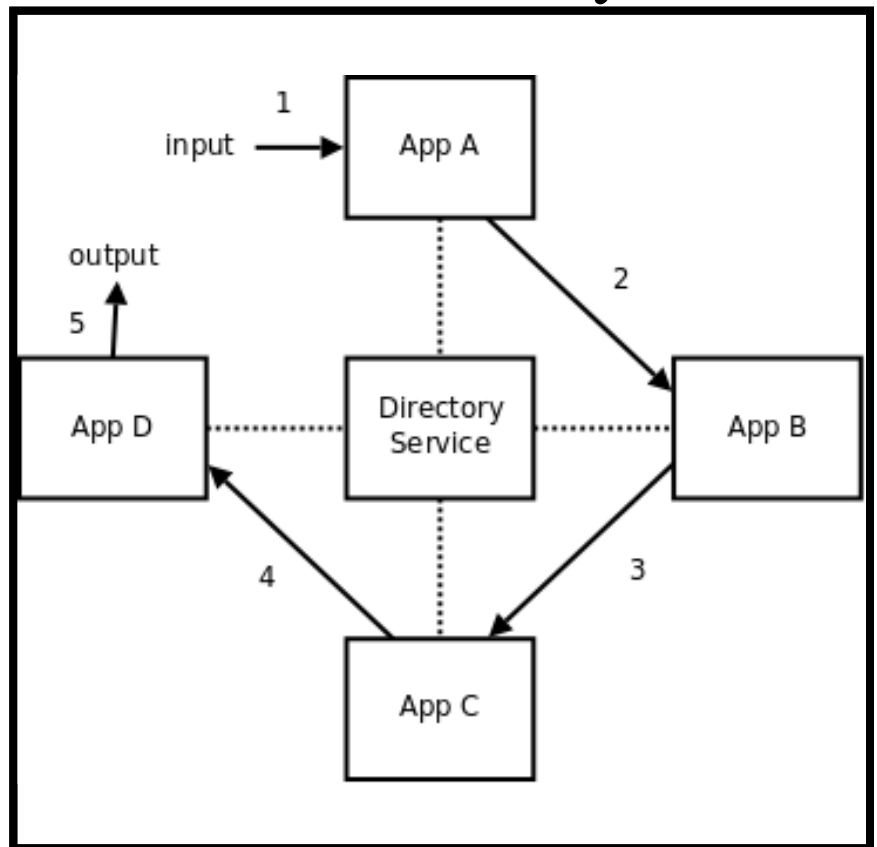
M. Al-Turany, Panda Collaboration Meeting, Goa

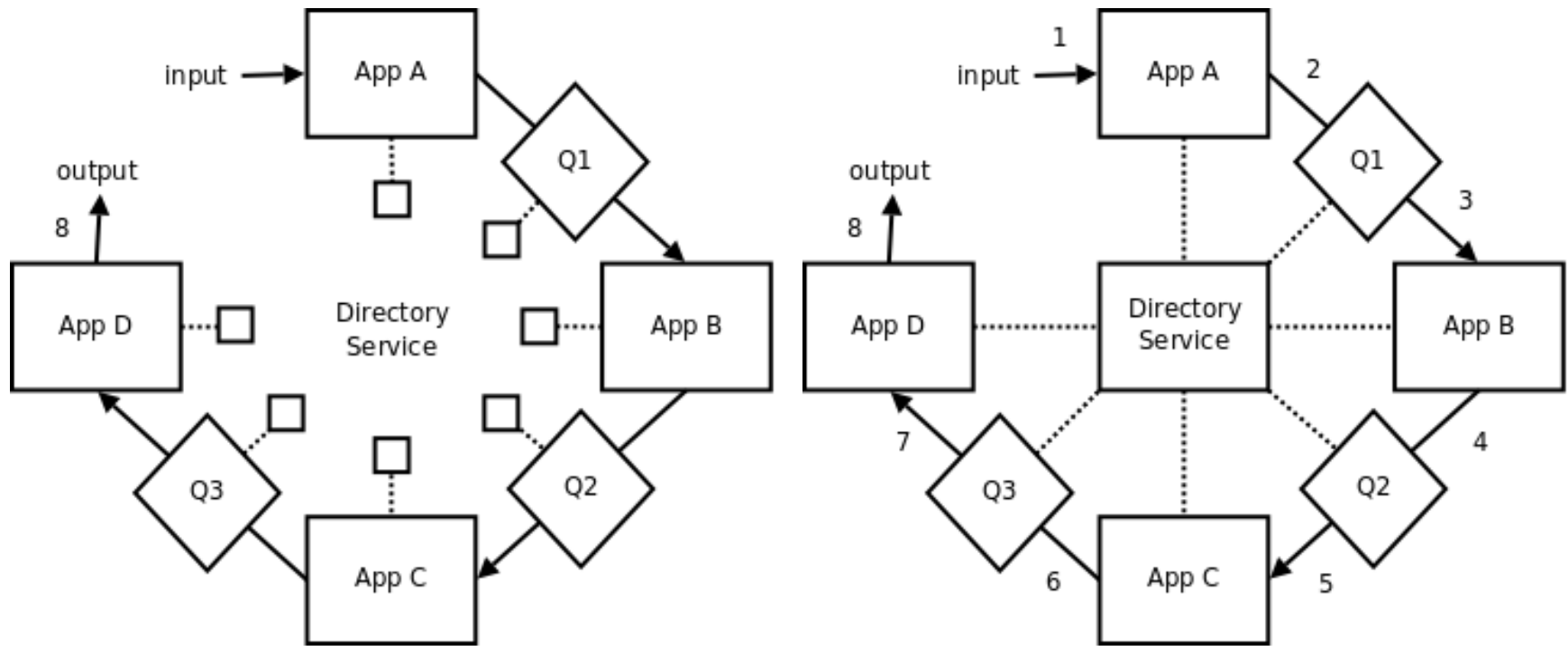# Examples of No Broker model in ZeroMQ

**No Broker**
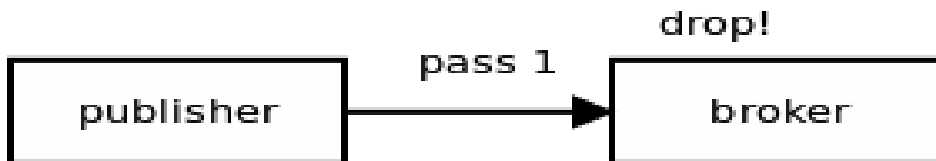


**Broker as a Directory Service**

# More Models

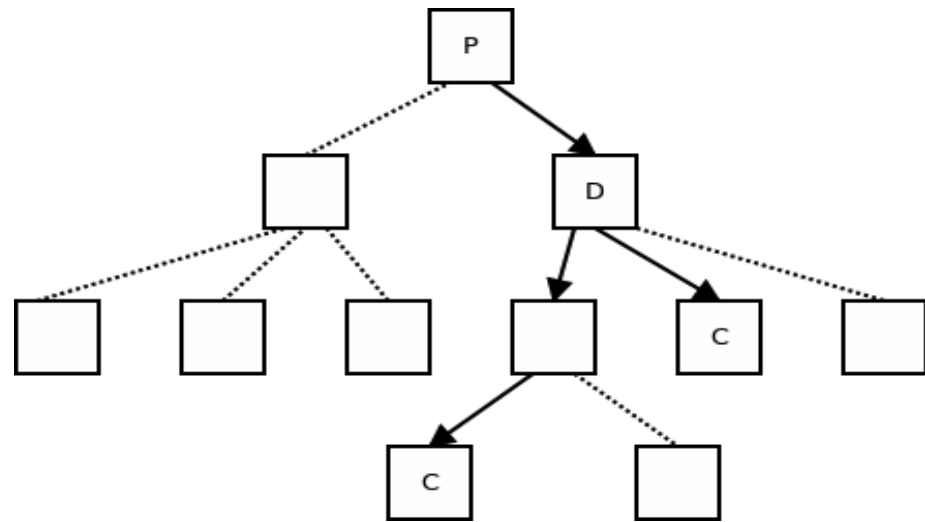**Distributed directory service**    **Distributed broker**

# PUB/SUB (publish/subscribe)



In most broker-based systems consumers subscribe for messages with the broker, however, there's no way for broker to subscribe for messages with the publisher. So, even if there is no consumer interested in the message it is still passed from the publisher to the broker, just to be dropped there

Messages should travel only through those lattices in the message distribution tree that lead to consumers interested in the message

M. Al-Turany, Panda Collaboration Meeting, Goa

# *Subscription Forwarding*

- XPUB is similar to a PUB socket, except that you can receive messages from it. The messages you receive are the subscriptions traveling upstream.
- XSUB is similar to SUB except that you subscribe by sending a subscription message to it.
- Subscription messages are composed of a single byte, 0 for un-subscription and 1 for subscription

M. Al-Turany,  Panda Collaboration Meeting, Goa