# Schedule your micro-services on Docker Swarm with a Sirius-based workflow designer

Olivier Barais, Benoit Combemale,
Cédric Brun, Johann Bourcier,
David Bromberg

INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

Heterogeneous Modeling

# The Eclipse Project

**Gemoc Studio**

**eclipse**

http://eclipse.org/gemoc

http://gemoc.org/studio

## Language Workbench

**Design and compose your executable DSMLs**

## Modeling Workbench

**Edit and debug your heterogeneous models**

UMR IRISA

INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

# Talk outline

1. Motivations

2. Approach overview

3. Quick introduction to Docker

4. A Sirius-based workflow designer

5. Runtime implementation overview

6. Docker and HPC

7. Lesson learnt and open questions

# Your own open-source and low-cost Netflix solution using micro-services

# For teaching SLE/MDE, Distributed computing, Operating System and Chaos Engineering



*"Chaos Engineering is the discipline of experimenting on a distributed system in order to build confidence in the system's capability to withstand turbulent conditions in production."*

INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

# Motivations

- Securely manage high-volume live and on demand video processing/encoding solutions in combination with the scale

- Elasticity of your own/private cloud (edge/fog computing, https://arrayofthings.github.io/)

- Automatically provision and dynamically scale worker instances, and can seamlessly integrate those resources with on site infrastructure to instantly expand video processing capacity

# Execute and isolate tasks using containers

- batch computing/cluster management tool using Docker as execution/isolation system

# Execute and isolate tasks using containers

- We are like Pirates, pillaging for resources instead of booty!
  - We want to run our jobs. We want to get results
  - And when we find available resources, we need to ensure application and environment compatibility

=> This is where containers can be a perfect fit…

- But as I mentioned, our use-case and needs are different from enterprise!

IRISA

INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

# What is Docker?

- Docker is a 'container technology'
  - Linux-specific
    - can't run Mac OSX, Windows *in* docker containers
    - But *can* run docker containers *on* Mac OSX & Windows
  - Shrink-wrap your software, run it on any Linux platform

- *Not* a virtual machine
  - Similar to virtual machines, but more lightweight
    - Smaller, faster to start, easier to maintain and manage
    - Lighter on system resources => vastly more scalable
  - VM-thinking will lead to poor results, avoid it!

IRISA

INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

# Why use it?

- Portability:
  - No need to rebuild your application for a new platform!
    - Build a container once, run it anywhere
      - AWS/GCP/…
    - Stable s/w versions across all platforms, no runtime glitches
  - Think of it as 'modules-to-go'
    - Instead of 'module load PQR' you 'docker pull PQR'
    - No waiting for modules to be built/deployed for you!

- Reproducibility:
  - Because your s/w is stable, your pipeline is reproducible
    - Run the exact same binaries again 10 years from now ☺ ☹

# What can you do with it?

- Computational workloads
  - Use applications without having to install them
  - Run your applications anywhere; clouds, HPC centres
  - **Reproducible pipelines**

- Services
  - Web portals/gateways (R/Shiny, Apache, Jupyter…)
  - Continuous build systems (Gitlab…)
  - For prototyping or for production running (databases etc)

# History

- Dotcloud, Inc creates PaaS service

- January 2013, work starts on docker internally

- March 2013, first public release

- Statistics:
    - 44 328 stars on github
    - 13 152 forks
    - 1693 contributors
    - 32 929 Commits

- Massive community interest

- Created by Solomon Hykes (French engineer ;)

- Open source project => Mobby for open innovation

**EVER TRIED.
EVER FAILED.
NO MATTER.
TRY AGAIN.
FAIL AGAIN.
FAIL BETTER.**

*Samuel Beckett (1906-1989)*

# Who uses Docker?

**Companies using Docker**



**And many more…**

# Who uses Docker?

**Docker PAAS Providers**

Microsoft Azure

Google Cloud Platform

amazon web services™

DigitalOcean

dotCloud

tutum

StackDock

Quay.io

**And many more…**

# Who uses Docker?

**As an Infrastructure Tool along side**

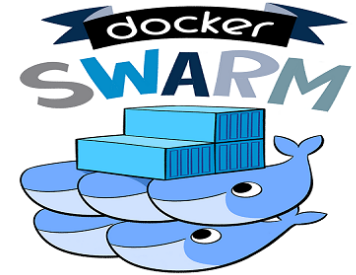INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

# Approach overview

# Architecture Overview

# Focus on docker swarm

- Native clustering for Docker

- Written in GO

- Swarm is a simple tool which controls a cluster of Docker hosts and exposes it as a single "virtual" host

- Swarm uses the standard Docker API as its frontend, which means any tool which speaks Docker can control swarm transparently

But, at this time, swarm mode is focused on long-running services, no real support for batch scheduling
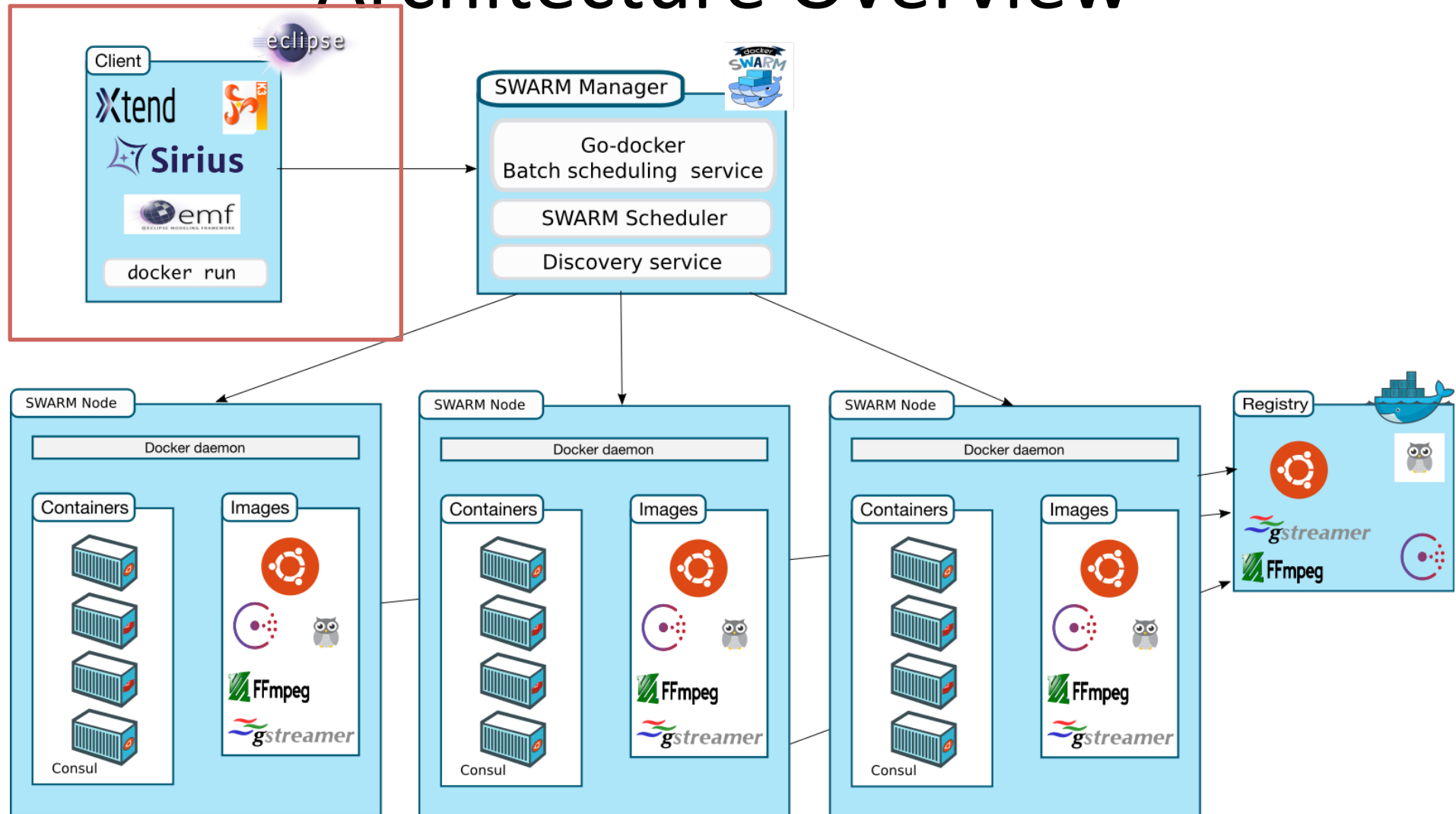
# Focus on Go-docker

- Batch computing/cluster management tool using Docker as execution/isolation system

- Written mostly in Python

- Open source

- like Sun Grid Engine/Torque/Slurm...

- Main contributors (GenOuest BioInformatics Platform)
  - Olivier Sallou [IRISA],
  - Cyril Monjeaud [IRISA]

But, at this time, it is a bit complex to extend for complex scheduling policies based on QoS priorities or resources availability
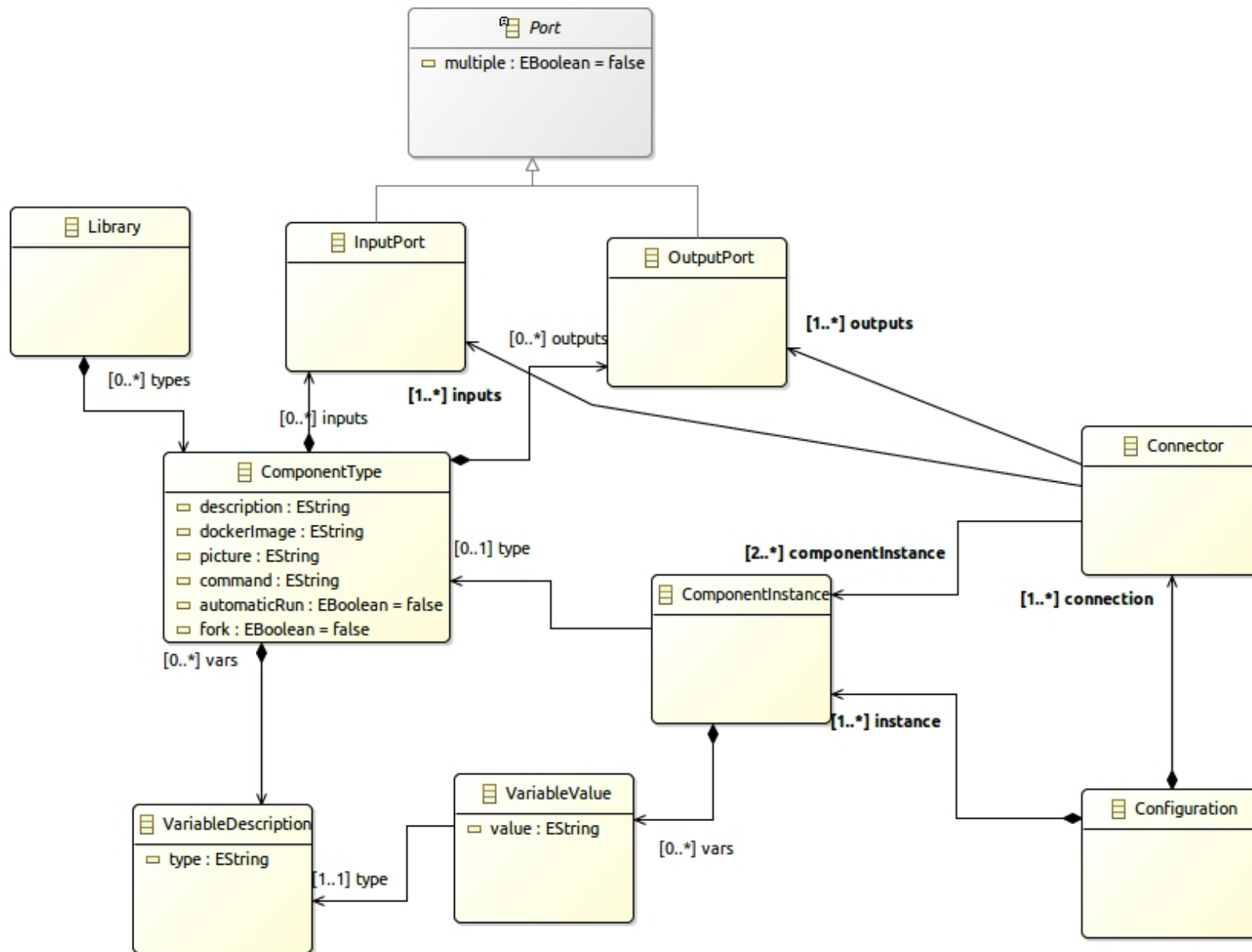
IRISA

Design
# Architecture Overview

# What else is missing?

- Tooling for non-experts
  - Graphical workflow designer to let an end-user defining its own video-processing workflow

    => Eclipse technologies to the rescue

- A language based on flow-based programming paradigm
  - With its meta-model in Ecore
  - With its static semantics in OCL
  - With its operational semantics in Xtend/K3
  - With its graphical concrete syntax in Sirius
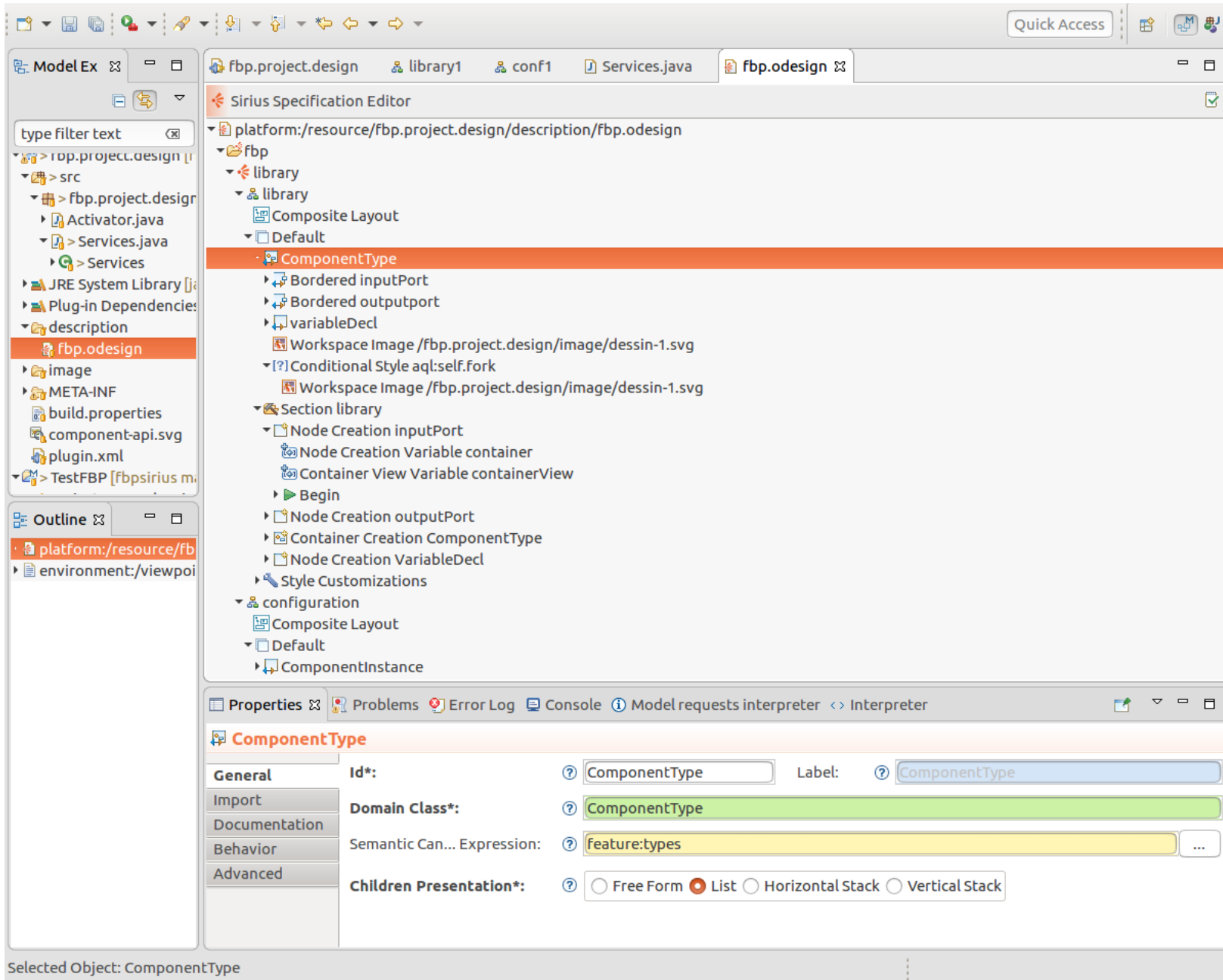  - With its animator with Sirius Animator (Release soon)

A metamodel

INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

Its graphical syntax description

# A flow-based programming editor

# Lesson learnt

- You can use directly Docker
  - You can! It works great for local and private resources. You can use it to develop and share your work with others using Docker-hub

- The basic value proposition is that they help to manage and run applications with complex dependencies easily and efficiently

- Specialize Docker swarm or SwarmKit is not "complex"

IRISA

# Cluster manager complexity

- For go-docker (28 lines of shell scripts)
- For Sirius and EMF stuffs (2 working days to get an initial version)

http://olivier.barais.fr/blog/posts/2015.12.01/GODocker_VideoEncoding.html

http://olivier.barais.fr/blog/posts/2016.03.24/GODocker_on_top_rpi.html

https://github.com/barais/swarm

https://github.com/barais/fbpsirius

IRISA

INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

# Scheduling policy

- Pinning on each raspberry pi
  - One core for managing the GPU encoding
  - One core for cluster management (consul agent, docker swarm agent) and chunk transfer management
  - Two cores for Software video encoding

# Cost

| Device | Number | Unit Price | Total |
|---|---|---|---|
| Rapsberry Pi 2 | 16 | 35$ | 560$ |
| Switch | 1 | 80$ | 80$ |
| Alimentation | 4 | 10$ | 40$ |
| 8Gb SDCARD class 10 | 16 | 5$ | 80$ |
| Ethernet cables | 16$ | 1$ | 16$ |
| Total for a small cluster | | | 776$ |

# Energy consumption

- 16 Raspberry Pis with the switch consume about 80 Watts when used at full power

## VS

- i7 5775C which consumes on average 99 watts when performing a x264 encoding task[1]

1. Taken from the hardware test conducted here

http://techreport.com/review/28751/intel-core-i7-6700k-skylake-processor-reviewed/5

# Performance evaluation

- Comparison of the encoding of a video using a workstation (featuring an Intel(R) Core(TM) i7-5600U CPU @ 2.60GHz, 16 Gb of memory, running on Linux Ubuntu) and our cluster of 16 Raspberry Pi 2

- Encoding a H264 video file into another H264 video file with the "High profile".

- Input video and output video resolution = 1280*688 Px.

# First performance evaluation

- time needed to encode a small video chunk of 2 mins and 30 seconds both on the workstation and on a single raspberry pi 2

| Device | Encoding | Time in second |
| --- | --- | --- |
| Rapsberry Pi 2 | Software | 1601.5 s |
| Rapsberry Pi 2 | Hardware | 554.6 s |
| Workstation (i7) | Software | 126.9 s |

# Second performance evaluation

- time needed to encode a small video chunk of 25 min both on the workstation and the farm of raspberry pi 2

| Device | Time in second |
| --- | --- |
| Farm of Rapsberry Pi 2 | 530,2 s |
| Workstation (i7) | 1281 s |

IRISA

INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

# The good news

- Docker helps to support such requirements
  - I need root!
  - Complex environments
  - Custom distros
  - Bringing your own application+stack
  - Preserving the stack for reproducibility
  - Sharing validated HPC stacks to users
  - Enhancing cluster management and testing

UMR IRISA

INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

# The good news

- HPC community is moving
    - HPCS Singularity
    - NeRSC Shifter
    - Keep compatibility with SLURM for resources allocation
- Jupyter
- R-studio

IRISA

INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

# Wrap up

- Data Intensive computing        often require complex software stacks

- Efficiently supporting "big software" in HPC environments offers many challenges

# Open questions

- Lots of domains such as videos editing requires their own cluster manager and scheduler

- Container scheduler such as apache Mesos or Docker swarm are extensible
  - clean framework for developing and integrating these extensions but Python, Java, or Go are low-level GPL for developing Scheduler

⇔ But why not providing a clean and safe DSL for **designing cluster scheduler ?**

# Discussion/Comments/ Questions

Towards micro-services architecture to transcode videos in the large at low costs - Olivier Barais, Johann Bourcier, David Bromberg, Christophe Dion, In Proceedings of the International conference on Telecommunications and Multimedia (TEMU), 2016

Greening the Video Transcoding Service with Low-Cost Hardware Transcoders - Peng Liu, Jongwon Yoon, Lance Johnson, Suman Banerjee. In proceedings of the 16th USENIX Annual Technical Conference (USENIX ATC '16). June 22–24, 2016, Denver, CO, USA
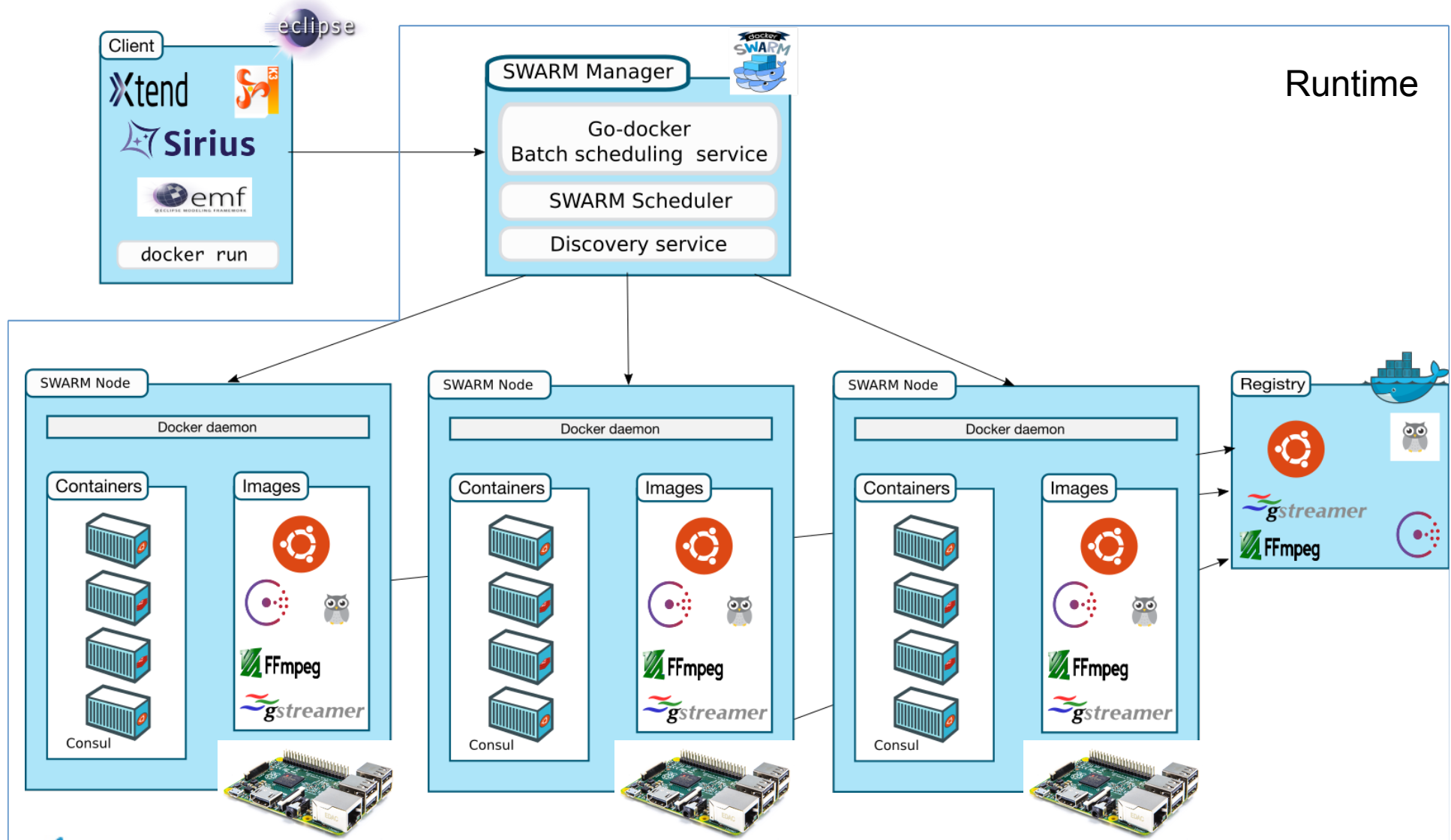
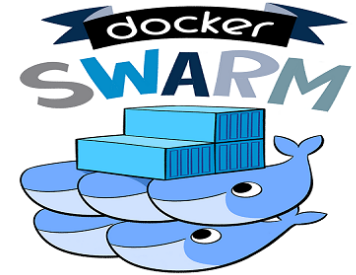# Current work, next steps and open questions

- Compare the performance in fixing the core used by each process

- Implement a version with nomad, mesos (IRT B-COM)

- Compare with the use of Hadoop and Apache Hadoop YARN (IRT B-COM)

# Architecture Overview

# Focus on docker swarm

- Native clustering for Docker

- Written in GO

- Swarm is a simple tool which controls a cluster of Docker hosts and exposes it as a single "virtual" host

- Swarm uses the standard Docker API as its frontend, which means any tool which speaks Docker can control swarm transparently

But, at this time, swarm mode is focused on long-running services, no real support for batch scheduling

# Focus on Go-docker

- Batch computing/cluster management tool using Docker as execution/isolation system

- Written mostly in Python

- Open source

- like Sun Grid Engine/Torque/...

- Main contributors (GenOuest BioInformatics Platform)
  - Olivier Sallou [IRISA],
  - Cyril Monjeaud [IRISA]

But, at this time, it is a bit complex to extend for complex scheduling policies based on QoS priorities or resources availability
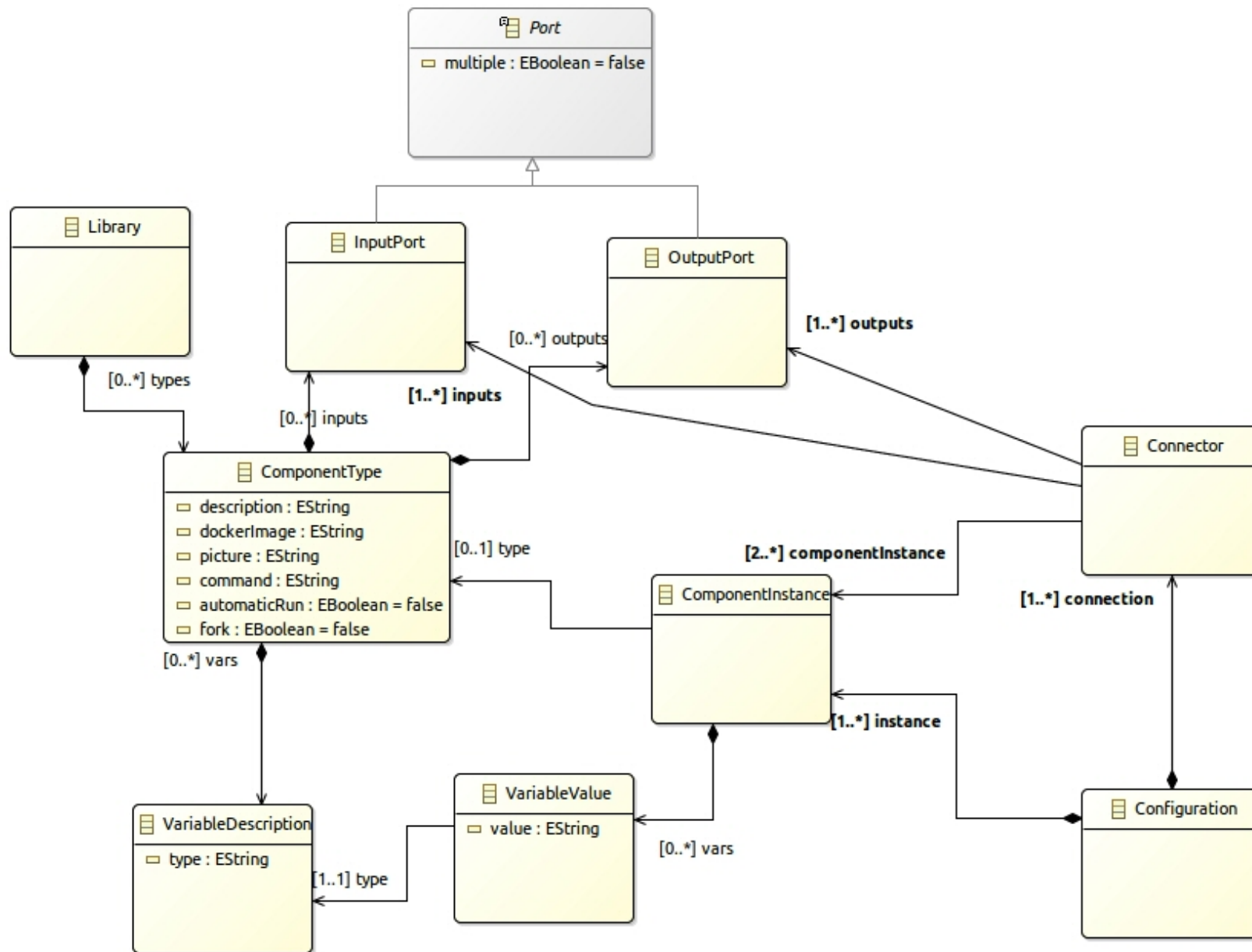
UMR IRISA

# What else is missing?

- Tooling for non-experts
  - Graphical workflow designer to let an end-user defining its own video-processing workflow

    => Eclipse technologies to the rescue

- A language based on flow-based programming paradigm
  - With its meta-model in Ecore
  - With its static semantics in OCL
  - With its operational semantics in Xtend/K3
  - With its graphical concrete syntax in Sirius
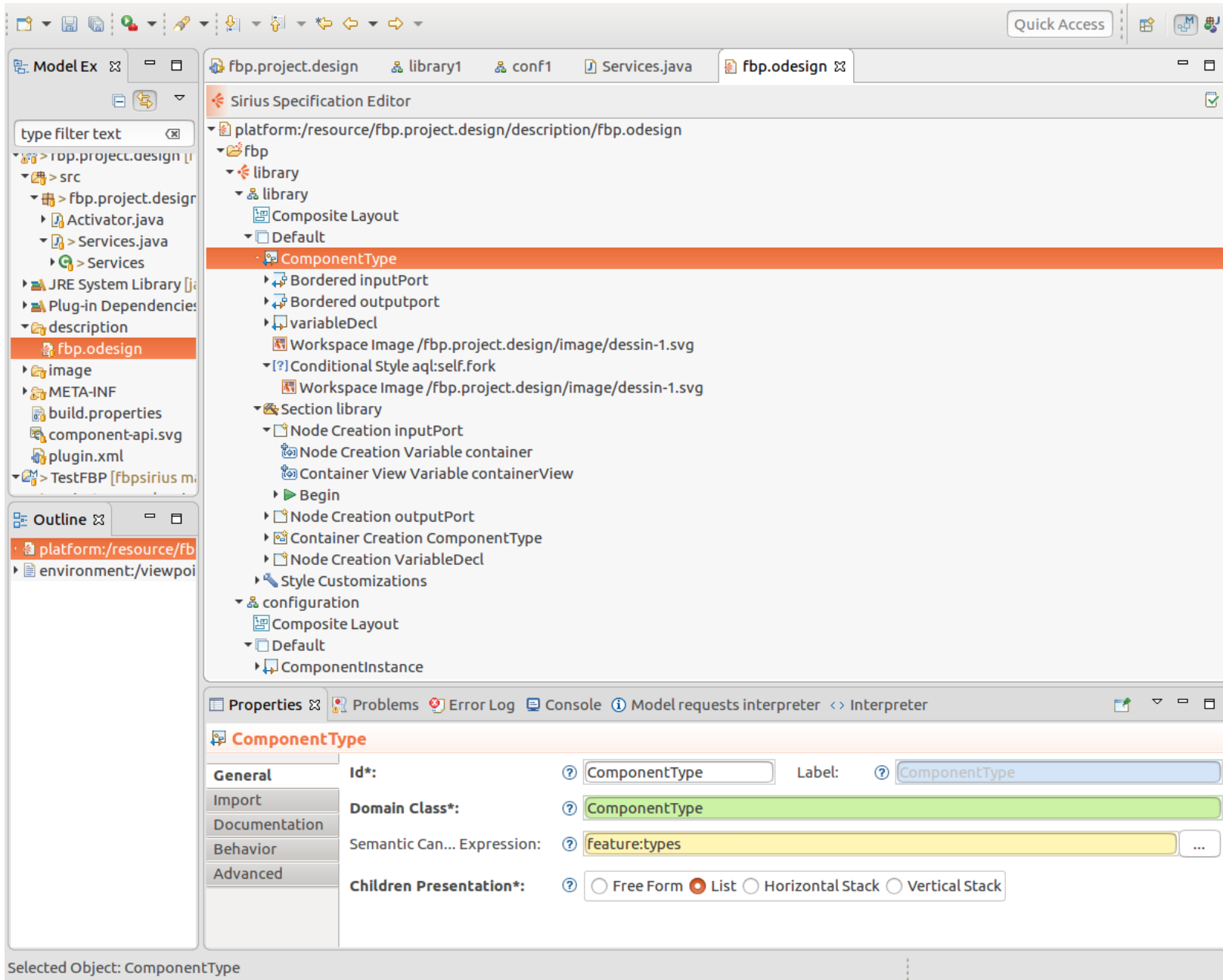  - With its animator with Sirius Animator (Release soon)

A metamodel

Its operational semantics

# A flow-based programming editor

# Demo

http://olivier.barais.fr/blog/posts/2015.12.01/GODocker_VideoEncoding.html
http://olivier.barais.fr/blog/posts/2016.03.24/GODocker_on_top_rpi.html
https://github.com/barais/swarm
https://github.com/barais/fbpsirius

IRISA

INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

# Running on open source hardware and software

- Low-cost open source hardware

- Videos editing open source solutions

- System containers solution

# Implementation          1/4

IRISA

INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

# Implementation 2/4

- A video transcoding workflow:
  - a *splitter* task
  - a *chunk transfer* task to transfer each chunk to a set of targeted host
  - a *scheduler* takes the decision to start the transcoding process on a specific node based on runtime information
  - a *video encoding* task with two different implementations: one for software encoding and another one for hardware encoding

IRISA

# Implementation 3/4

- A video transcoding workflow:
  - an *encoded chunk transfer* to transfer the encoded chunk back

  - a *merger* task which gather all encoded video chunk and assemble them incrementally

  - a *streamer* task which takes the output of the merger task to stream the newly encoded video

IRISA

INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

# Implementation                    4/4

- Splitter and merger: *MKVToolNix* or *FFMpeg*

- Chunk transfer: *udpcast, nfs, glusterfs*

- Videos encoding: *FFMpeg, OpenMAX, Gstreamer*

- Scheduler: Go-Docker or a modified version of swarm

- Key/Value data store. Consul

- Performance analysis: *Cadvisor, Grafana, InfluxDB*

# Lesson learnt

- The development effort required to setup such platform for video transcoding

- The cost of our solution, both regarding financial investment and energy consumption

- The intrinsic encoding performance of our specific deployment setup

INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES

# Current work, next steps and open questions

- Use of Sirius animator to view the status of the running workflow

- Tooling for ensuring the correctness of the resulting videos

- Improve the rate control management between encoding tasks

IRISA