# Extending the FairRoot framework to allow for simulation and reconstruction of free streaming data

Mohammad Al-Turany
Dennis Klein
Anar Manafov
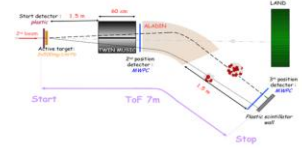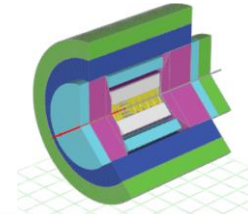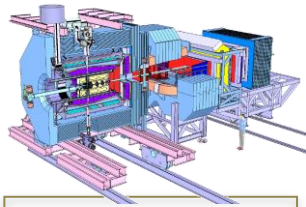Alexey Rybalchenko
Florian Uhlig
(GSI Darmstadt)

# What Does it Mean?

- Introduce pipelined data processing to the current Framework. (This talk!)
- Introduce time based simulation instead of event wise one. (already shown in CHEP 2012)
- Ideally: Keep compatibility to the current offline scheme.

# FairRoot



**Start testing the VMC concept for CBM**

**Panda decided to join->**
**FairRoot: same Base package for different experiments**

**R3B joined**

**EIC (Electron Ion Collider BNL)**
**EICRoot**

**SOFIA (Studies On Fission with Aladin)**

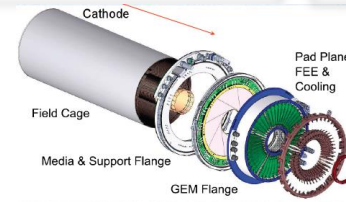**2004**   **2006**   **2010**   **2011**   **2012**   **2013**
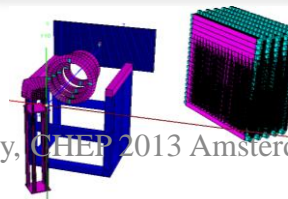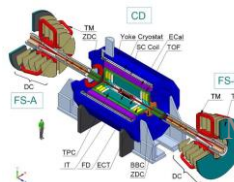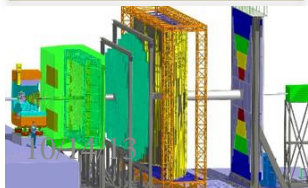
**First Release of CbmRoot**
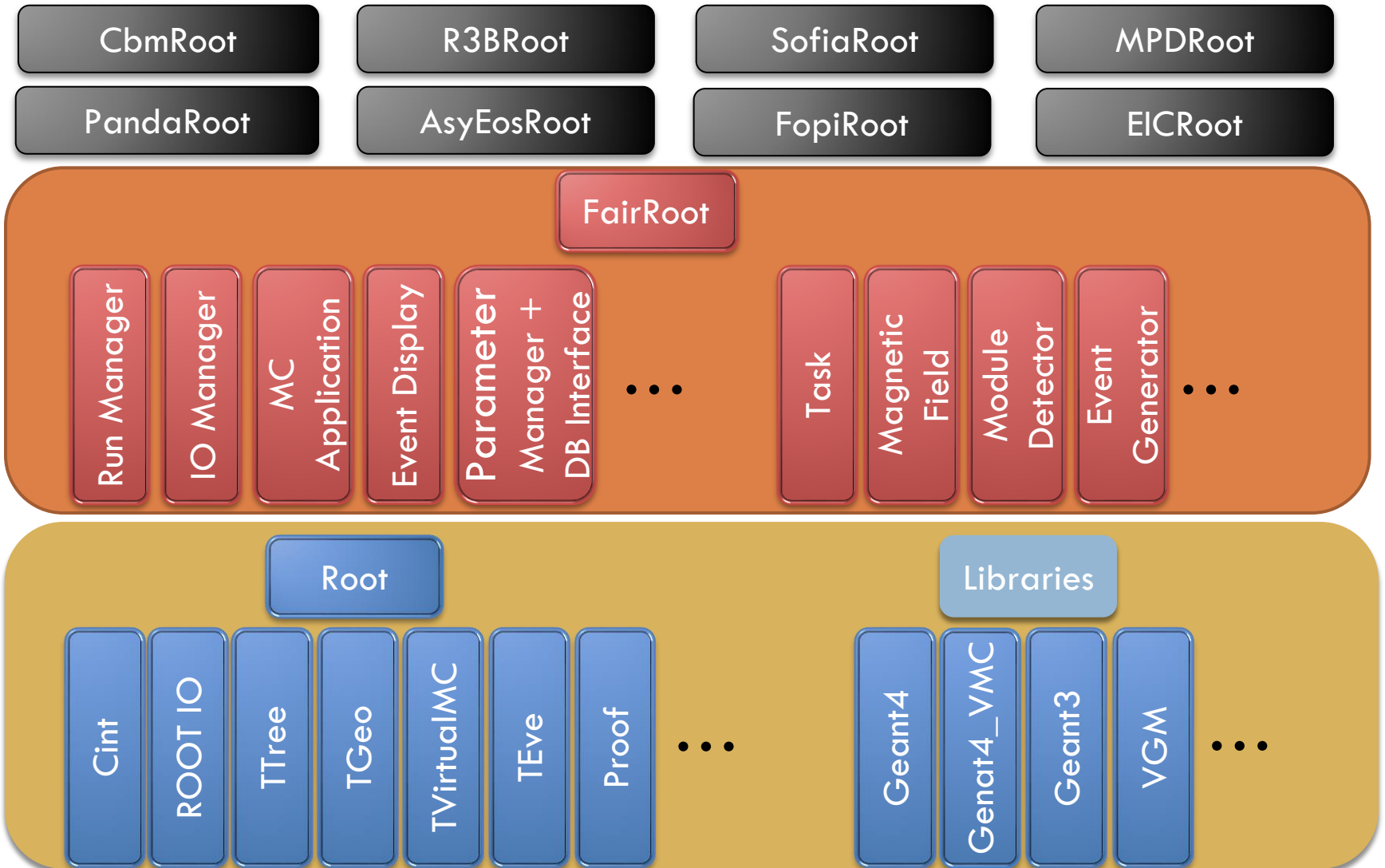
**MPD (NICA) start also using FairRoot**

**ASYEOS joined (ASYEOSRoot)**

**GEM-TPC seperated from PANDA branch (FOPIRoot)**

**ENSAR-ROOT**
**Collection of modules used by structural nuclear phsyics exp.**

# FairRoot: Implementation

CbmRoot   R3BRoot   SofiaRoot   MPDRoot

PandaRoot   AsyEosRoot   FopiRoot   EICRoot

FairRoot

Run Manager | IO Manager | MC Application | Event Display | Parameter Manager + DB Interface | ... | Task | Magnetic Field | Module Detector | Event Generator | ...

Root

Libraries

Cint | ROOT IO | TTree | TGeo | TVirtualMC | TEve | Proof | ... | Geant4 | Genat4_VMC | Geant3 | VGM | ...

# Next challenge is:  Online vs. Offline   or
# Online + Offline ?



**300 GB/s
20M Evt/s**

**> 60 000
CPU-core
or Equivalent
GPU, FPGA, …**
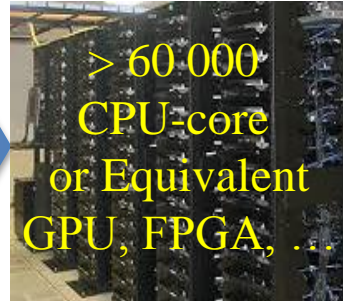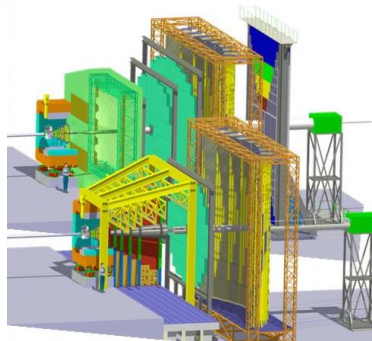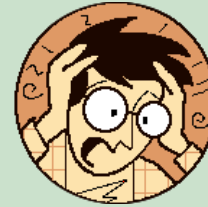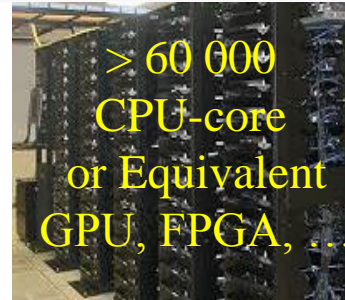
**< 1 GB/s
25K Evt/s**

How to distribute the processes?
How to manage the data flow?
How to recover processes when they crash?
How to monitor the whole system?
……

**1 TB/s**

**> 60 000
CPU-core
or Equivalent
GPU, FPGA, …**

**1 GB/s**

# Design constrains

- Highly flexible:
    - different data paths should be modeled.
- Adaptive:
    - Sub-systems 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.

# FairRoot: Where we are now?

- ROOT event loop

- User code in Task hierarchy

- Task hierarchy runs sequentially in one process

- Tasks implement only algorithms (can be exchanged/replaced)

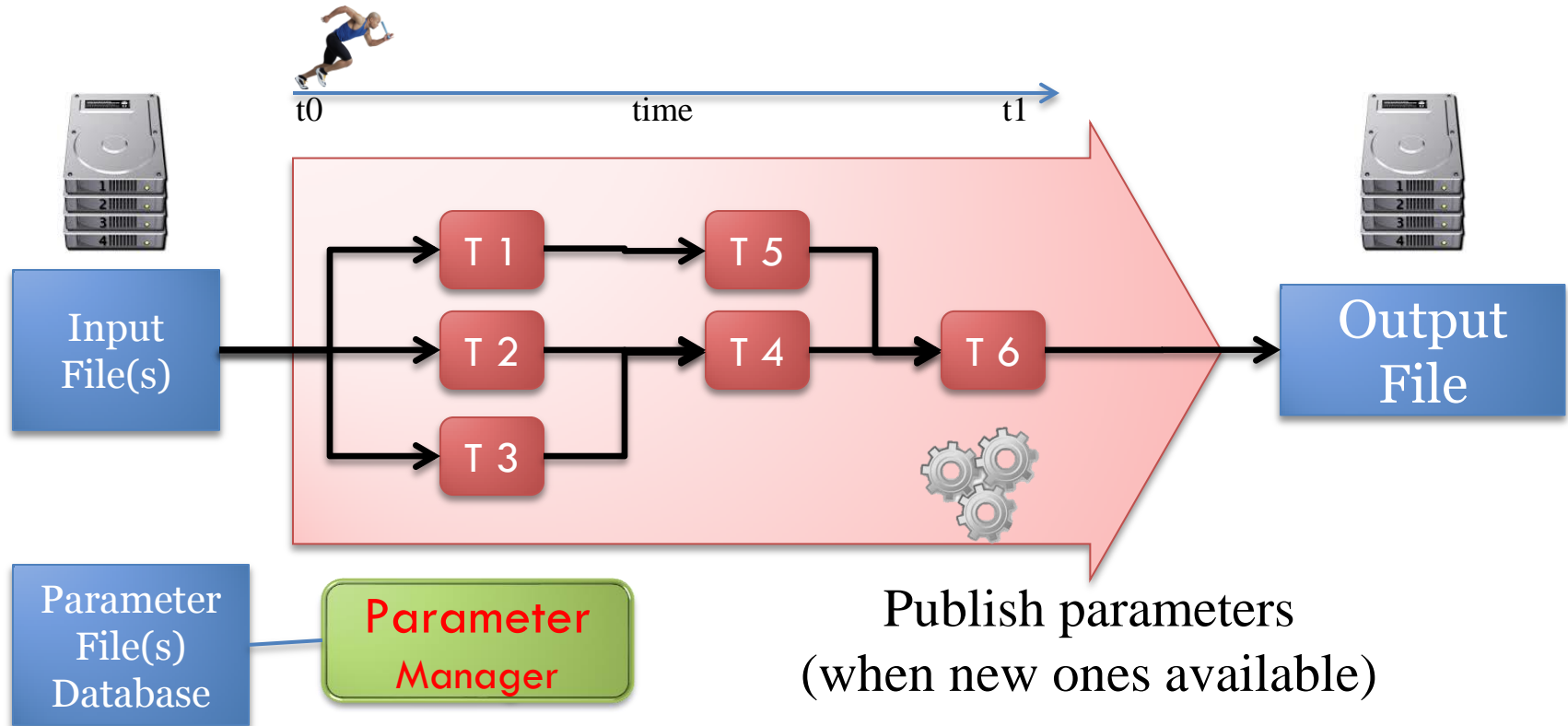# FairRoot: Where we are going ? (almost there!)

- Each Task is a process (can be Multi-threaded)

- Message Queues for data exchange

- Support multi-core and multi node



Publish parameters
(when new ones available)

# Before Re-inventing the Wheel

- What is available on the market and in the community?
  - A very promising package: ZeroMQ is available since 2011
- Do we intend to separate online and offline? NO
- 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.

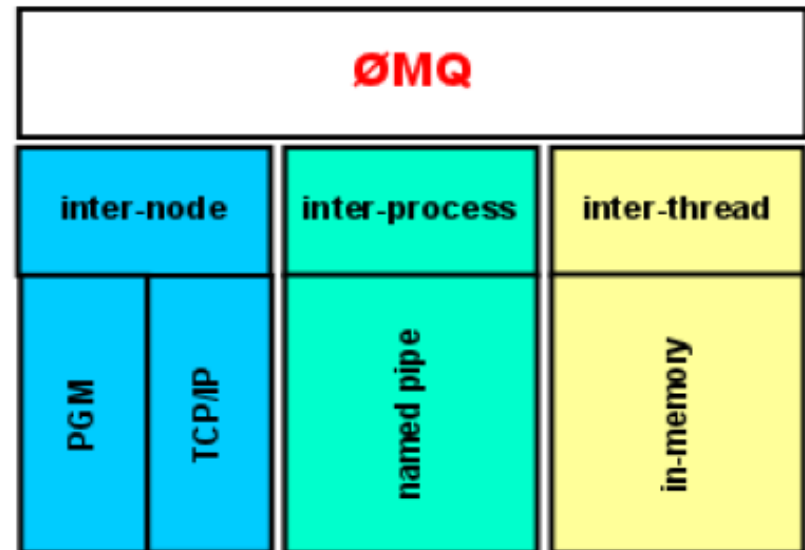- A messaging library, which allows you to design a complex communication system without much effort

- Abstraction on higher level than MPI (programming model is easier )

- Is suitable for loosely coupled and more general distributed systems

- Multiplatform, multi-language (+30)

- Small (20K lines of C++ code)

- Large and active open source community.

- Open source LGPL free software (large community)

# ZeroMQ sockets provide efficient transport options

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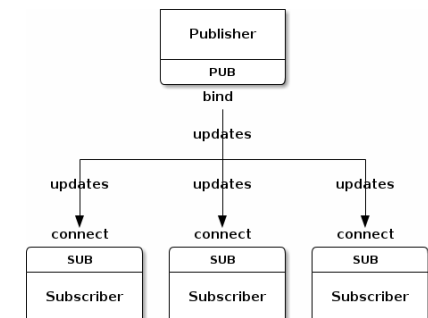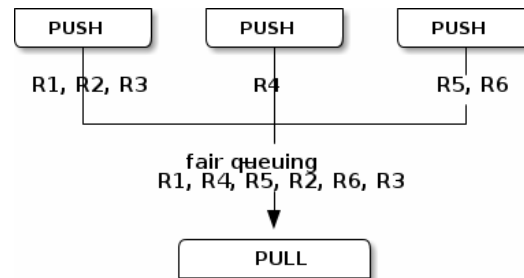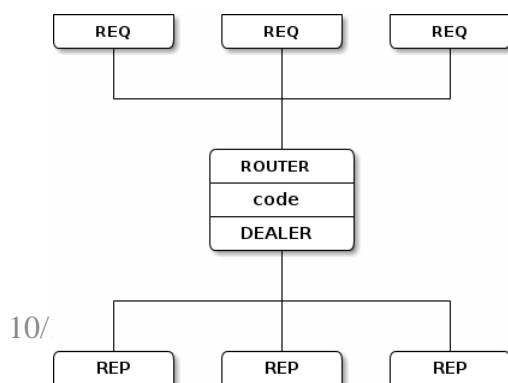  - which is really just inter-process across nodes communication



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. (remote procedure call and task distribution pattern)

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

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

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

# Current Status
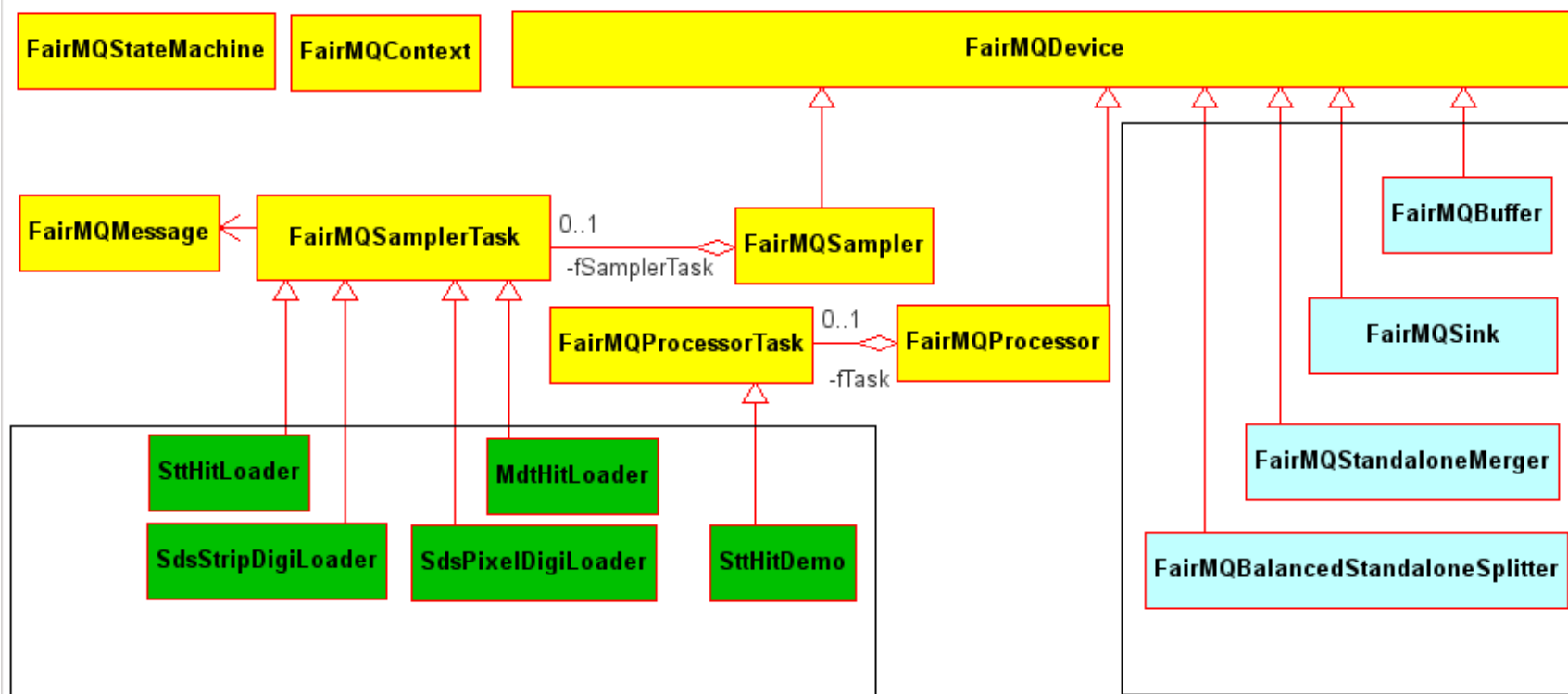
- The Framework delivers some components which can be connected to each other in order to construct a processing pipeline(s).
- All components share a common base called Device (ZeroMQ Class).
- Devices are grouped by three categories:
  - Source:
    - Data Sampler
  - Message-based Processor:
    - Sink, Splitter, Merger, Buffer, Proxy
  - Content-based Processor:
    - Processor

# Design

# Integrating the existing software:

ROOT Files, Lmd Files, Remote event server, ...

| FairRootManager | FairRunAna | FairTasks<br><br>Init()<br>Re-Init()<br>Exec()<br>Finish() | FairMQProcessorTask<br><br>Init()<br>Re-Init()<br>Exec()<br>Finish() |

Root (Event loop)

ZeroMQ

# FairRoot: Example 3

4 -Tracking stations with
a dipole field

Simulation:
A) 10k event: 10 Protons/ev
B) 20k event: 300 Protons/ev

Digitization

Reconstruction:
Hit/Cluster Finder

# From digits to hits with ROOT:



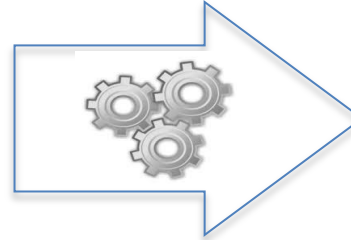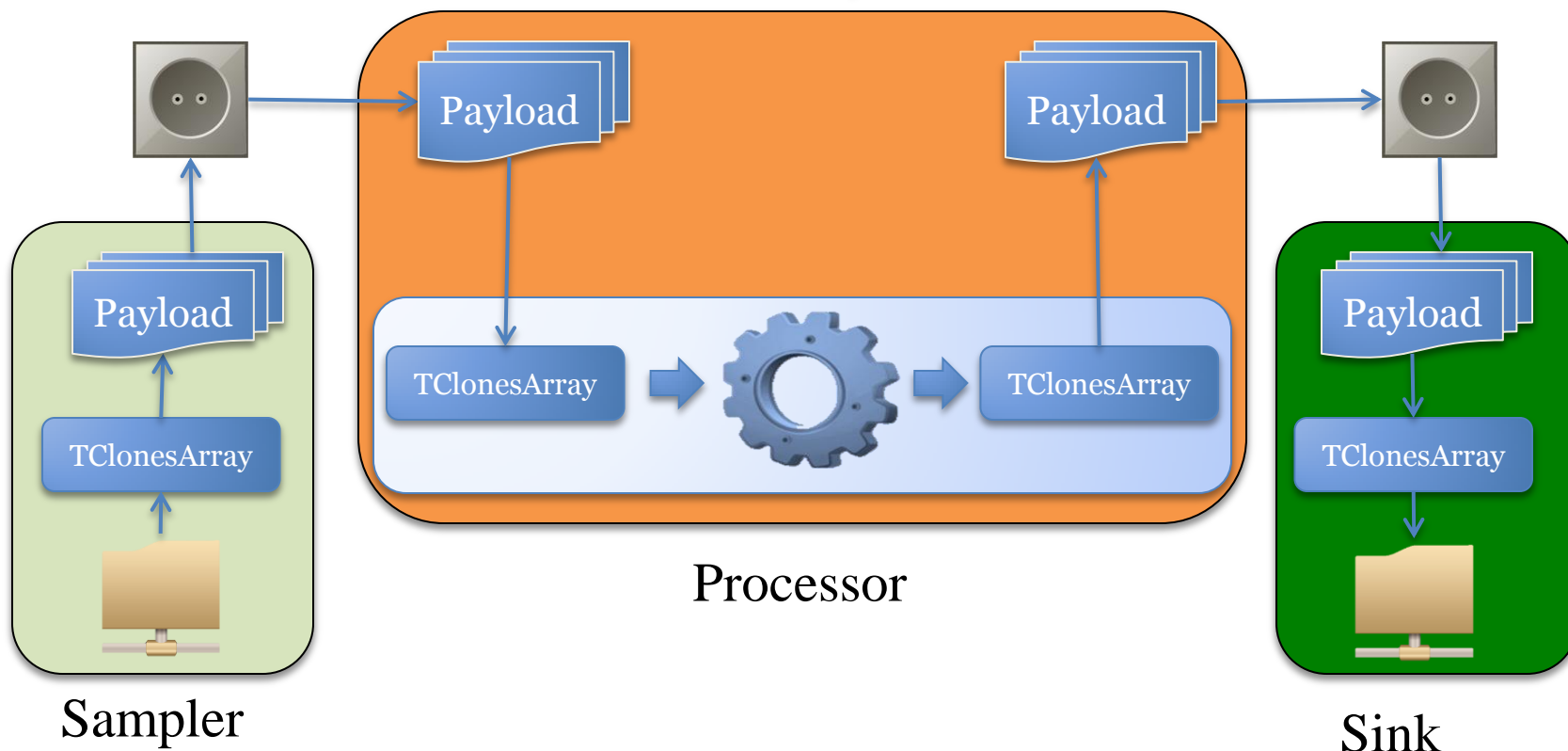| RUN | CPU Time (s) (Wall time) | Memory (Mbyte) |
|---|---|---|
| 10k Events, 10 Protons/event | 12 | 143 |
| 20k Events, 300 Protons/event | 162 | 241 |

# From digits to hits with **ØMQ** :



Digits → Hits

Sampler

Processor

Sink

# Test 1: Reconstruction
## 10k Event    10 Tracks/event



**root    11.85s    143MB**



| sample r | Push | processor | Push | sink |
|---|---|---|---|---|
| 6.51 s<br>135 MB | | 8.78 s<br>38 MB | | 3.24 s<br>38 MB |

CPU Time 25% less time

48 % more memory

| sample r | Push | processor | Push | sink |
|---|---|---|---|---|
| 6.96 s<br>135 MB | | 5.59 s<br>35MB | | 1.89 s<br>39 MB |
| | | 5.56 s<br>34MB | | 1.89 s<br>39 MB |

CPU Time 40% less time

97 % more memory

# Test 1: Reconstruction
## 20k Event 300 Tracks/event

root    162 s    241MB

CPU Time 14% less time

10 % more memory

| sampler | Push | processor | Push | sink |
|---------|------|-----------|------|------|
| 93 s | | 142 s | | 15 s |
| 151 MB | | 71 MB | | 41 MB |

CPU Time 36 % less time

18 % more memory

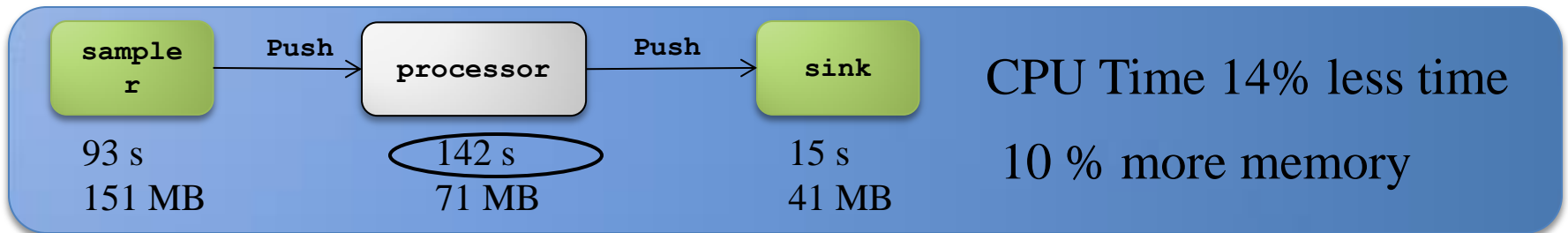| sampler | | processor | Push | sink |
|---------|------|-----------|------|------|
| | Push | 93 s | | 10 s |
| | | 35MB | | 40 MB |
| 103 s | | | | |
| 135MB | Push | processor | Push | sink |
| | | 94 s | | 10 s |
| | | 35MB | | 40 MB |

# Test 1: Reconstruction
## 20k Event   300 Tracks/event

root    162 s    241MB



| sampler | Pub Sub | processor | Pub Sub | sink |

94 s
150 MB

105 s
57 MB

11 s
40 MB

CPU Time 35% less time
2 % more memory

**32 %  events lost !**

sampler  splitter  →(Pub Sub)→ processor →(Pub sub)→ sink

104 s
135MB

3 s
33 MB

93 s
35MB

10 s
40 MB

93 s
36 MB

10 s
40 MB

CPU Time 36 % less time

32 % more memory

**0 %  events lost !**

ØMQ

Digits

Hits

Payload

Payload

Payload

Payload

TClonesArray → ⚙ → TClonesArray

TClonesArray

TClonesArray

Processor

Sampler

Sink

Overhead: Copy data from STL to TClonesArray and back

ØMQ  Digits ⟹ Hits

Payload → ⚙ → Payload

Processor

Payload

TClonesArray

Sampler

Payload

TClonesArray

Sink

M. Al-Turany, CHEP 2013 Amsterdam
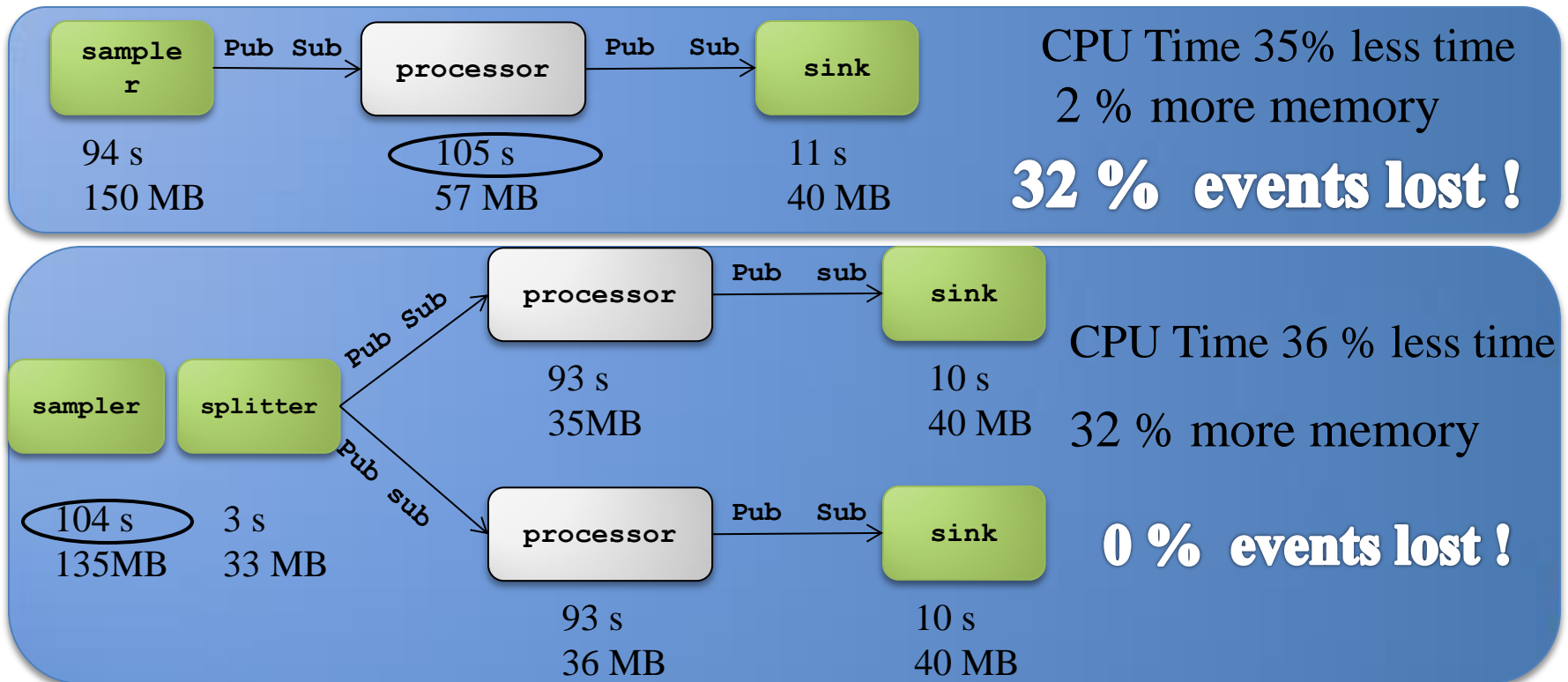
# Test 1: Reconstruction
## 20k Event 300 Tracks/event

root    162 s    241MB



sampler → Push → processor → Push → sink

73 s
135 MB

13 s
34 MB

14 s
82 MB

CPU Time 55% less time

4 % more memory

sampler → Push → processor → Push → sink

74 s
135MB

7 s
34MB

8 s
40 MB

sampler → Push → processor → Push → sink

7 s
34 MB

8 s
40 MB

CPU Time 55 % less time

17 % more memory

# Test 1: Reconstruction
# 20k Event  300 Tracks/event

**root   162 s   241MB**

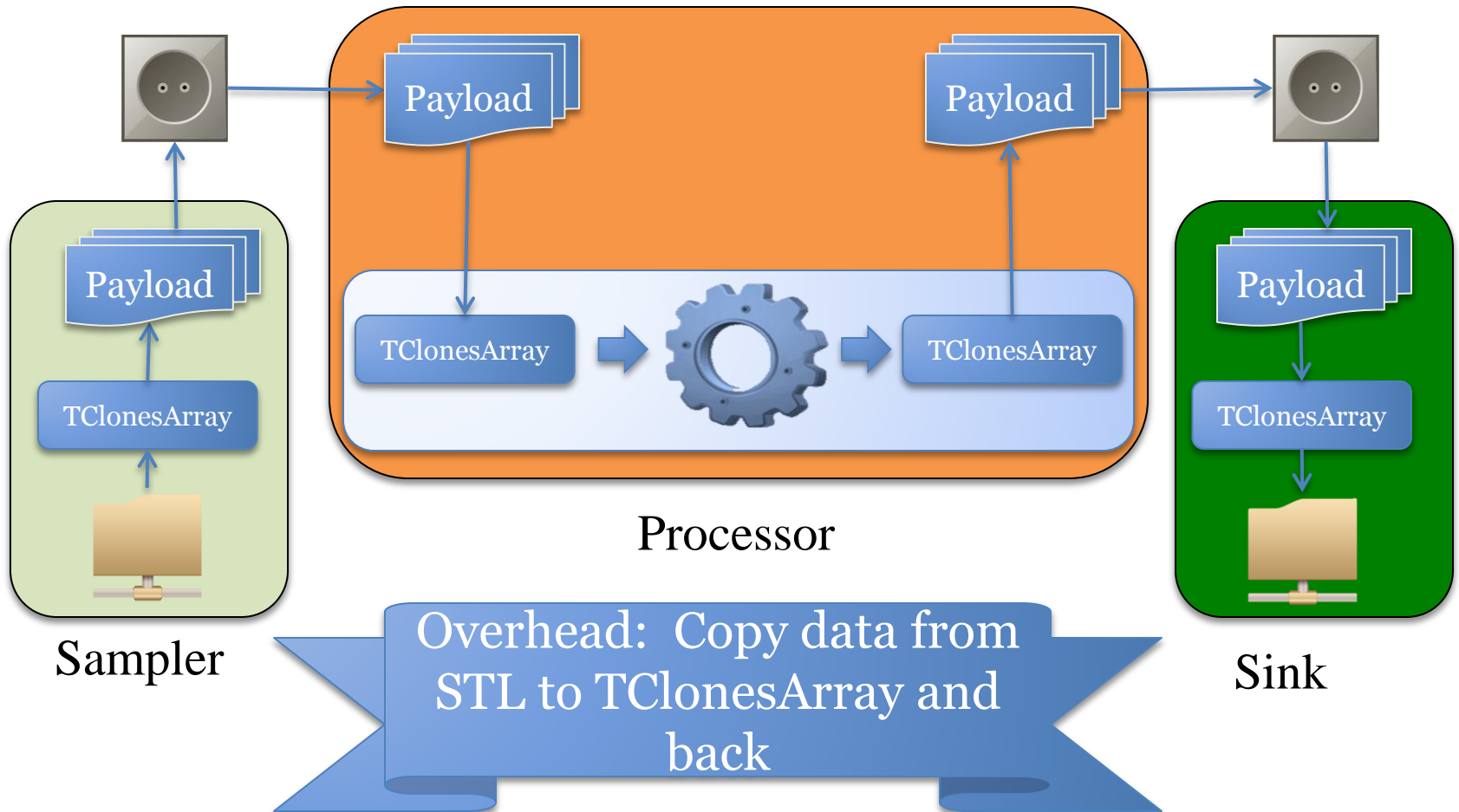

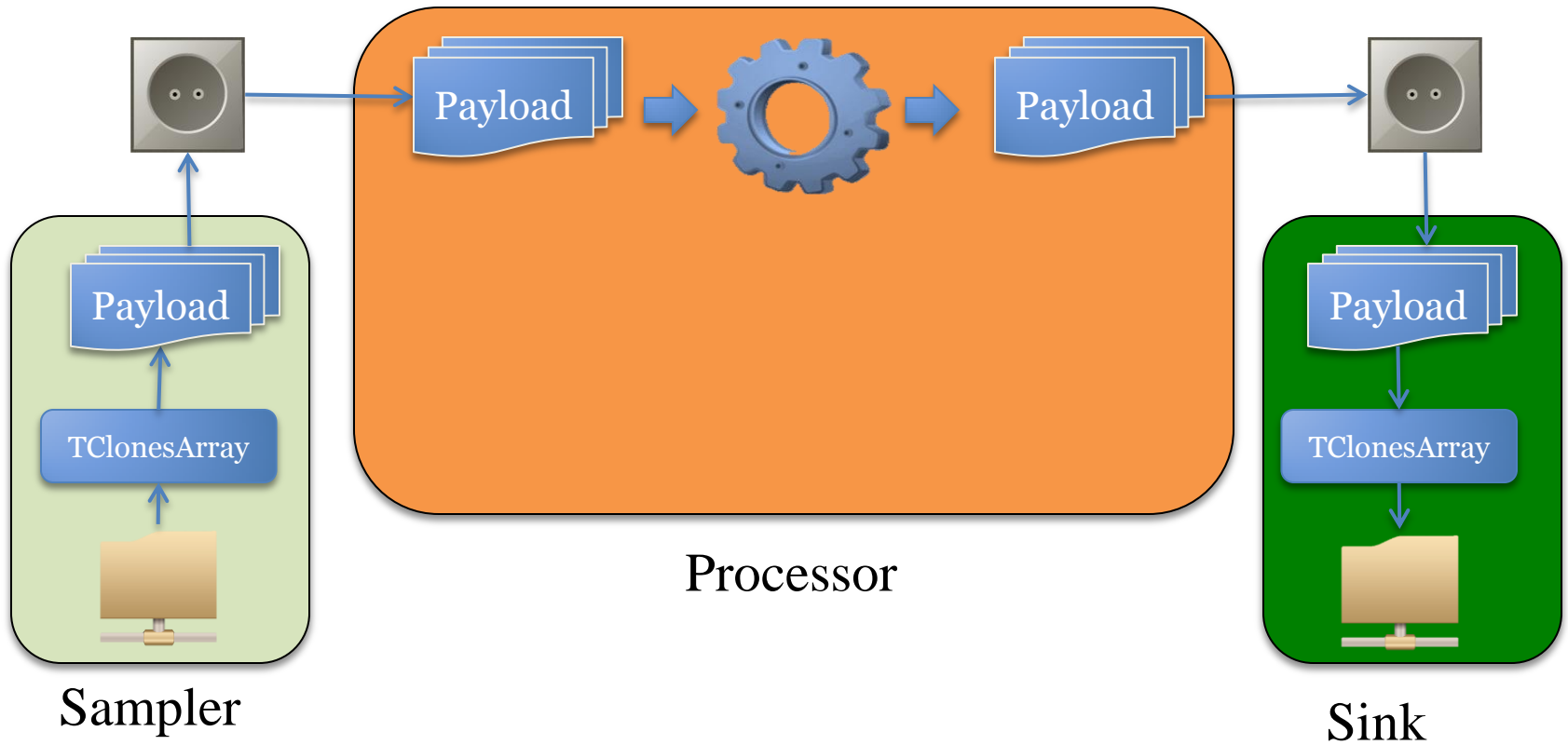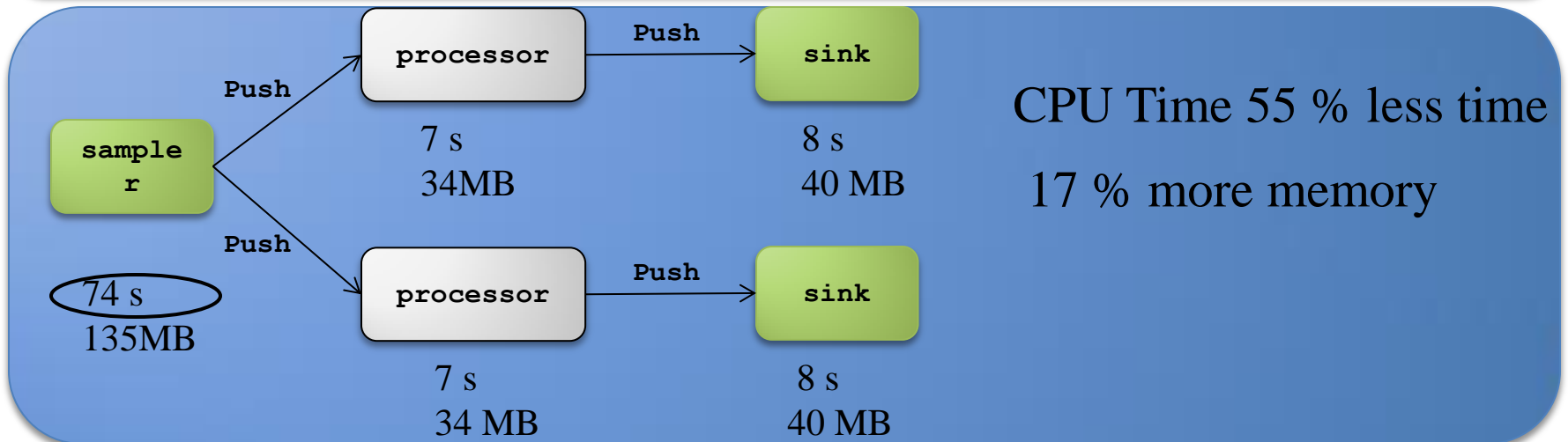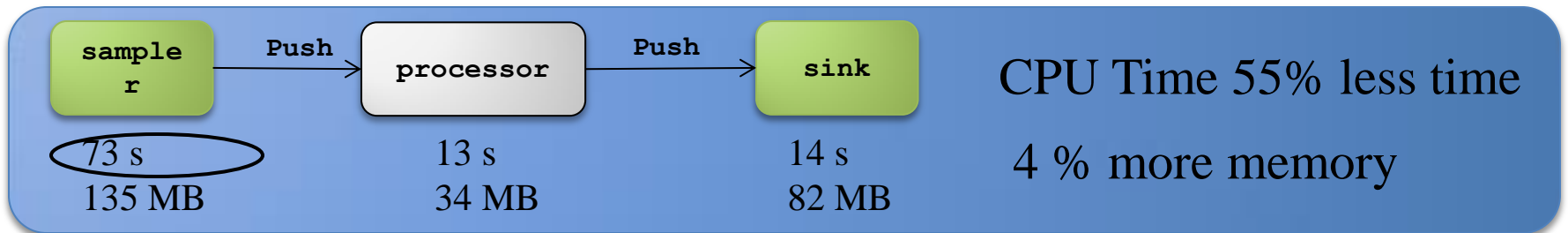| sampler | Pub Sub | processor | Pub Sub | sink |
|---------|---------|-----------|---------|------|
| 74 s | | 13 s | | 14 s |
| 135 MB | | 34 MB | | 83 MB |

CPU Time 54% less time
2 % more memory

**0 %  events lost !**

| sampler | splitter | Pub Sub | processor | Pub sub | sink |
|---------|----------|---------|-----------|---------|------|
| | | | 7 s | | 10 s |
| | | | 34 MB | | 39 MB |
| 78s | 3 s | Pub sub | processor | Pub Sub | sink |
| 135MB | 33 MB | | 7 s | | 10 s |
| | | | 34MB | | 39MB |

CPU Time 52 % less time

30 % more memory

**0 %  events lost !**

# Summary

- ZeroMQ communication layer is integrated into our offline framework (FairRoot).

- On the short term we will keep both options: ROOT based event loop and concurrent processes communicating with each other via ZeroMQ.

- On long term we are moving away from single event loop to distributed processes.

# Next Step: Design and development of a dynamic deployment system (DDS)

- STORM is very attractive but no native support for C++ !

- We need to utilize any RMS (Resource Management system)

- Support different topologies and process dependencies

- Device (process) is a single entity of the system

    o Each device has its own watchdog process

    o Devices are defined by a set of props and rules,

    o All devices are statically inherited (should support) 3 interfaces:
      IDDSConfig, IDDSStatus, and IDDSLog
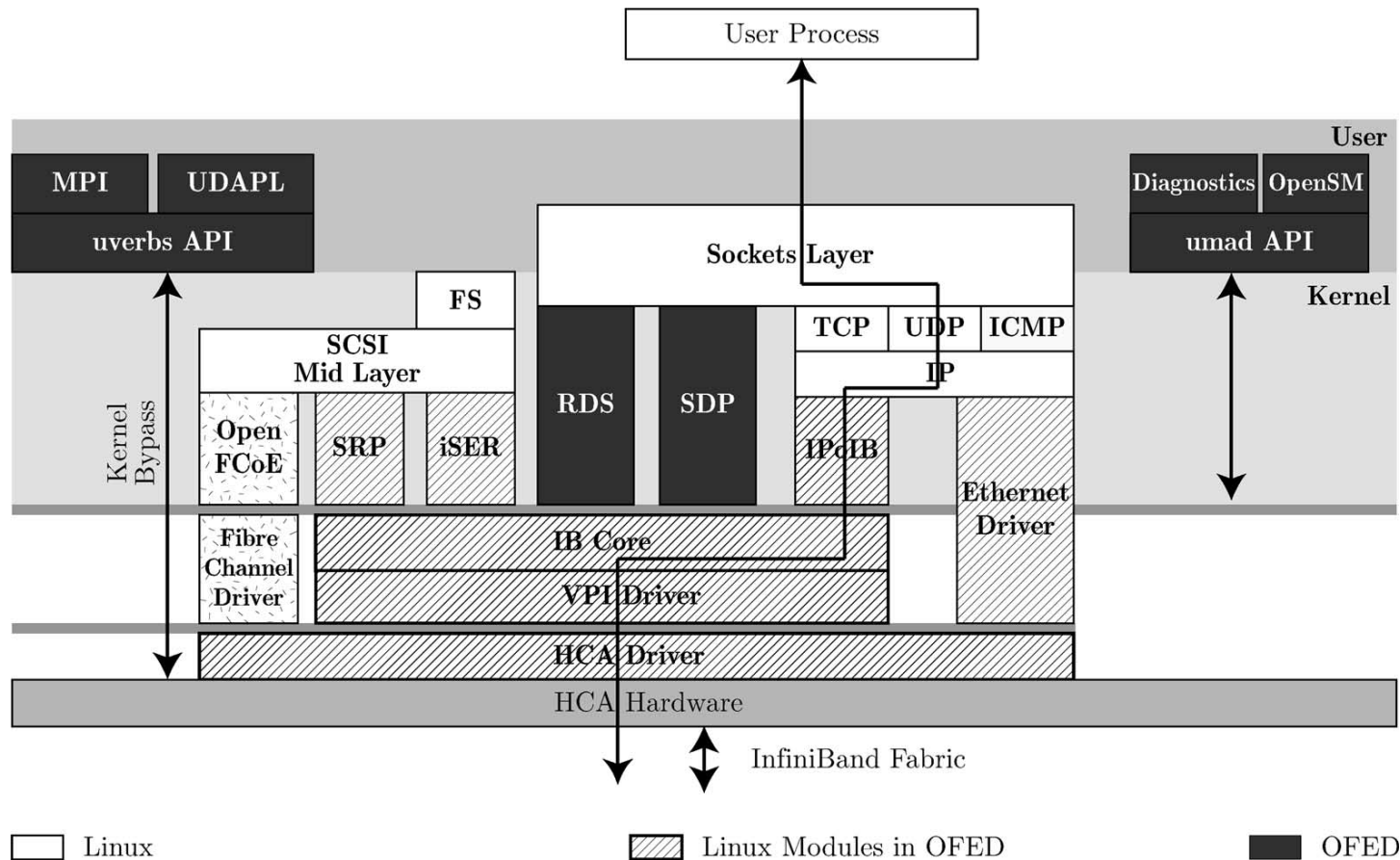
- .....



Thank you

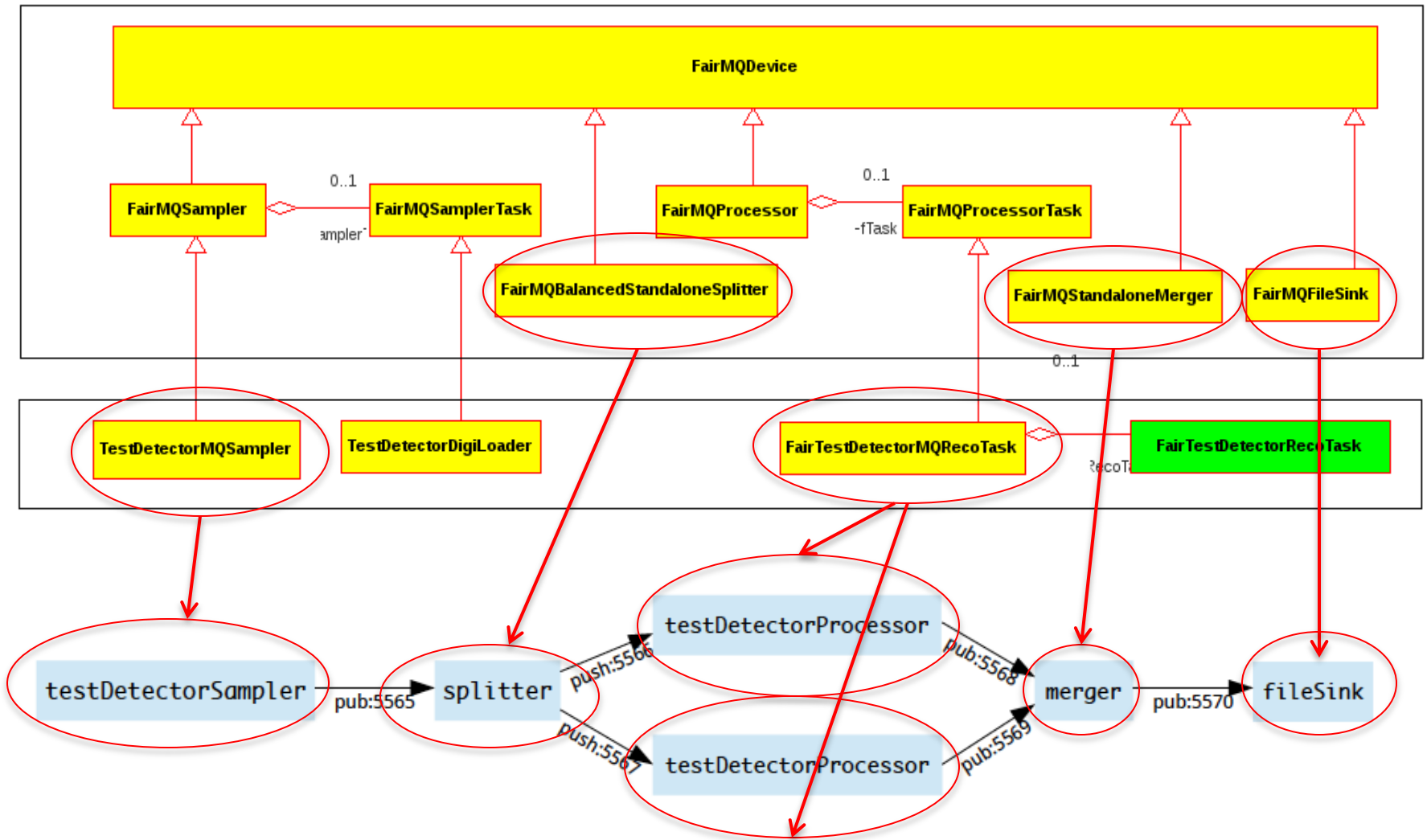# Backup

# Message format (Protocol)

- Potentially any content-based processor or any source can change the application protocol. Therefore, the framework provides a generic Message class that works with any arbitrary and continuous junk of memory (FairMQMessage).

- One has to pass a pointer to the memory buffer, the size in bytes, and can optionally pass a function pointer to a destructor, which will be called once the message object is discarded.
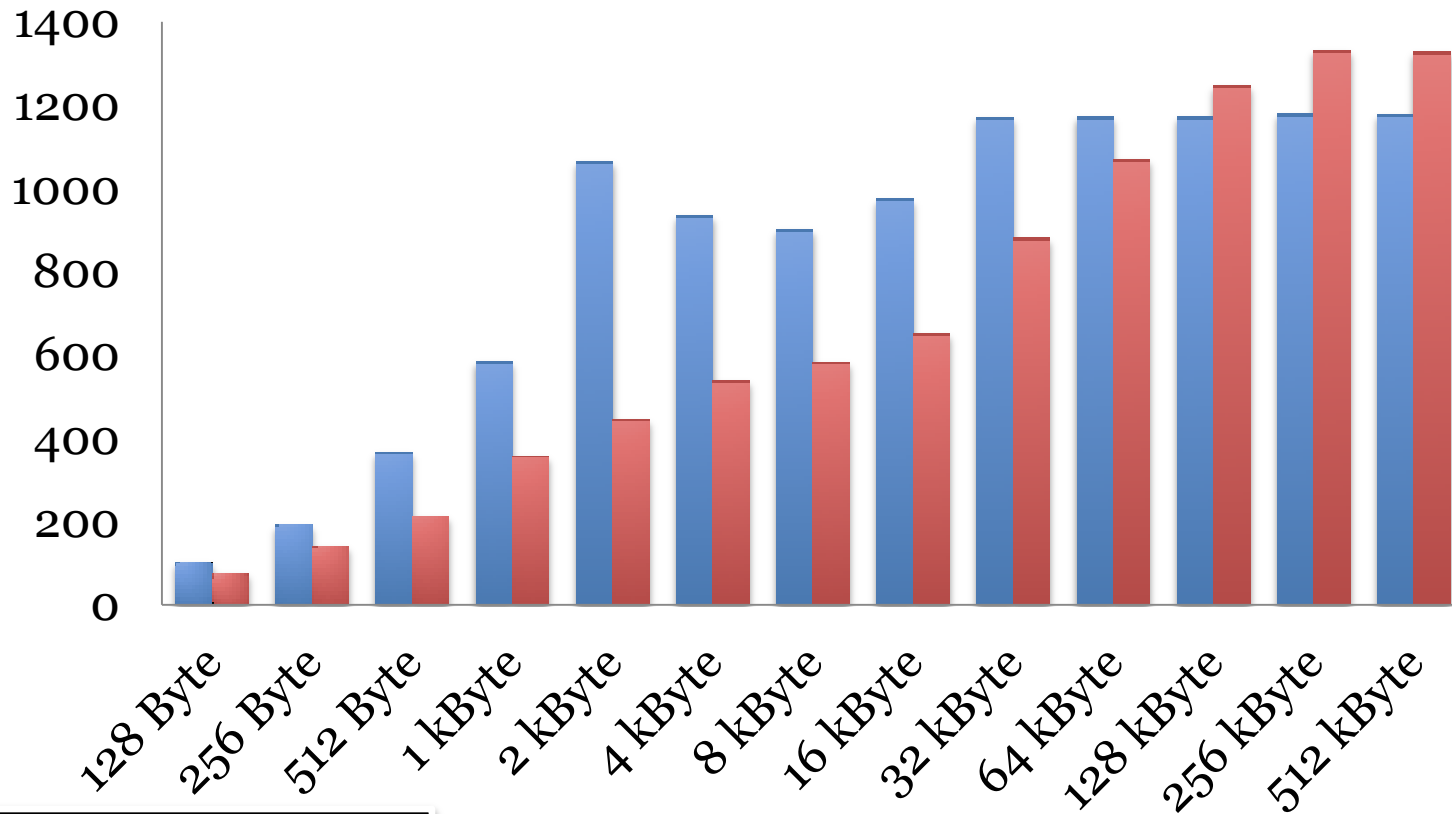
# Native InfiniBand/RDMA is faster than IP over IB



Implementing ZeroMQ over IB verbs will improve the performance.

# Fairbase/example/Tutorial3

# Payload in Mbyte/s as function of message size



Chart showing payload in Mbyte/s (y-axis 0 to 1400) as a function of message size (x-axis: 128 Byte, 256 Byte, 512 Byte, 1 kByte, 2 kByte, 4 kByte, 8 kByte, 16 kByte, 32 kByte, 64 kByte, 128 kByte, 256 kByte, 512 kByte). Legend: 10 Gbit (blue), 56 Gbit IB (red).

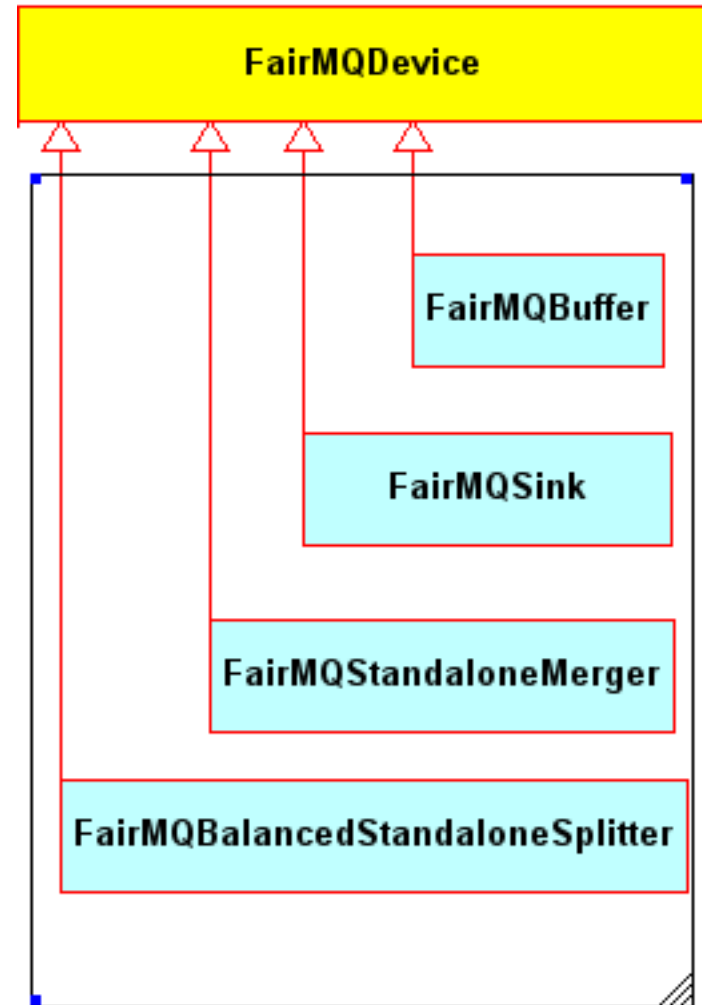**ZeroMQ works on InfiniBand but using IP over IB**

# Results

- Throughput of <span style="color:red">940 Mbit/s</span> was measured which is very close to the theoretical limit of the TCP/IPv4/GigabitEthernet

- 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>

Each message consists of digits in one event for one detector, with size of few kBytes
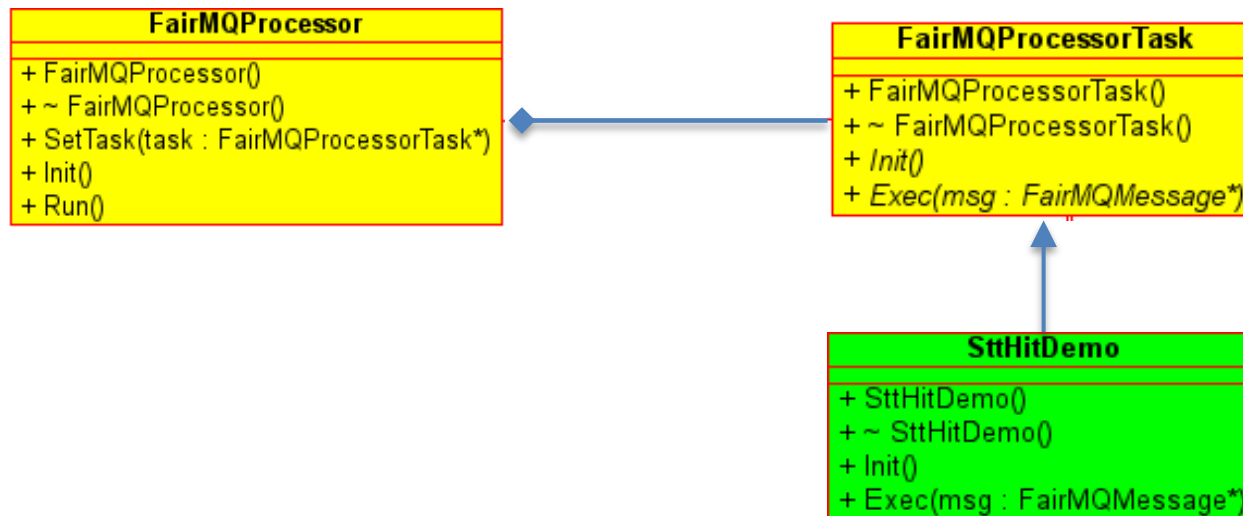
# Message-based Processor



- All message-based processors inherit from Device and operate on messages <span style="color:red">without interpreting their content</span>.

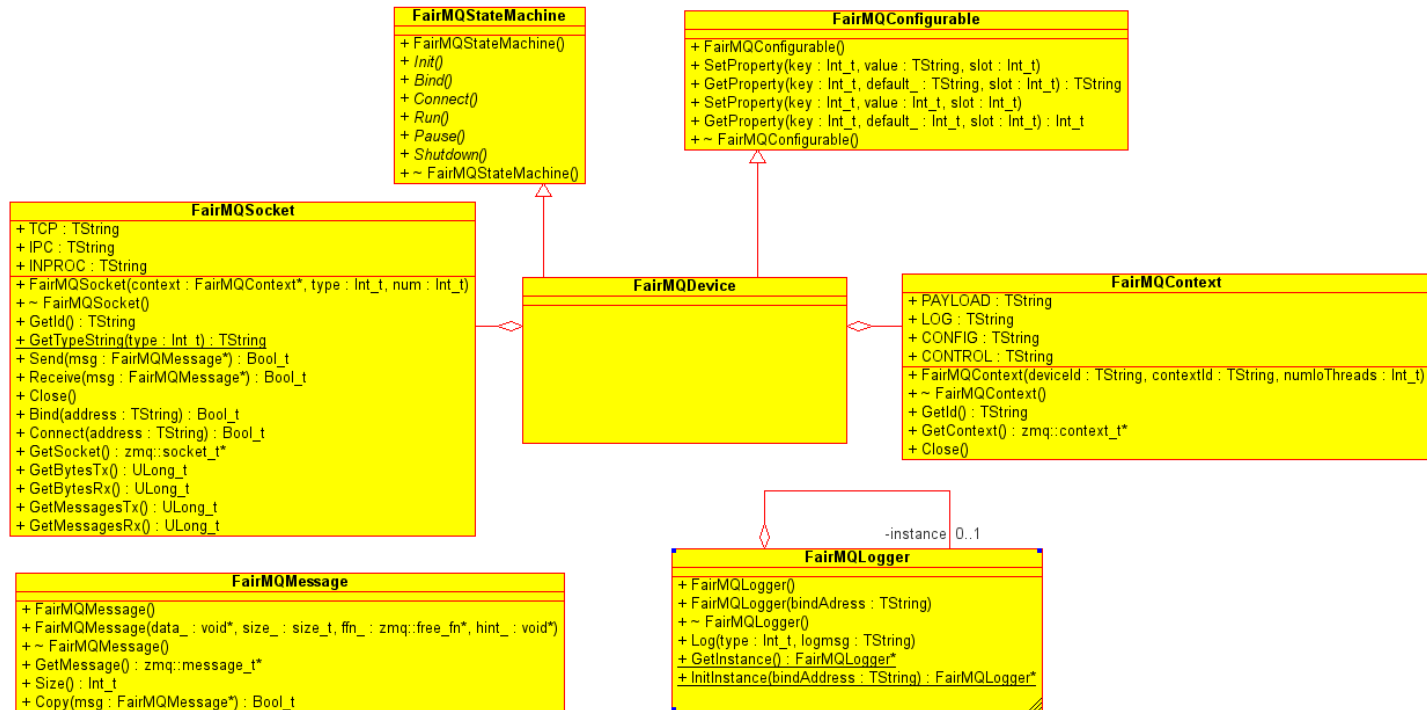- Four message-based processors have been implemented so far

# Content-based Processor

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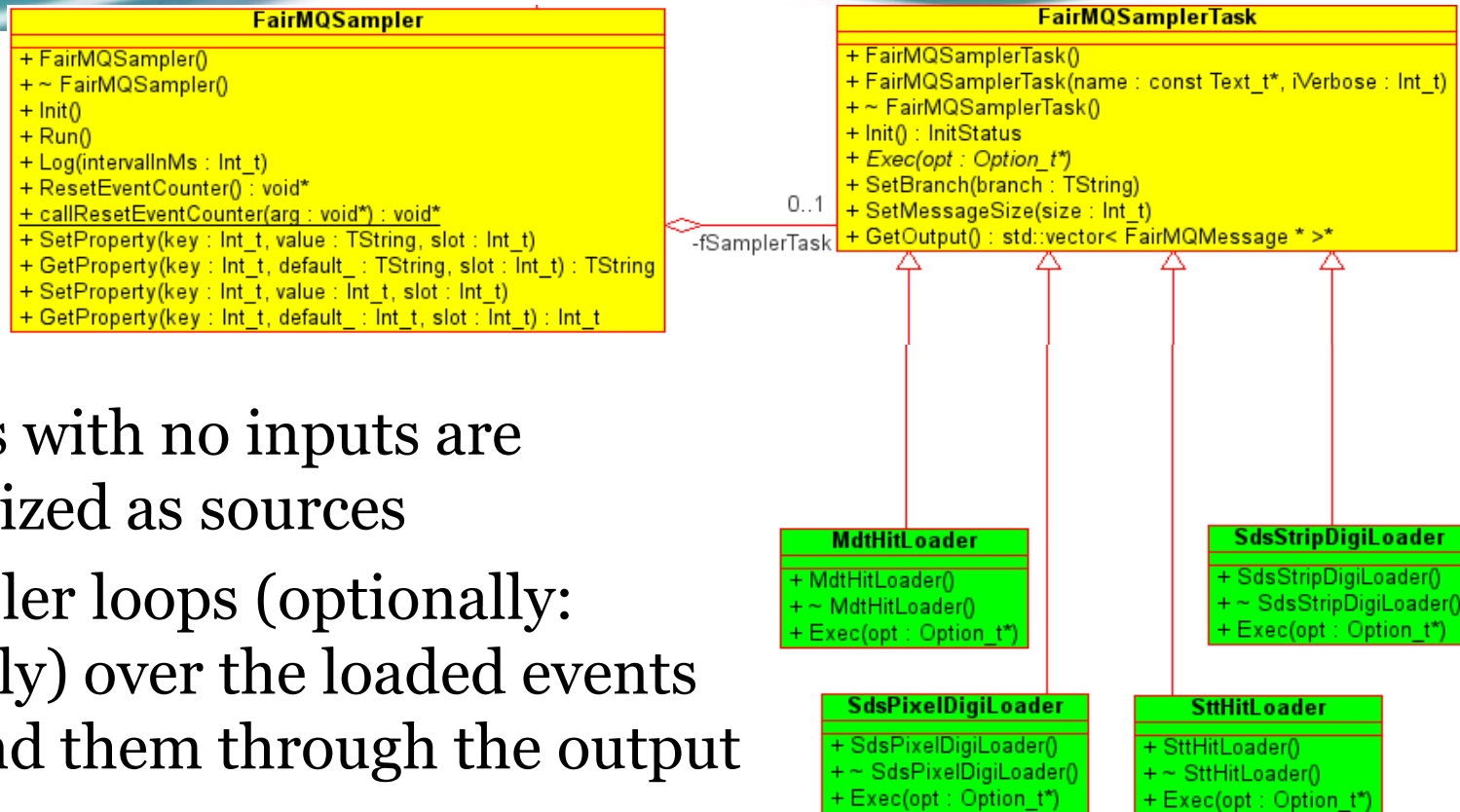
# Device

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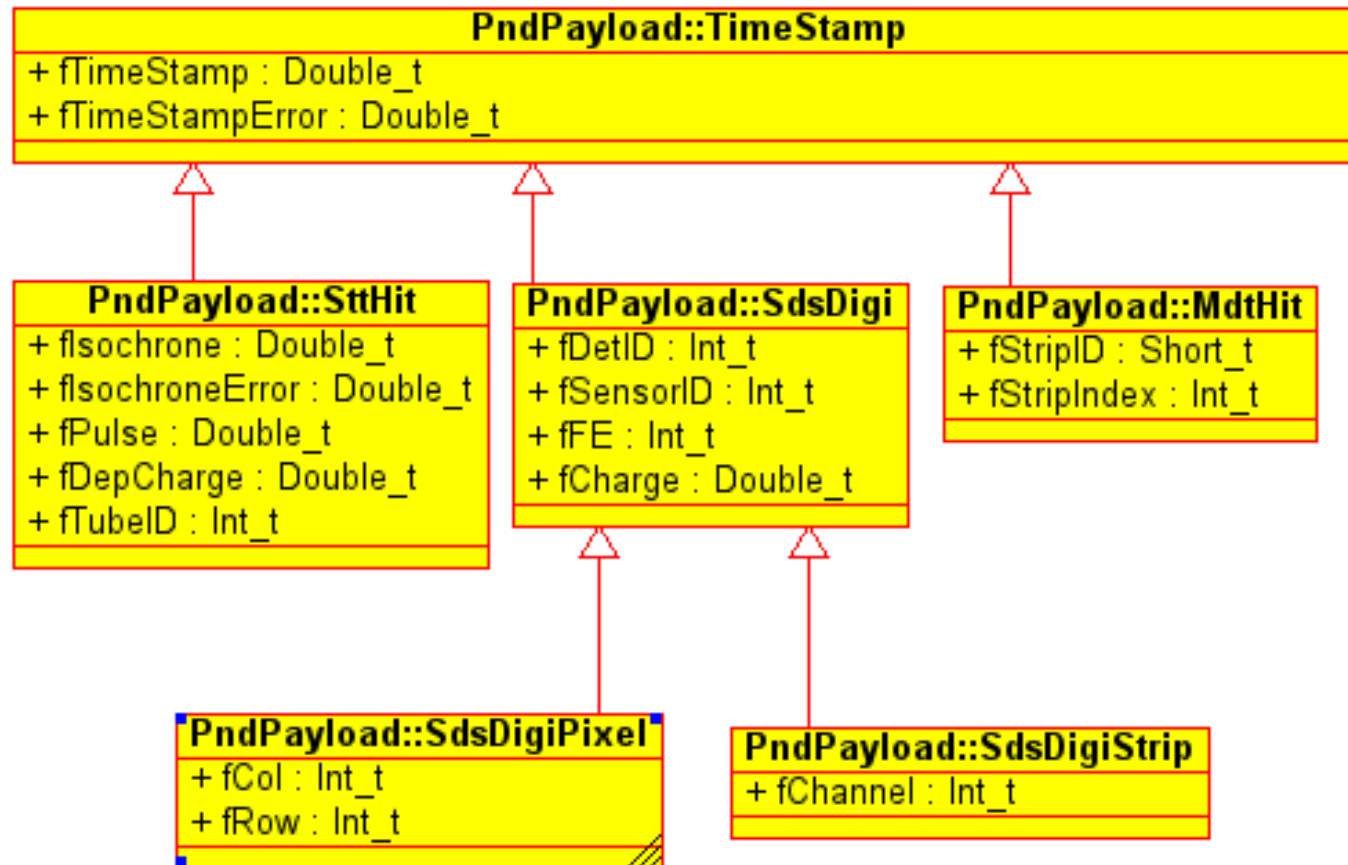
# Sampler

- Devices with no inputs are categorized as sources

- A sampler loops (optionally: 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

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

# ØMQ Features

- Message blobs of Zero to N bytes

- One socket connect to many sockets

- Queuing sender and receiver

- Automatic TCP (re)connect

- Zero-copy for large messages