

Evolution of the ALICE Software Framework for LHC Run 3

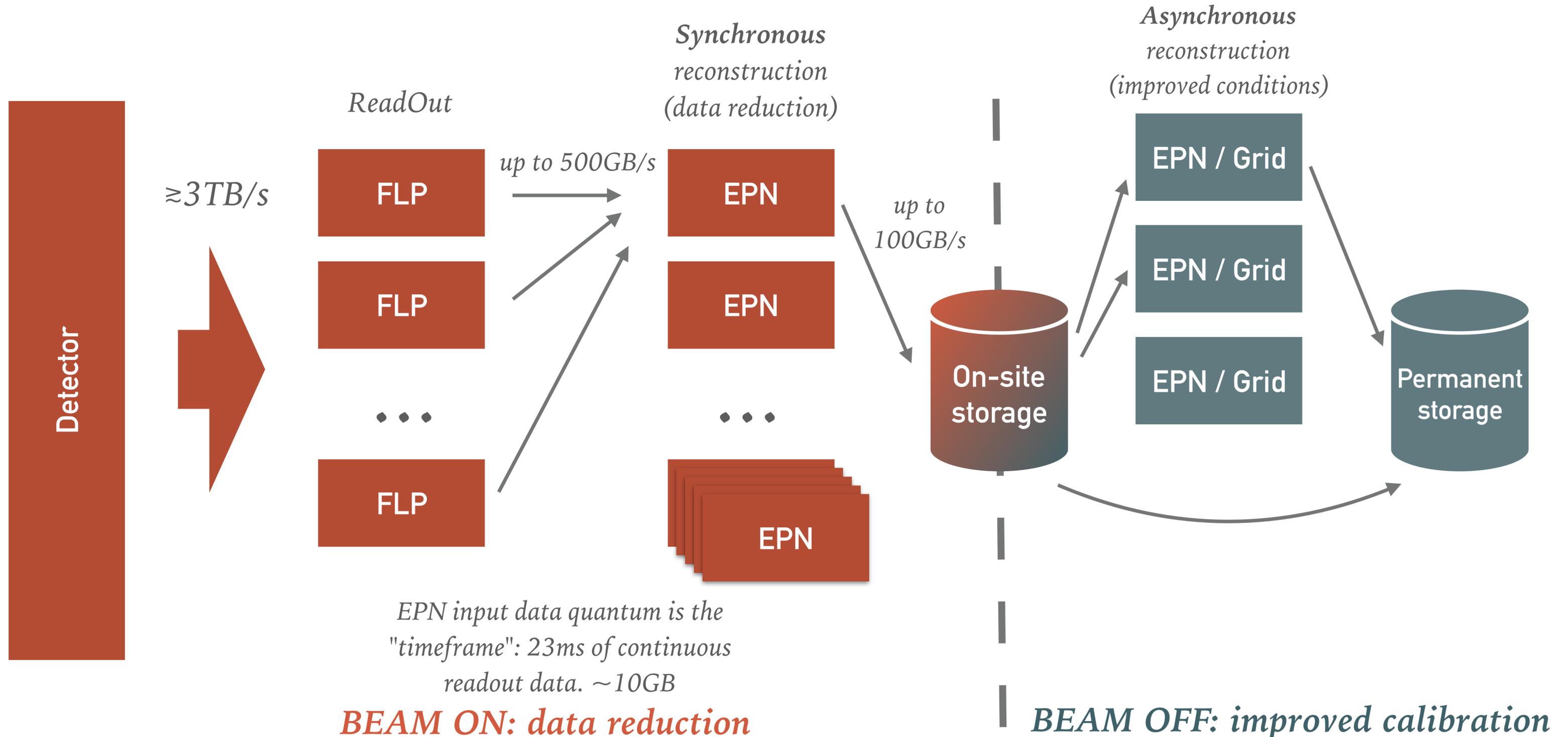


Giulio Eulisse (CERN), for the ALICE Collaboration



ALICE

ALICE IN RUN 3: POINT2



ALICE 02 SOFTWARE FRAMEWORK IN ONE SLIDE

Transport Layer: ALFA / FairMQ¹

- Standalone processes *for deployment flexibility.*
- Message passing *as a parallelism paradigm.*
- Shared memory *backend for reduced memory usage and improved performance.*

¹See "[ALFA: ALICE-FAIR new message queuing based framework](#)" by Mohammad Al Turani

ALICE 02 SOFTWARE FRAMEWORK IN ONE SLIDE

Data Layer: O2 Data Model

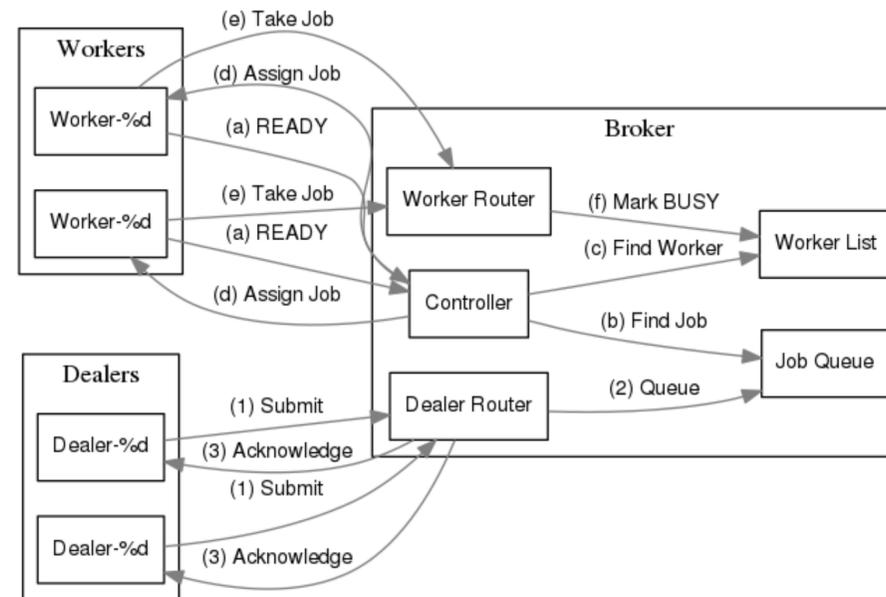
Message passing aware data model. Support for multiple backends:

- **Simplified, zero-copy format** optimised for performance and direct GPU usage. Useful e.g. for TPC reconstruction on the GPU.
- **ROOT based serialisation.** Useful for QA and final results.
- **Apache Arrow based.** Useful as backend of the analysis ntuples and for integration with other tools.

Transport Layer: ALFA / FairMQ¹

- **Standalone processes** for deployment flexibility.
- **Message passing** as a parallelism paradigm.
- **Shared memory backend** for reduced memory usage and improved performance.

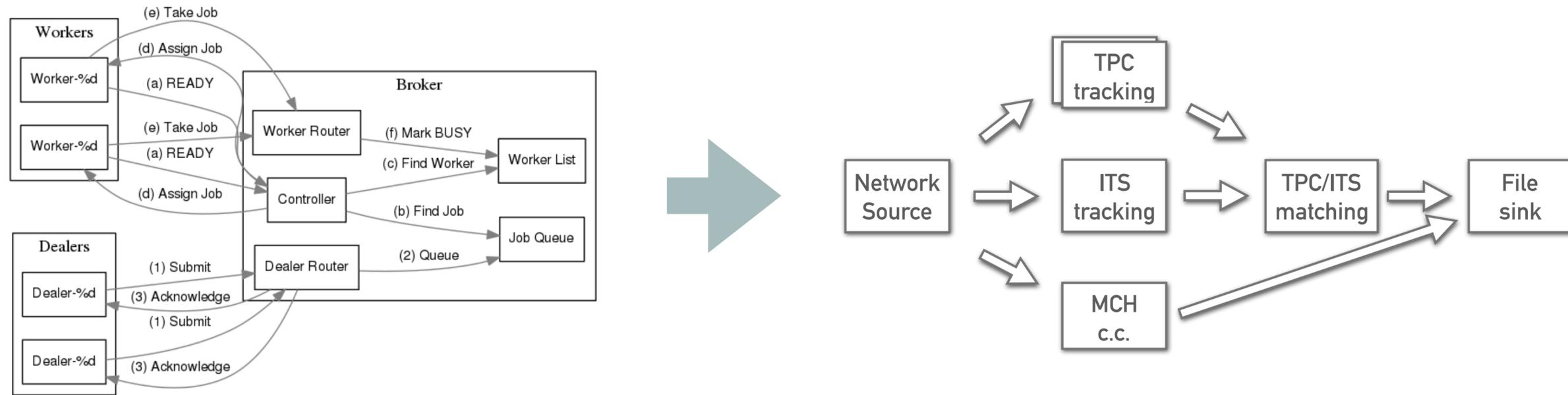
DISTRIBUTED SYSTEMS ARE HARD



There are only **two hard problems** in distributed systems:

2. Exactly-once delivery
1. Guaranteed order of messages
2. Exactly-once delivery

DISTRIBUTED SYSTEMS ARE HARD



There are only **two hard problems** in distributed systems:

2. Exactly-once delivery
1. Guaranteed order of messages
2. Exactly-once delivery

Since too many people did not get the joke, we started thinking how to simplify this for the user, as a result we decided to build a **data flow engine (pipelines!)** on top of our distributed system backend.

ALICE O2 SOFTWARE FRAMEWORK IN ONE SLIDE

Data Processing Layer (DPL)

Abstracts away the hiccups of a distributed system, presenting the user a familiar "Data Flow" system.

- *Reactive-like design (push data, don't pull)*
- *Declarative Domain Specific Language for topology configuration (C++17 based).*
- *Integration with the rest of the production system, e.g. Monitoring, Logging, Control.*
- *Laptop mode, including graphical debugging tools.*

Data Layer: O2 Data Model

Message passing aware data model. Support for multiple backends:

- **Simplified, zero-copy format optimised for performance and direct GPU usage.** Useful e.g. for TPC reconstruction on the GPU.
- **ROOT based serialisation.** Useful for QA and final results.
- **Apache Arrow based.** Useful as backend of the analysis ntuples and for integration with other tools.

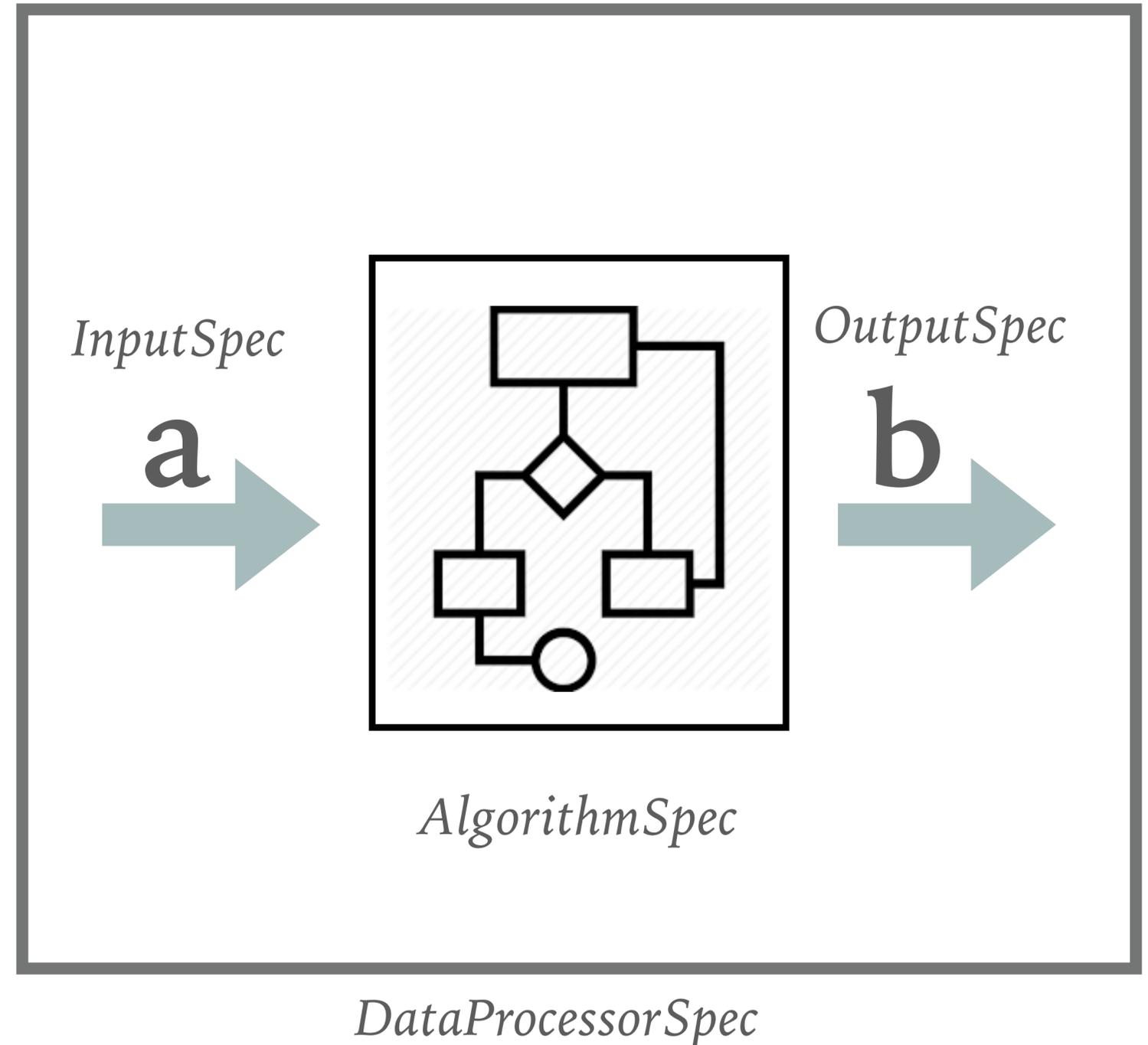
Transport Layer: ALFA / FairMQ¹

- *Standalone processes for deployment flexibility.*
- *Message passing as a parallelism paradigm.*
- *Shared memory backend for reduced memory usage and improved performance.*

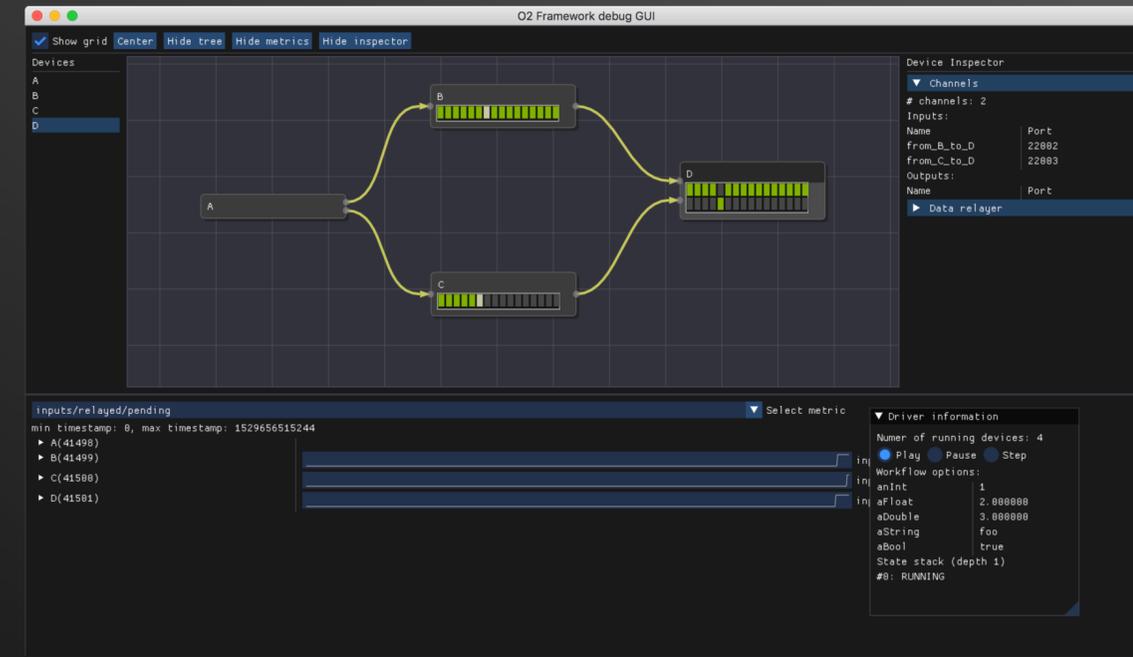
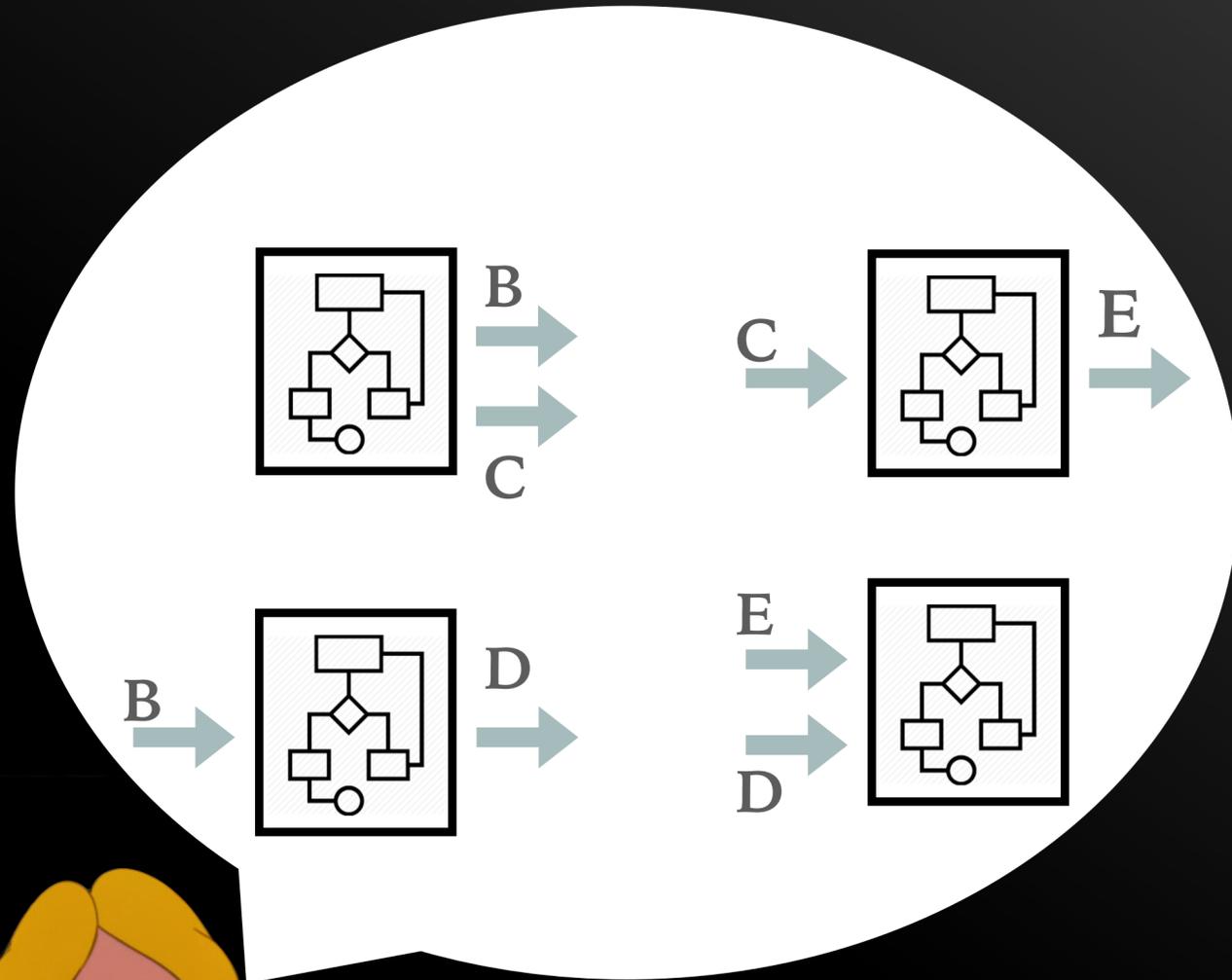
DATA PROCESSING LAYER: BUILDING BLOCK

A `DataProcessorSpec` defines a pipeline stage as a building block.

- Specifies *inputs and outputs* in terms of the O2 Data Model descriptors.
- Provide an implementation of how to act on the inputs to produce the output.
- Advanced user can express possible data or time parallelism opportunities.



DATA PROCESSING LAYER: IMPLICIT TOPOLOGY



Data Processing Layer

Topology is defined implicitly.

Topological sort ensures a viable dataflow is constructed (no cycles!).

Laptop users gets immediate feedback through the debug GUI.

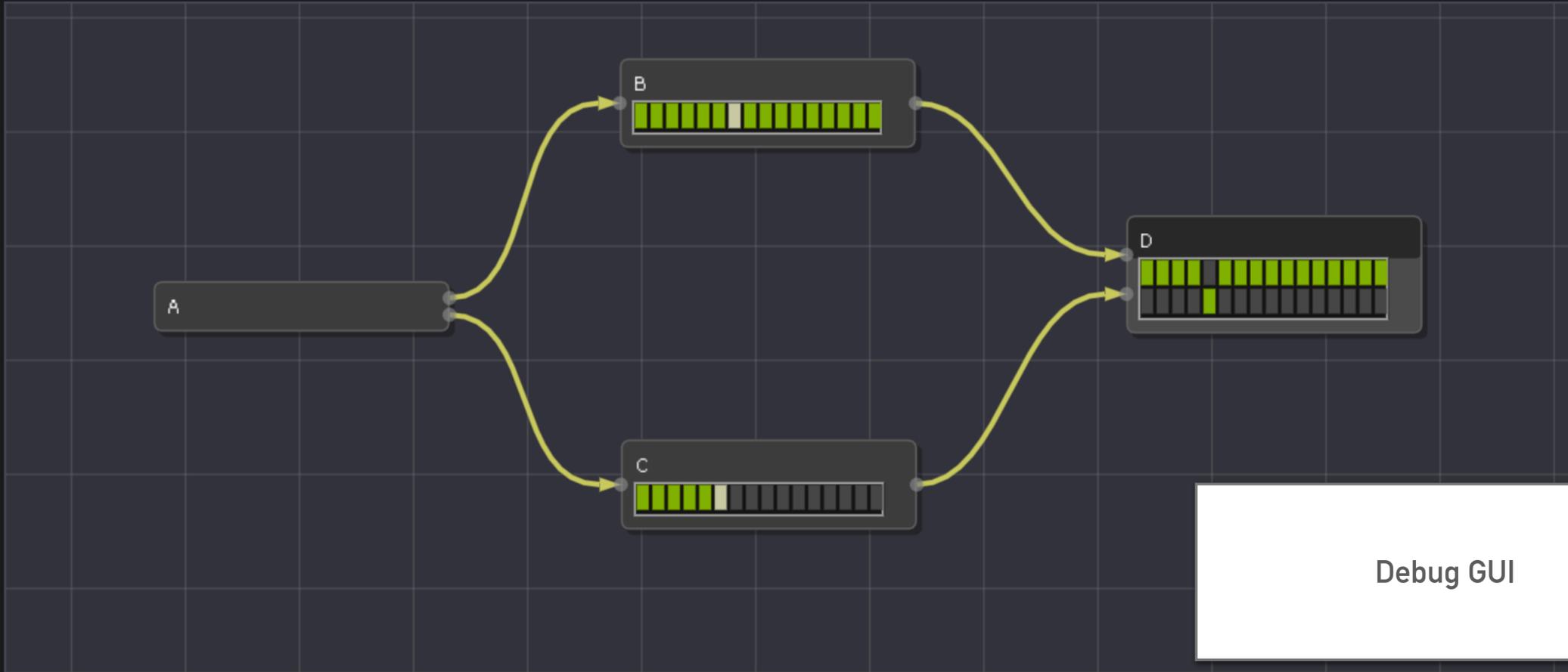
Service API allows integration with non data flow components (e.g. Control)



Show grid Center Hide tree Hide metrics Hide inspector

Devices

- A
- B
- C
- D



Device Inspector

Channels

channels: 2

Inputs:

Name	Port
from_B_to_D	22002
from_C_to_D	22003

Outputs:

Name	Port
------	------

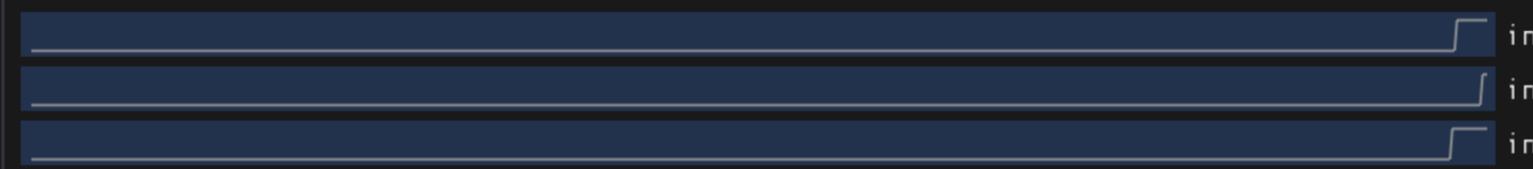
Data relayer

Debug GUI

inputs/relayed/pending Select metric

min timestamp: 0, max timestamp: 1529656515244

- ▶ A(41498)
- ▶ B(41499)
- ▶ C(41500)
- ▶ D(41501)



Driver information

Number of running devices: 4

Play Pause Step

Workflow options:

aInt	1
aFloat	2.000000
aDouble	3.000000
aString	foo
aBool	true

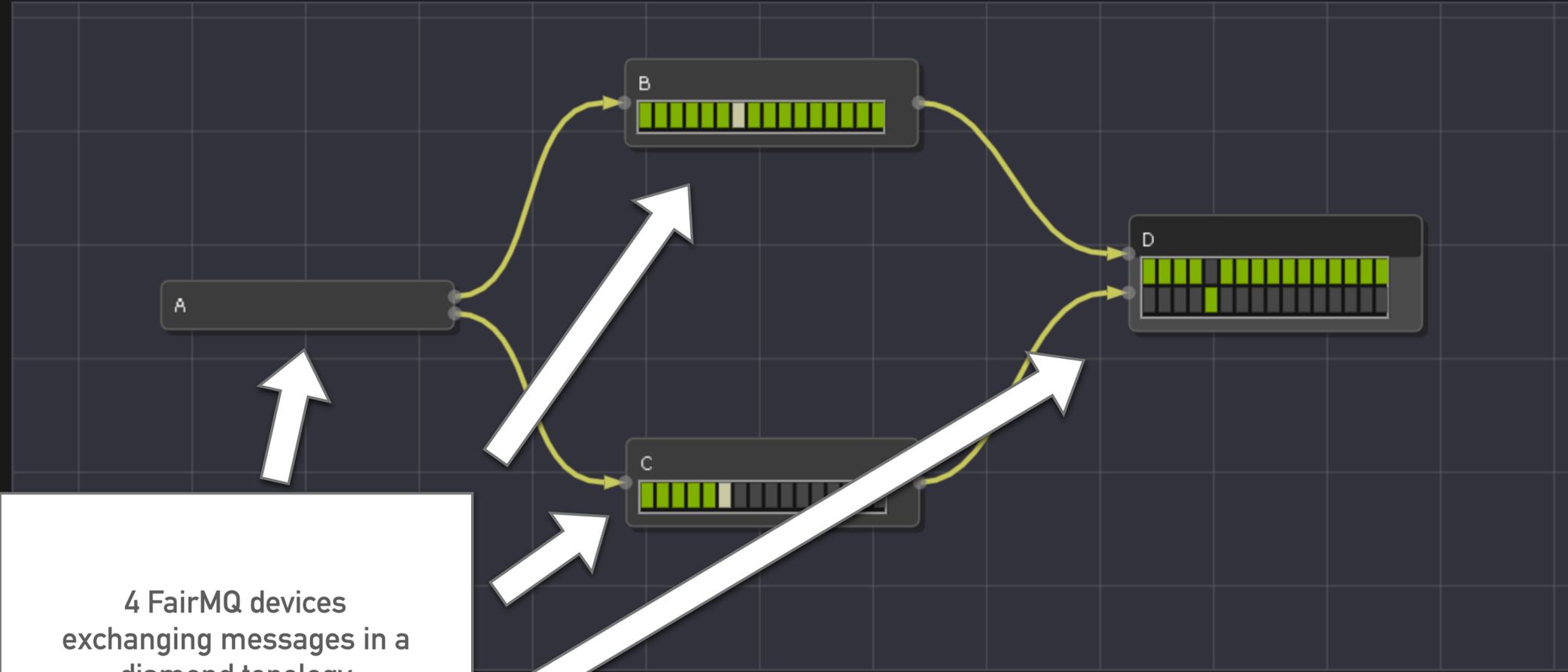
State stack (depth 1)

#0: RUNNING

Show grid Center Hide tree Hide metrics Hide inspector

Devices

- A
- B
- C
- D



4 FairMQ devices exchanging messages in a diamond topology

Device Inspector

Channels

channels: 2

Inputs:

Name	Port
from_B_to_D	22002
from_C_to_D	22003

Outputs:

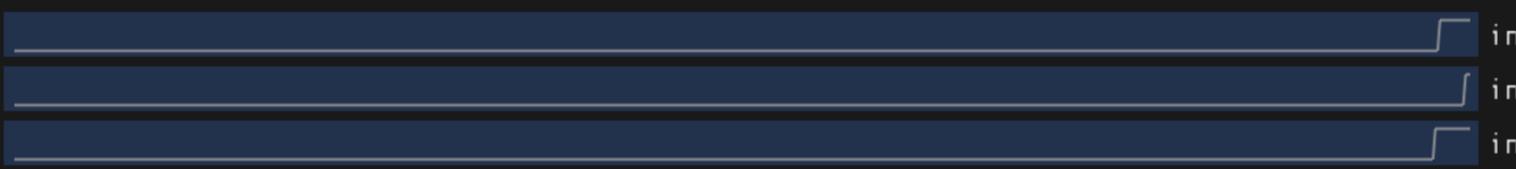
Name	Port

Data relayer

inputs/relayed

- min timestamp:
- A(41498)
 - B(41499)
 - C(41500)
 - D(41501)

Select metric



Driver information

Numer of running devices: 4

Play Pause Step

Workflow options:

aInt	1
aFloat	2.000000
aDouble	3.000000
aString	foo
aBool	true

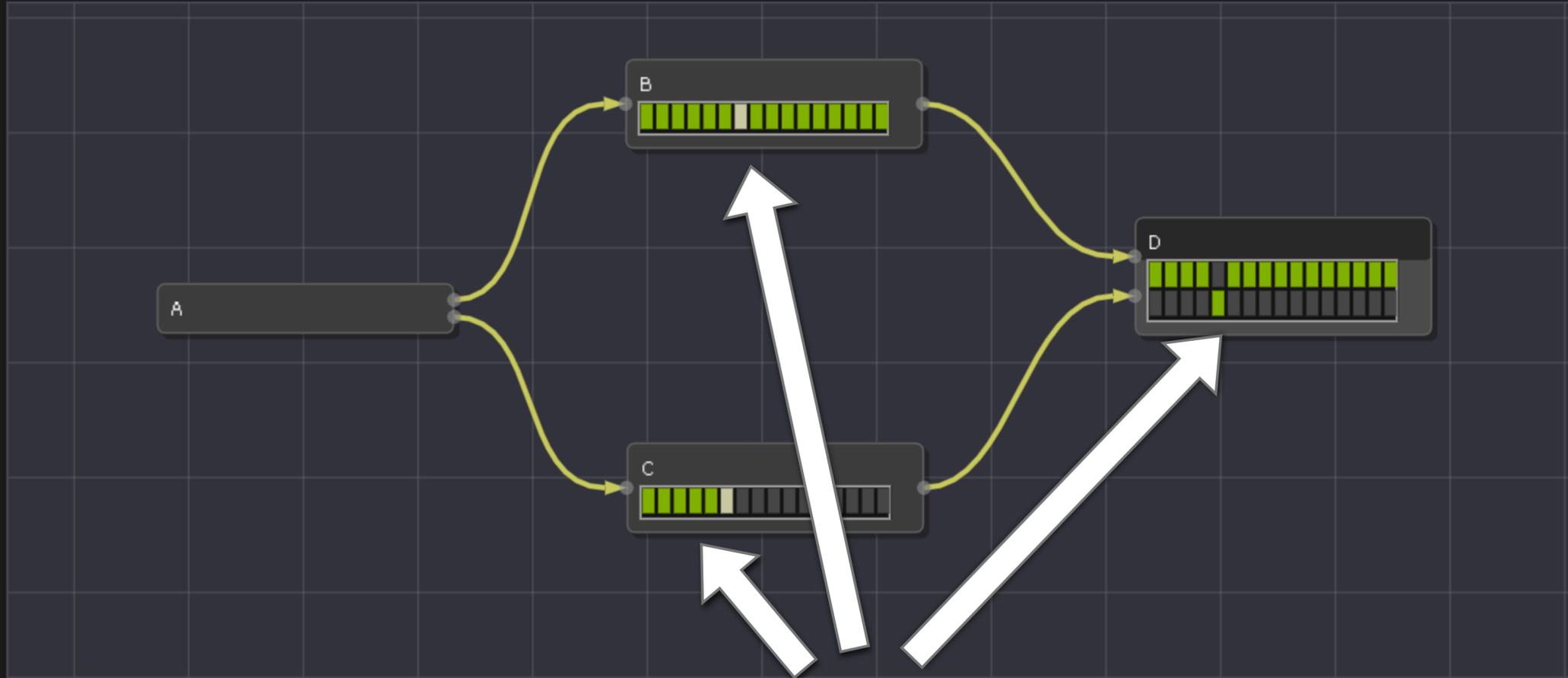
State stack (depth 1)

#0: RUNNING

Show grid Center Hide tree Hide metrics Hide inspector

Devices

- A
- B
- C
- D



Device Inspector

▼ Channels

channels: 2

Inputs:

Name	Port
from_B_to_D	22002
from_C_to_D	22003

Outputs:

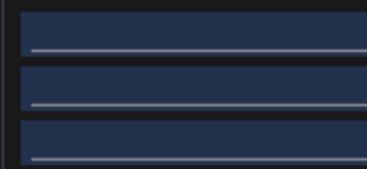
Name	Port
------	------

▶ Data relayer

inputs/relayed/pending

min timestamp: 0, max timestamp: 1529656515244

- ▶ A(41498)
- ▶ B(41499)
- ▶ C(41500)
- ▶ D(41501)



GUI shows state of the various message queues in realtime. Different colors mean different state of data processing.

▼ Select metric

▼ Driver information

Number of running devices: 4

Play Pause Step

Workflow options:

aInt	1
aFloat	2.000000
aDouble	3.000000
aString	foo
aBool	true

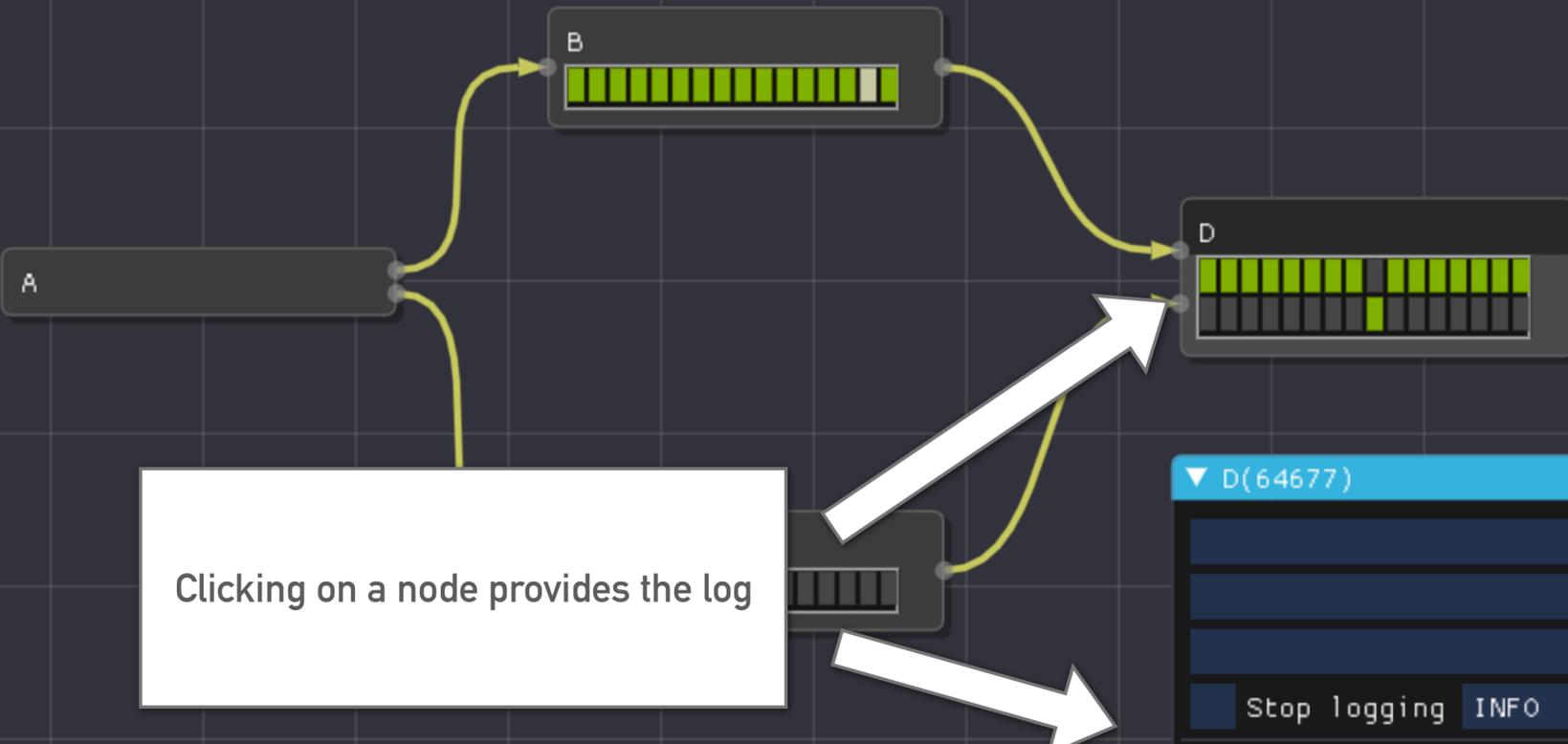
State stack (depth 1)

#0: RUNNING

Show grid Center Hide tree Hide metrics Hide inspector

Devices

- A
- B
- C
- D



Device Inspector

▼ Channels

channels: 2

Inputs:

Name	Port
from_B_to_D	22002
from_C_to_D	22003

Outputs:

Name	Port
▶ Data relayer	

▼ D(64677)

Log filter

Log start trigger

Log stop trigger

Stop logging INFO ▼ Log level

```
[10:53:30][INFO] from_C_to_D[0]: in: 0 (0 MB) out: 0 (0 MB)
[10:53:30][INFO] from_B_to_D[0]: in: 0.999001 (0.000131868 MB) out: 0 (0 MB)
[10:53:31][INFO] from_C_to_D[0]: in: 0 (0 MB) out: 0 (0 MB)
[10:53:31][INFO] from_B_to_D[0]: in: 0 (0 MB) out: 0 (0 MB)
[10:53:32][INFO] from_C_to_D[0]: in: 1 (0.000132 MB) out: 0 (0 MB)
[10:53:32][INFO] from_B_to_D[0]: in: 0 (0 MB) out: 0 (0 MB)
[10:53:33][INFO] from_C_to_D[0]: in: 0 (0 MB) out: 0 (0 MB)
[10:53:33][INFO] from_B_to_D[0]: in: 1 (0.000132 MB) out: 0 (0 MB)
[10:53:34][INFO] from_C_to_D[0]: in: 0 (0 MB) out: 0 (0 MB)
[10:53:34][INFO] from_B_to_D[0]: in: 0 (0 MB) out: 0 (0 MB)
[10:53:35][INFO] from_C_to_D[0]: in: 0 (0 MB) out: 0 (0 MB)
[10:53:35][INFO] from_B_to_D[0]: in: 0 (0 MB) out: 0 (0 MB)
[10:53:36][INFO] from_C_to_D[0]: in: 0 (0 MB) out: 0 (0 MB)
[10:53:36][INFO] from_B_to_D[0]: in: 1 (0.000132 MB) out: 0 (0 MB)
[10:53:37][INFO] from_C_to_D[0]: in: 0.995025 (0.000131343 MB) out: 0 (0 MB)
[10:53:37][INFO] from B to D[0]: in: 1.99005 (0.000262687 MB) out: 0 (0 MB)
```

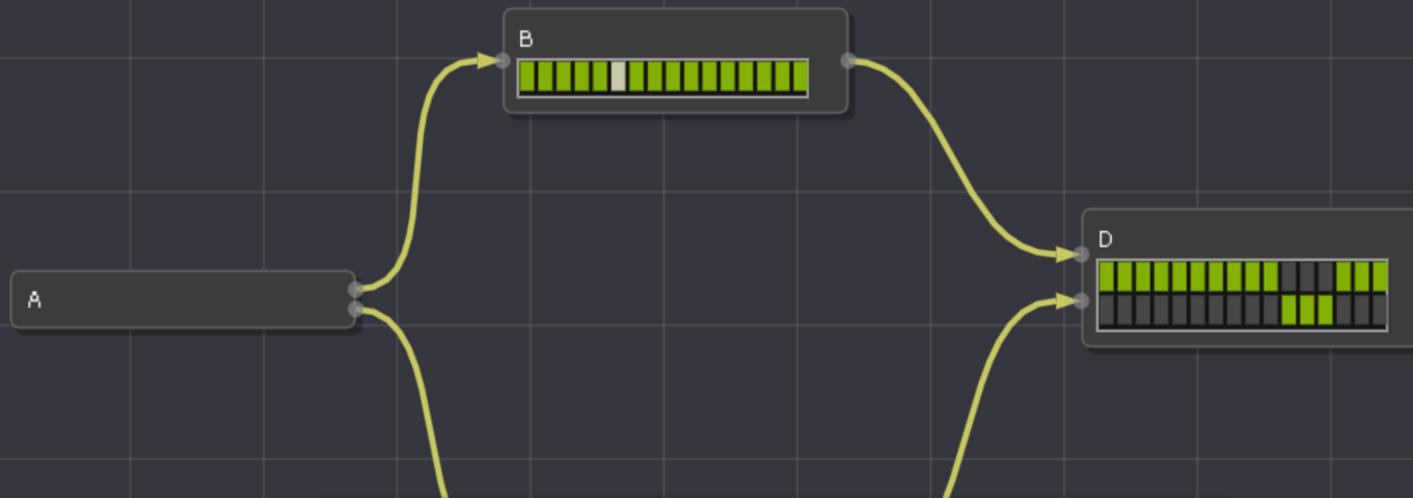
- ▶ A(64674)
- ▶ B(64675)
- ▶ C(64676)
- ▶ D(64677)

Workflow options:

Show grid
 Center
 Hide tree
 Hide metrics
 Hide inspector

Devices

A
B
C
D



An embedded metrics viewer provides in GUI feedback on DPL & user defined metrics. Multiple backends supported, including of course InfluxDB (i.e. for ALICE data taking) and Monalisa (Grid deployments). See "Towards the integrated ALICE Online-Offline (O2) monitoring subsystem", by Adam Wegrzynek

Device Inspector

Channels

channels: 2

Inputs:

Name	Port
from_A_to_C	22001

Outputs:

Name	Port
from_C_to_D	22003

Driver information

Numer of running devices: 4

Play
 Pause
 Step

Workflow options:

anInt	1
aFloat	2.000000
aDouble	3.000000
aString	foo
aBool	true

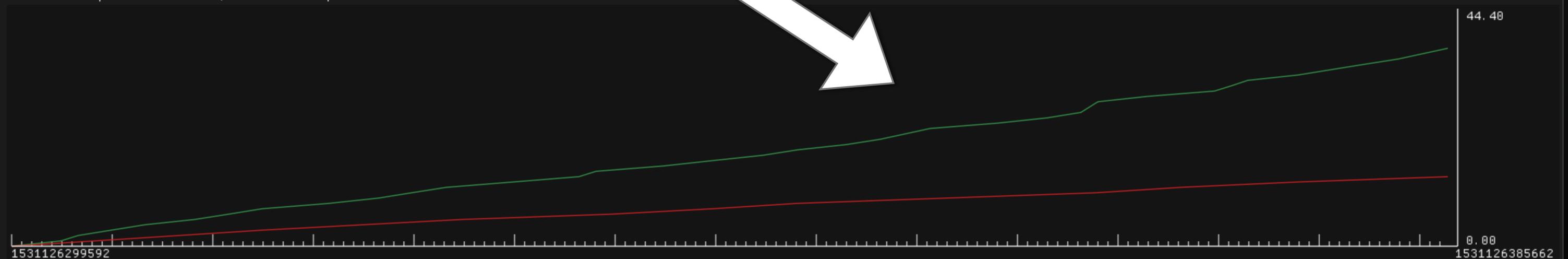
State stack (depth 1)

#0: RUNNING

dpl/stateful_process_count

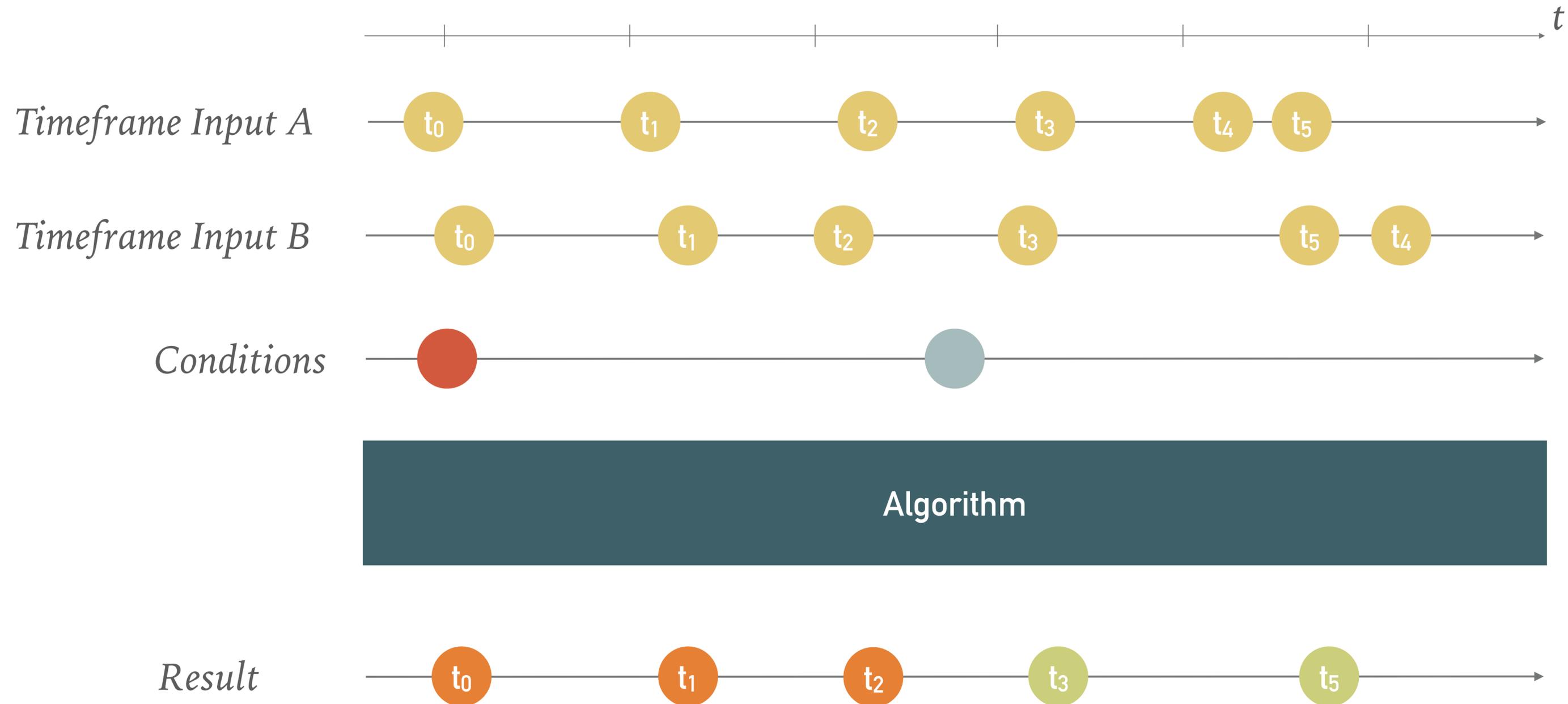
lines

min timestamp: 1531126299592, max timestamp: 1531126385662

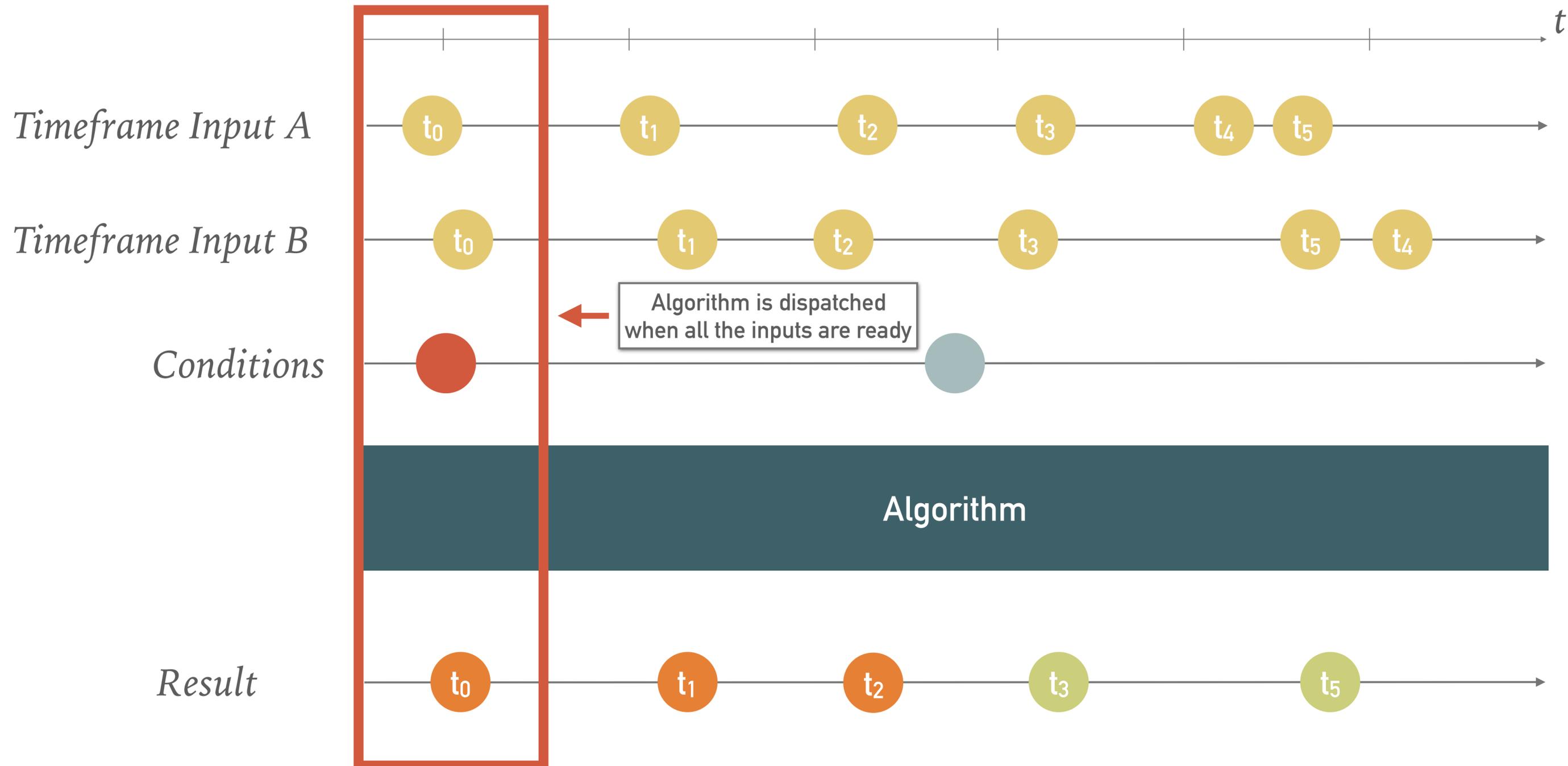


REACTIVE DESIGN

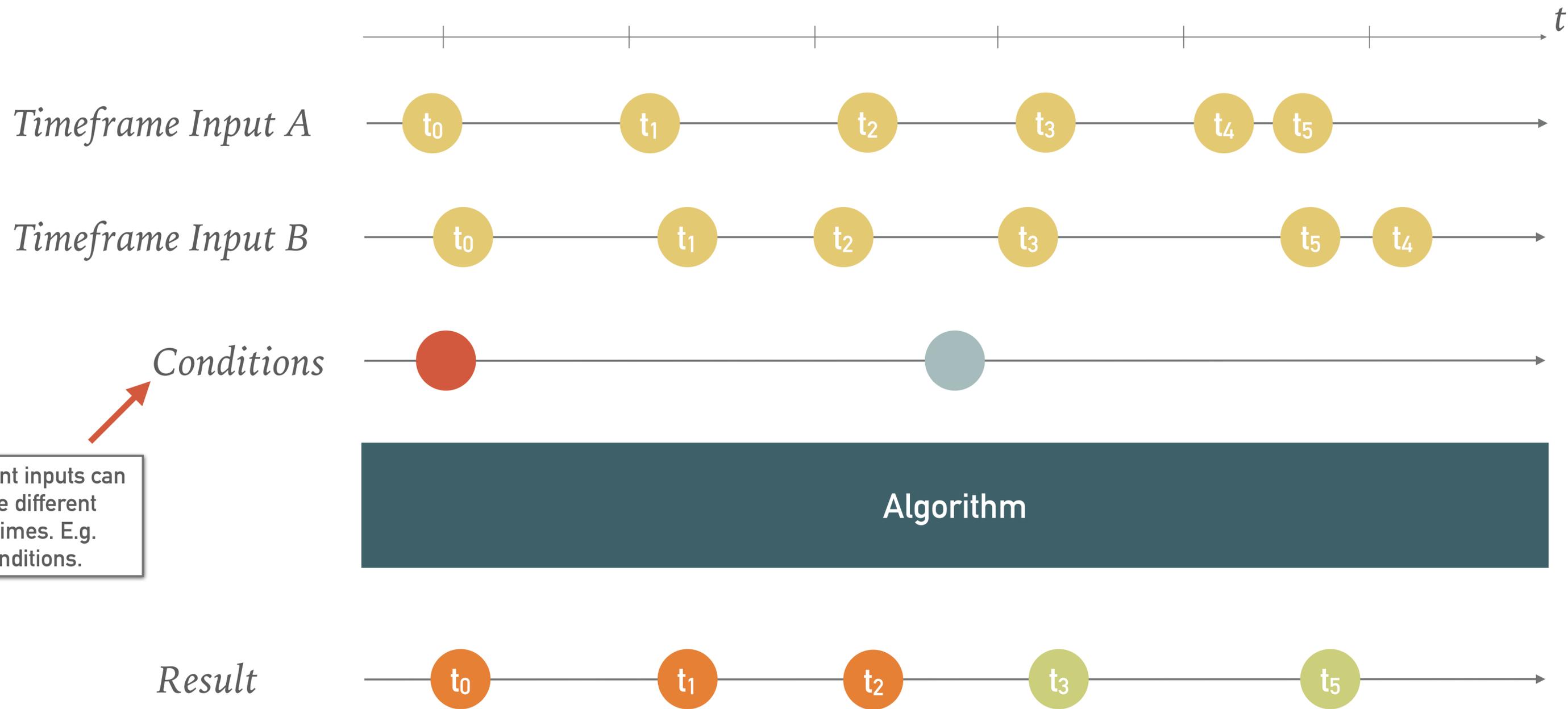
Data is described as pushed through the pipeline.



REACTIVE DESIGN



REACTIVE DESIGN

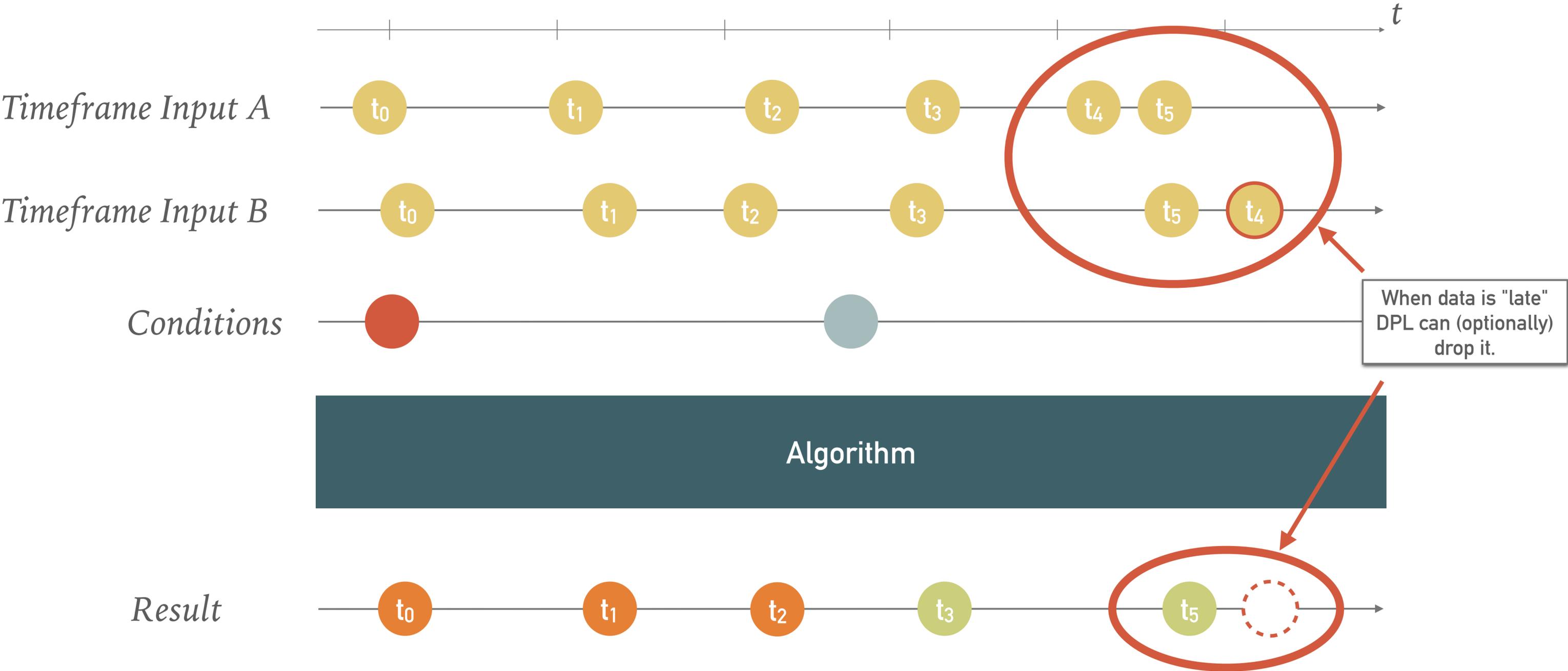


Conditions

Different inputs can have different lifetimes. E.g. conditions.

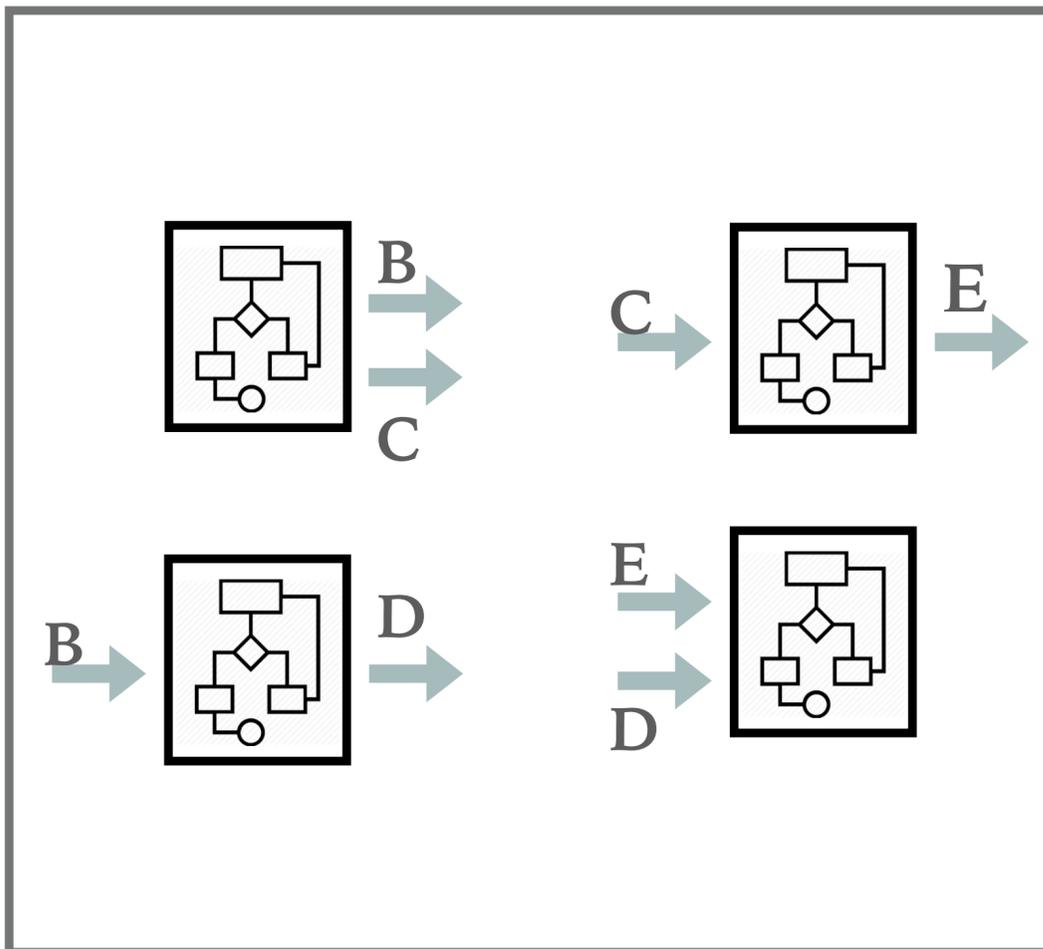
Result

REACTIVE DESIGN



Support for multiple deployment strategies.

Compiles into a single executable for the laptop user.



```
<topology id="o2-dataflow">
  <decltask id="A">
    <exe reachable="true">../bin/o2DiamondWorkflow --id A ...</exe>
  </decltask>
  <decltask id="B">
    <exe reachable="true">../bin/o2DiamondWorkflow --id B ...</exe>
  </decltask>
  <decltask id="C">
    <exe reachable="true">../bin/o2DiamondWorkflow --id C ...</exe>
  </decltask>
  <decltask id="D">
    <exe reachable="true">../bin/o2DiamondWorkflow --id D ...</exe>
  </decltask>
</topology>
```

Generates DDS¹ configuration for deployment on a farm.



Integration with O2 Control system ongoing.²

¹see "DDS – The Dynamic Deployment System" poster by Andrey Lebved

²see "Towards the ALICE Online-Offline (O2) control system" by Teo Mrnjavac

```

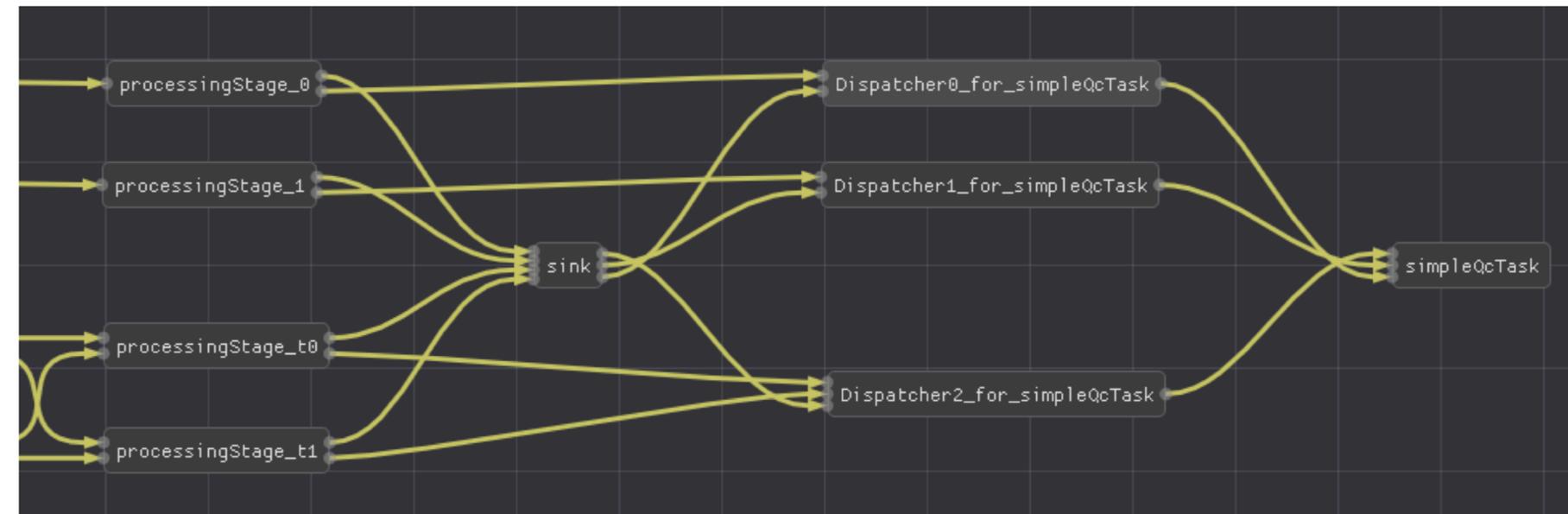
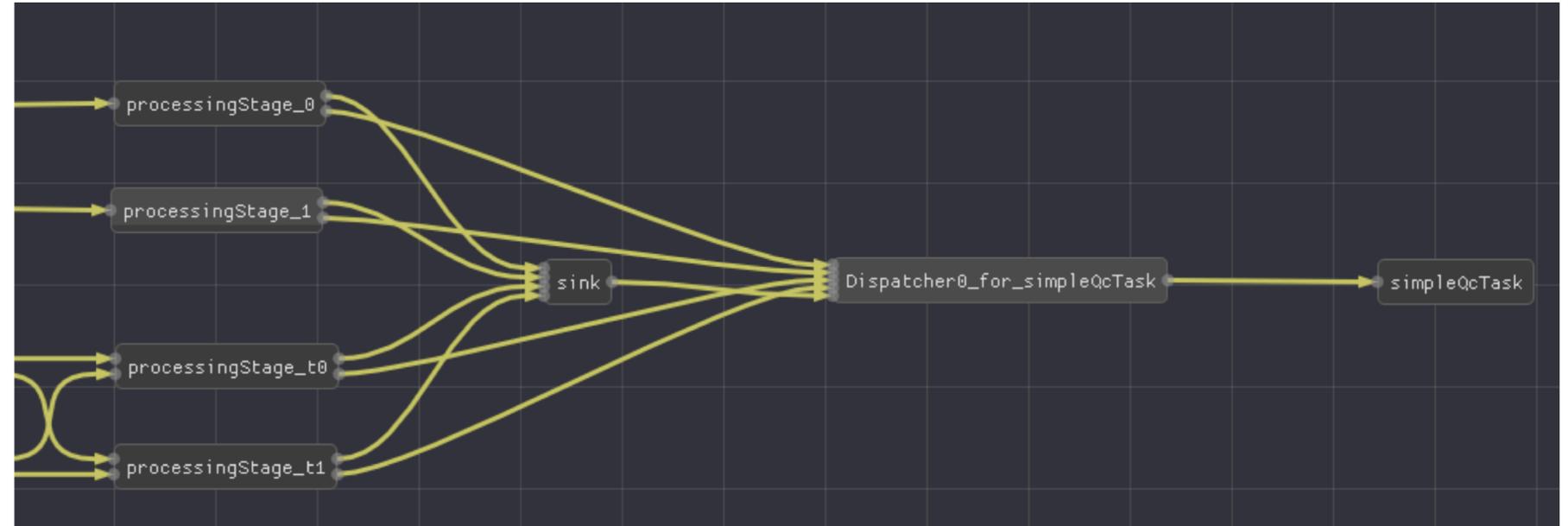
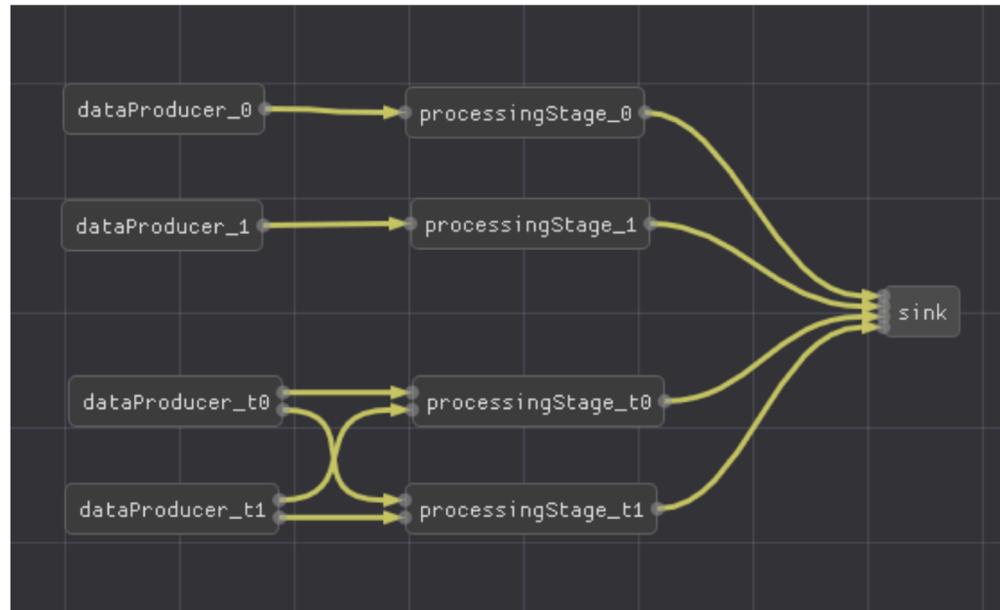
1 #include "Framework/runDataProcessing.h"
2
3 using namespace o2::framework;
4
5 AlgorithmSpec simplePipe(std::string const &what) {
6     return AlgorithmSpec{ [what](ProcessingContext& ctx) {
7         auto bData = ctx.outputs().make<int>(OutputRef{what}, 1);
8     } };
9 }
10
11 WorkflowSpec defineDataProcessing(ConfigContext const&specs) {
12     return WorkflowSpec{
13         {"A", Inputs{}, {OutputSpec{"a1"}, "TST", "A1"}, OutputSpec{"a2"}, "TST", "A2"}},
14         AlgorithmSpec{
15             [](ProcessingContext &ctx) {
16                 auto aData = ctx.outputs().make<int>(OutputRef{ "a1" }, 1);
17                 auto bData = ctx.outputs().make<int>(OutputRef{ "a2" }, 1);
18             }
19         },
20         {"B", {InputSpec{"x", "TST", "A1"}}, {OutputSpec{"b1"}, "TST", "B1"}, simplePipe("b1")},
21         {"C", {InputSpec{"x", "TST", "A2"}}, {OutputSpec{"c1"}, "TST", "C1"}, simplePipe("c1")},
22         {"D", {InputSpec{"b", "TST", "B1"}, InputSpec{"c", "TST", "C1"}}, Outputs{},
23         AlgorithmSpec{[](ProcessingContext &ctx) {}}
24     };
25 }
26 };
27 }

```

The previous example (GUI included) requires 27 user's SLOC.

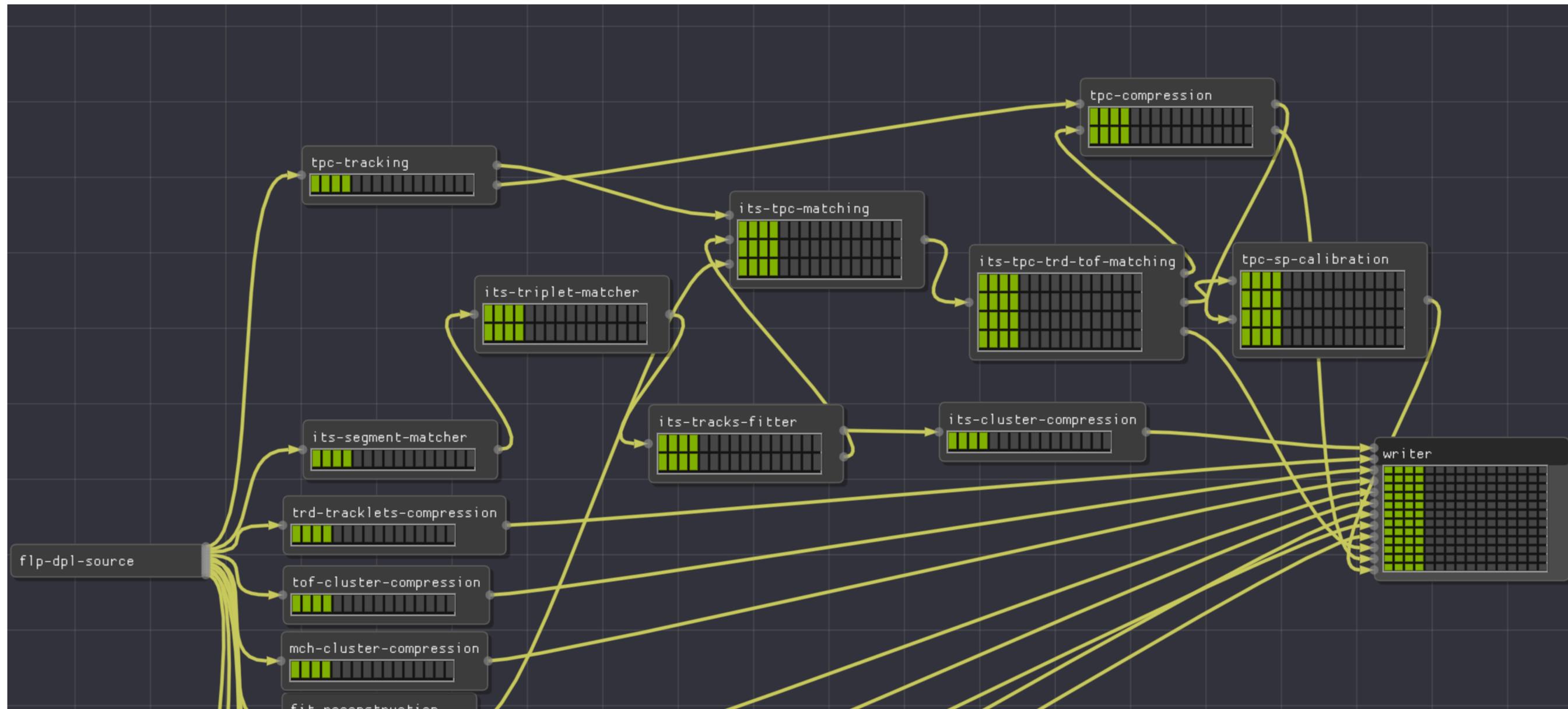
A FEW EXAMPLES

COMPOSABLE WORKFLOWS



Declarative configuration allows for easy customisation: e.g. adding a (one or more) dispatchers for QA.

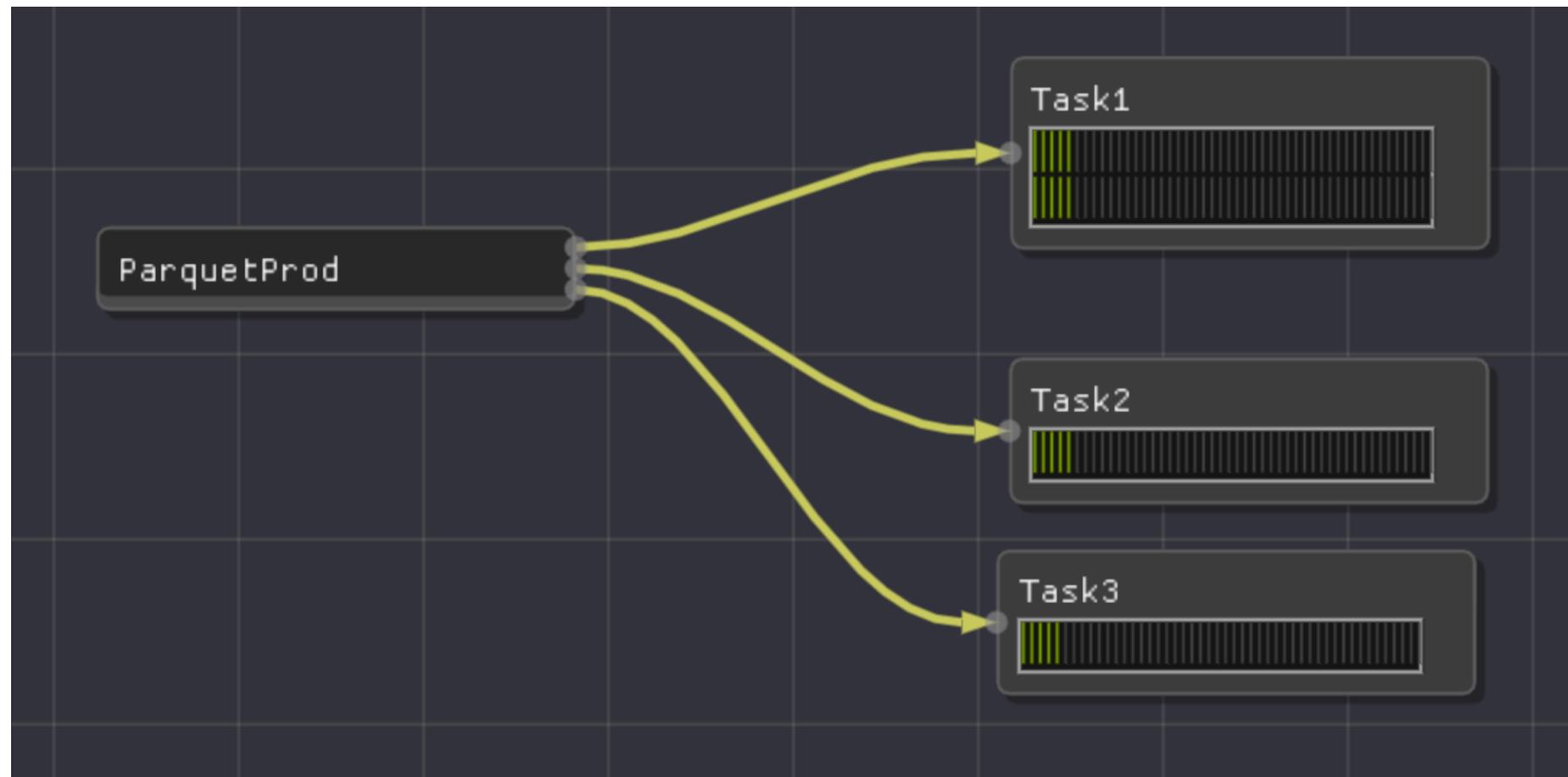
RECONSTRUCTION & GENERAL DATAFLOW



See "Data handling in the ALICE O2 event processing" by Matthias Richter

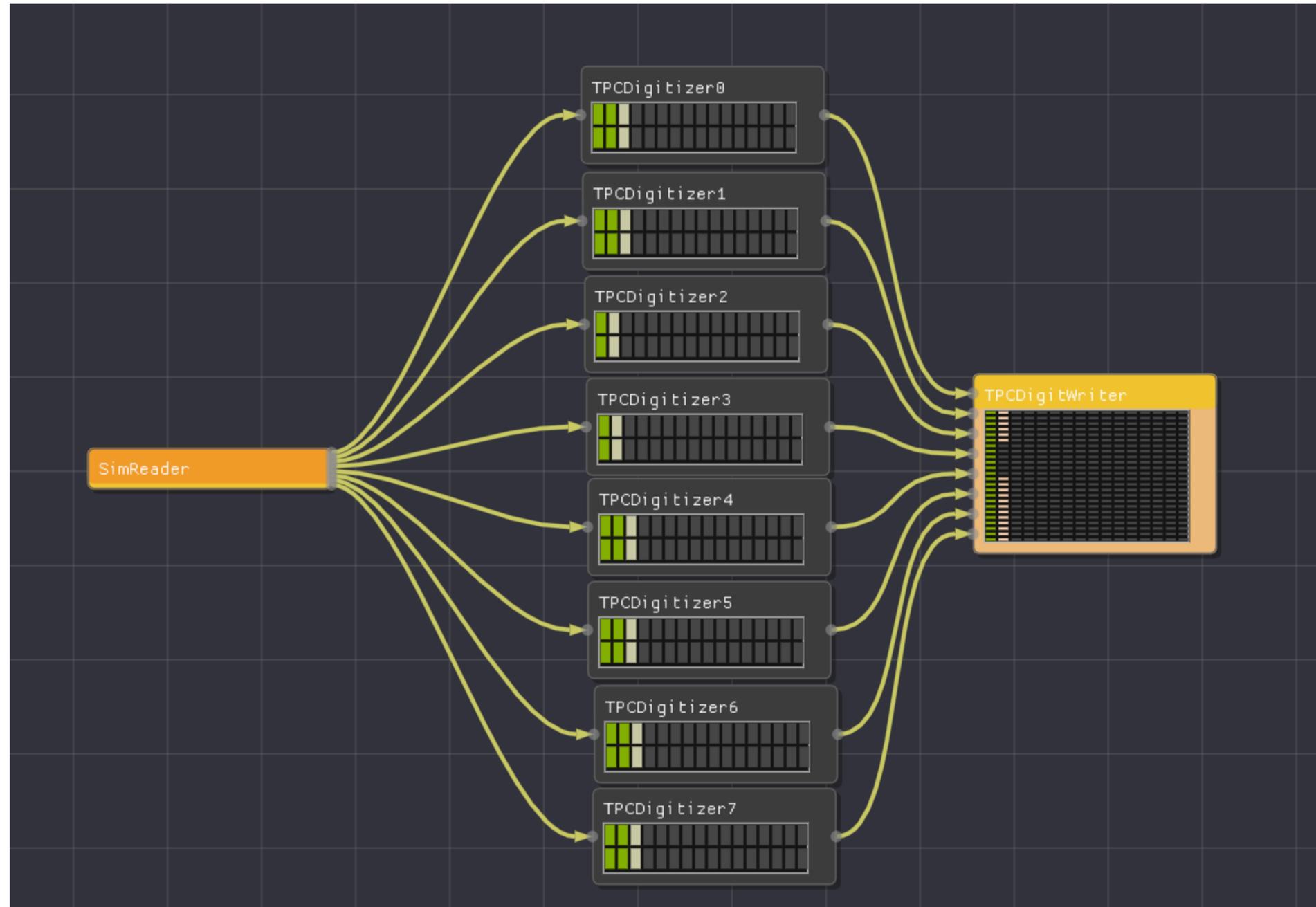
DPL AND ANALYSIS

We are investigating about using the Data Processing Layer also for Analysis.



See ["The ALICE Analysis Framework for LHC Run 3 "](#) by Dario Berzano

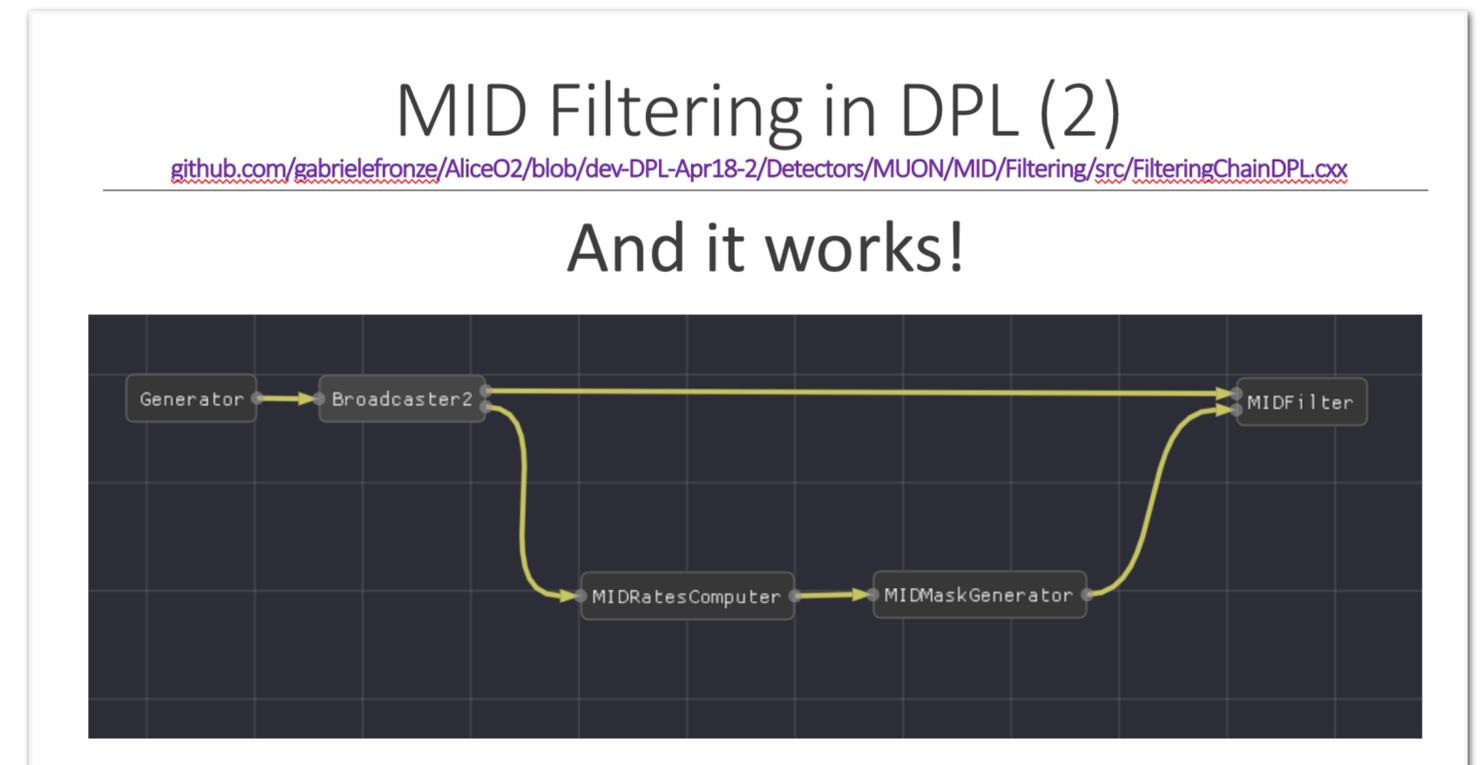
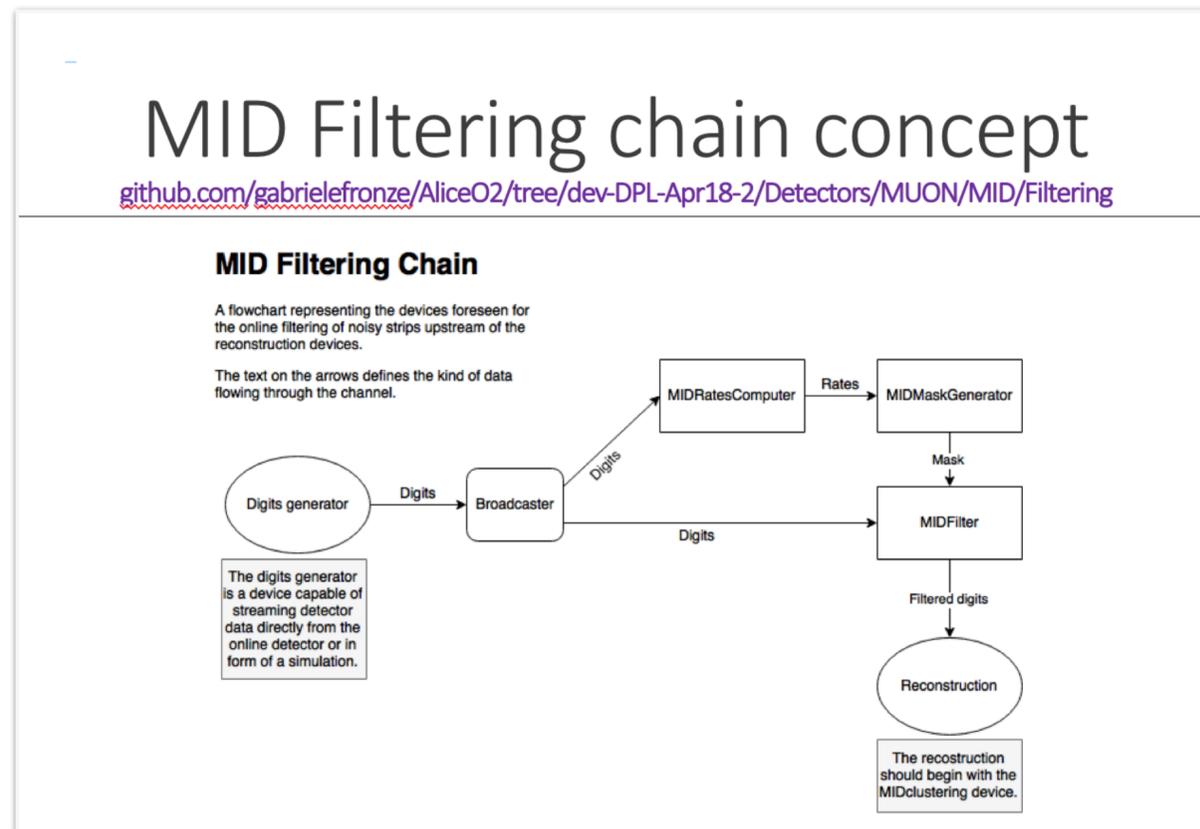
PARALLEL DIGITIZATION



See "A scalable and asynchronous detector simulation system based on ALFA" by Sandro Wenzel

DPL USAGE: MUON IDENTIFIER FILTERING CHAIN

Nice demonstrator by Gabriele Fronzè for Muon Identifier (MID) filtering.



SUMMARY

- The challenges posed by Run 3 imposed to rethink ALICE Computing Architecture, **blending the traditional Online and Offline roles.**
- The **message passing ALFA Framework** is the foundation of ALICE O2 Software Framework.
- We built a message passing / shared memory friendly data model which minimises copy and (de-)serialisation.
- Taking advantage of the O2 Data Model we build a **data flow engine on top of ALFA** to reduce user code and abstract away common hiccups of distributed systems.

BACKUP

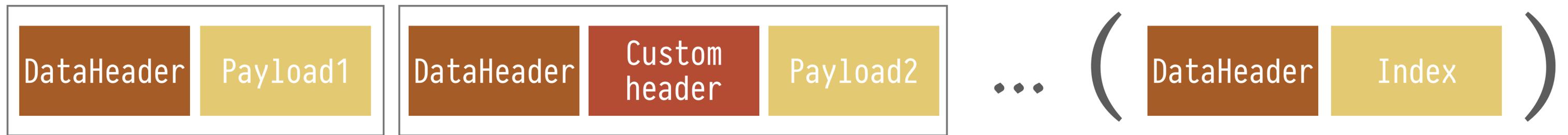
TIMEFRAME

Data quantum will not be the event, but the "Timeframe".

- *~23ms worth of data taking in continuous readout. Equivalent to 1000 collisions. Atomic unit.*
- *~10GB after timeframe building. Vast majority in TPC clusters.*
- *Compressed to ~2GB after asynchronous reconstruction, thanks to track-model-compression, storing clusters instead of ADC values, tailored fixed point integer format, logarithmic precision, entropy encoding.*
- *50x the number of collisions of RUN2.*
- *All MinBias. We need to (lossly) compress information, not filter it.*

02 DATA MODEL

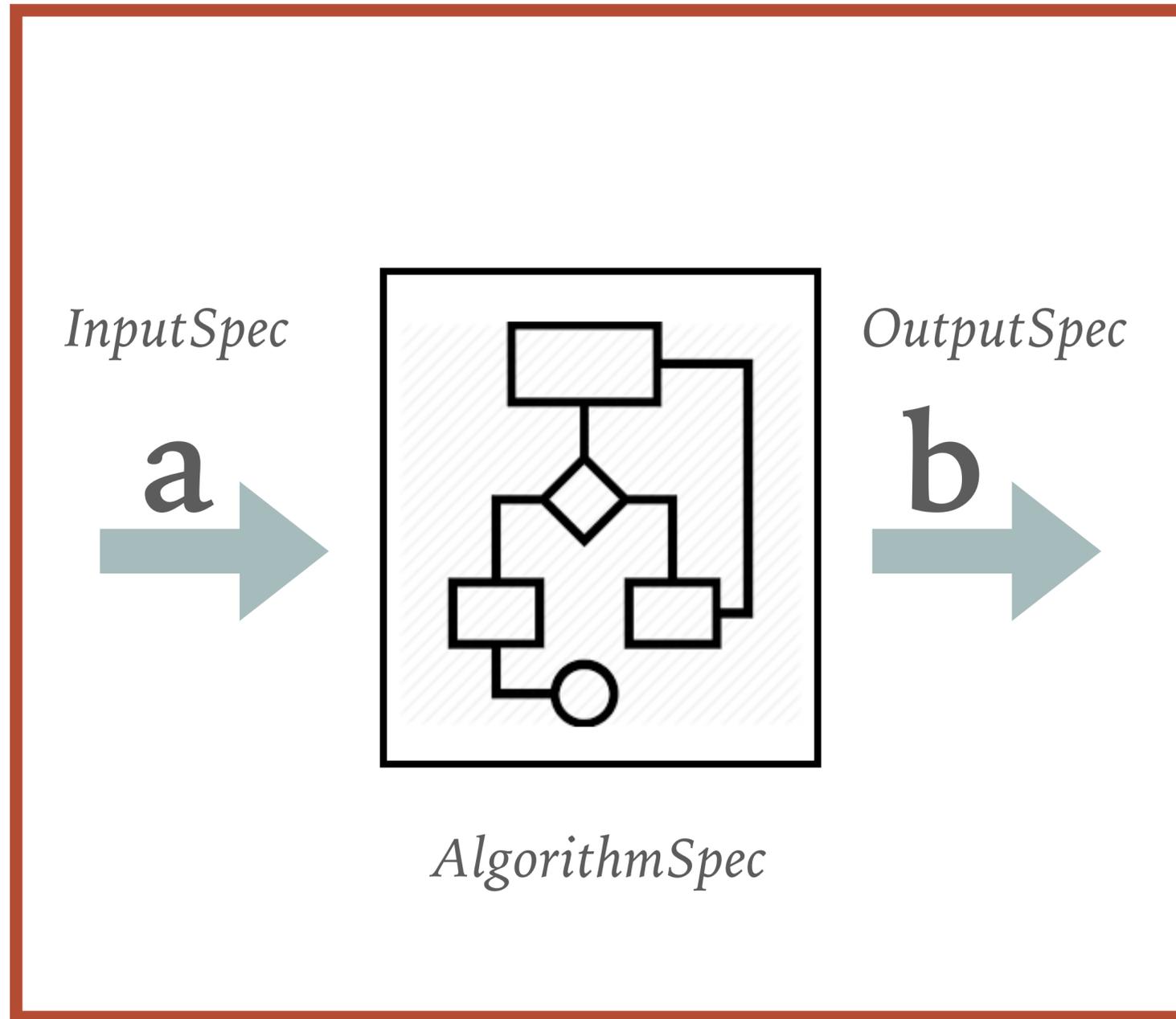
A timeframe is a collection of (header, payload) pairs. Headers defines the type of data. Different header types can be stacked to store extra metadata (mimicking a Type hierarchy structure). Both header and payloads should be usable in a message passing environment.



Different payloads might have different serialisation strategies. E.g.:

- *TPC clusters / tracks: flat POD data with relative indexes, well suitable for GPU processing.*
- *QA histograms: serialised ROOT histograms.*
- *AOD: some columnar data format. Multiple solutions being investigated.*

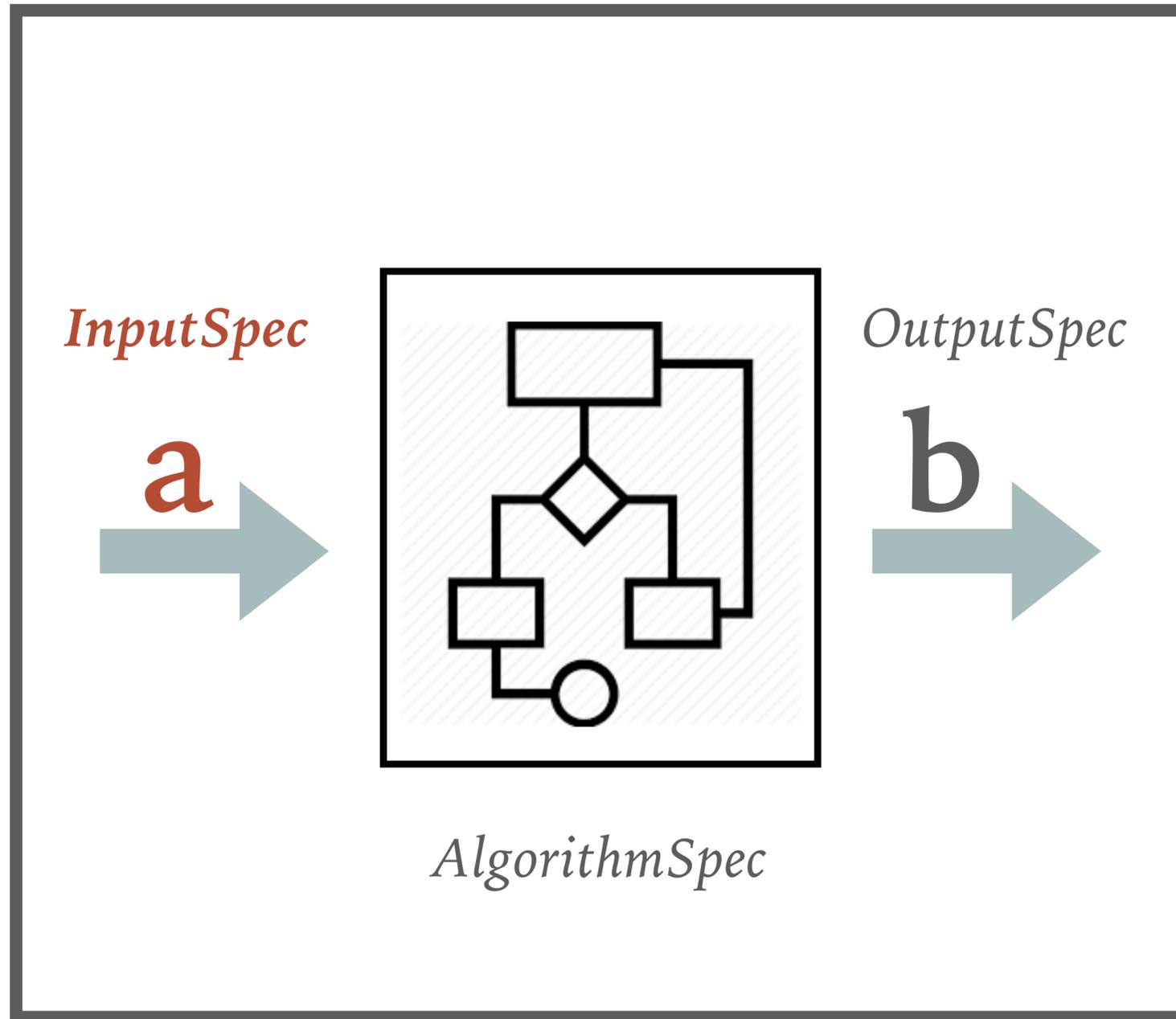
DATA PROCESSING LAYER: HOW



DataProcessorSpec

```
DataProcessorSpec{
  "A",
  Inputs{
    InputSpec{"a", "TPC", "CLUSTERS"}
  },
  Outputs{
    OutputSpec{"b", "TPC", "TRACKS"}
  },
  AlgorithmSpec{
    [](ProcessingContext &ctx) {
      auto track = ctx.outputs().make<Track>(OutputRef{ "b" }, 1);
    }
  }
}
```

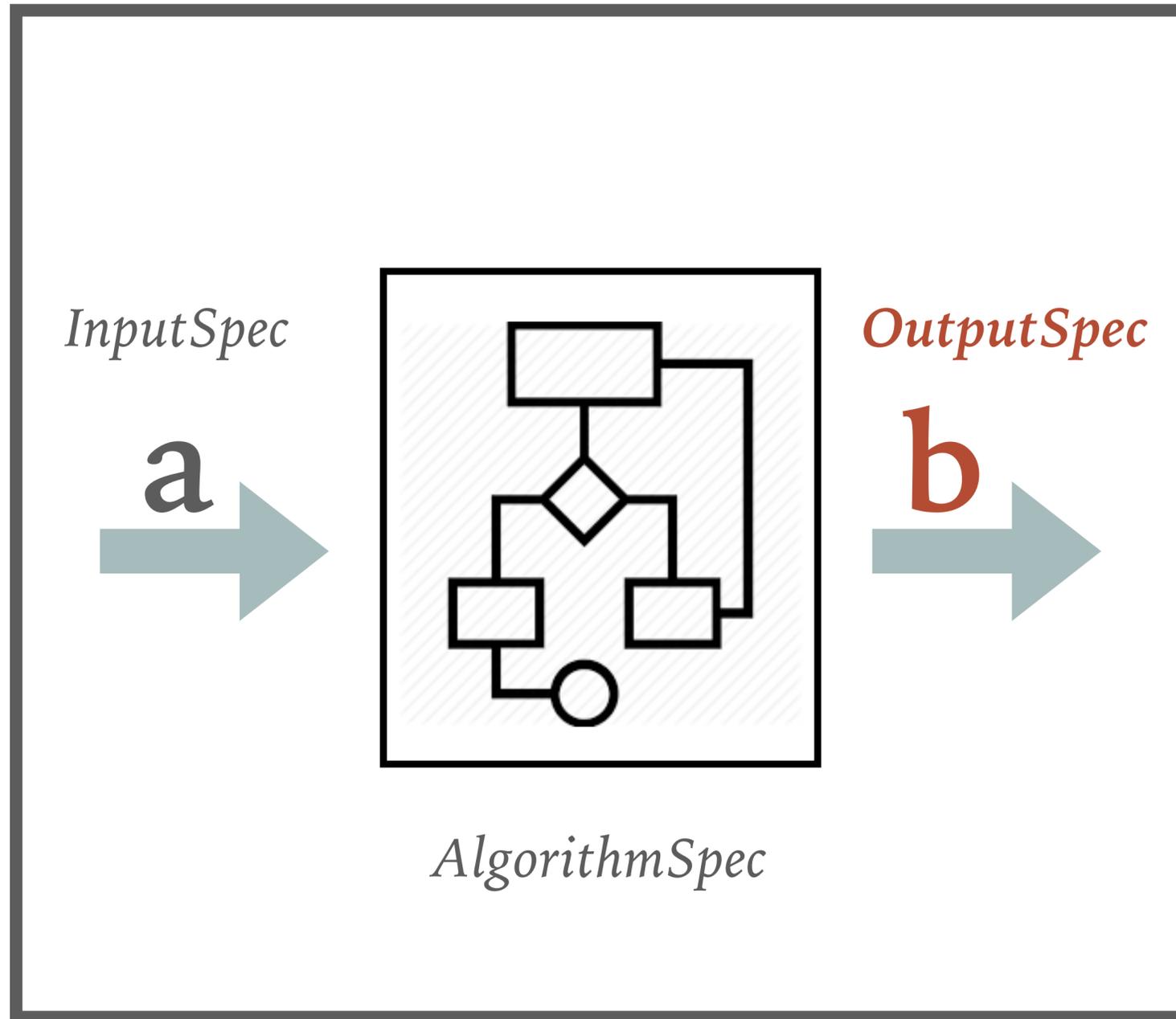
DATA PROCESSING LAYER: HOW



DataProcessorSpec

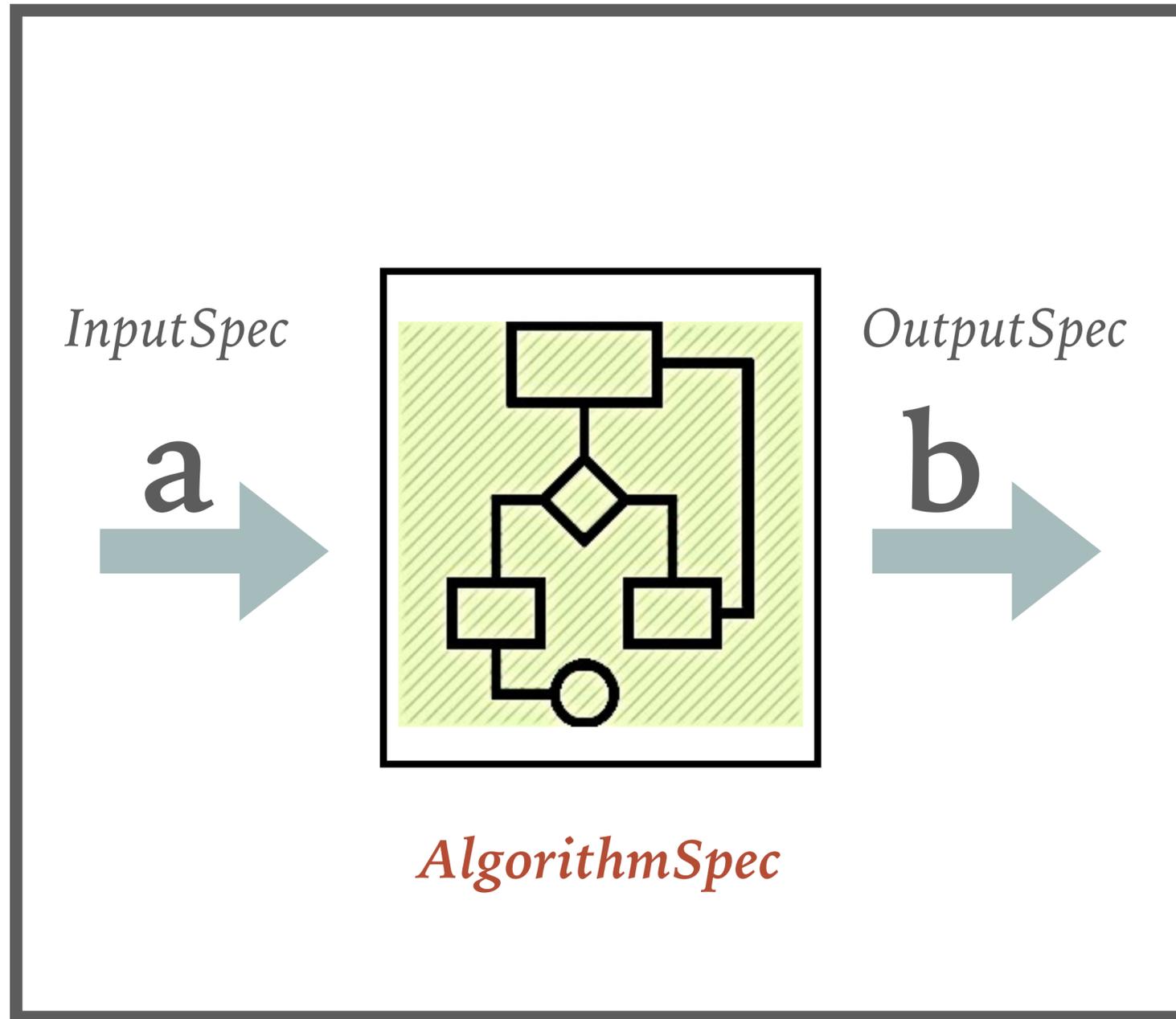
```
DataProcessorSpec{
  "A",
  Inputs{
    InputSpec{"a", "TPC", "CLUSTERS"}
  },
  Outputs{
    OutputSpec{"b", "TPC", "TRACKS"}
  },
  AlgorithmSpec{
    [](ProcessingContext &ctx) {
      auto track = ctx.outputs().make<Track>(OutputRef{ "b" }, 1);
    }
  }
}
```

DATA PROCESSING LAYER: HOW



```
DataProcessorSpec{
  "A",
  Inputs{
    InputSpec{"a", "TPC", "CLUSTERS"}
  },
  Outputs{
    OutputSpec{"b", "TPC", "TRACKS"}
  },
  AlgorithmSpec{
    [](ProcessingContext &ctx) {
      auto track = ctx.outputs().make<Track>(OutputRef{ "b" }, 1);
    }
  }
}
```

DATA PROCESSING LAYER: HOW



```
DataProcessorSpec{
  "A",
  Inputs{
    InputSpec{"a", "TPC", "CLUSTERS"}
  },
  Outputs{
    OutputSpec{"b", "TPC", "TRACKS"}
  },
  AlgorithmSpec{
    [](ProcessingContext &ctx) {
      auto track = ctx.outputs().make<Track>(OutputRef{ "b" }, 1);
    }
  }
}
```

HOW DO YOU LIMIT CONTEXT SWITCH COSTS?

We will have a number of **running** processes which is \approx the number of cores.

Our tasks **take long on a CPU scale** (seconds) thanks to the fact we treat one timeframe at the time (~ 1000 collisions). User code runs lock free.

By describing our computation in terms of **composable pipeline stages** we keep door open for (eventually dynamic) NxM mapping between data processors and actual processes.

We are **willing to pay an extra price** for the sake of:

- *Ease of deployment (microservices!)*
- *Crash resilience (data taking!)*
- *Ability to distribute over multiple nodes (HPC!)*
- *Flexibility (run GEANT3 + GEANT4 + FLUKA!)*

Limiting factor is in any case the GPU for TPC tracking (at least for the synchronous phase).

