# Scientific Workflow reusing and long term big  data preservation

Salima Benbernou

Université Paris Descartes

Salima.benbernou@parisdescartes.fr

# Outline

- Projet Square Predict
- Scientific workflows
    - Examples
    - Existing systems and limitations
- Using conventionnel workflow technologies in simulation/experiments
    - Introduction
    - Modeling using  BPEL
- Swf and Bog data
- Reusing in scientific workflow
    - Fragment reusing
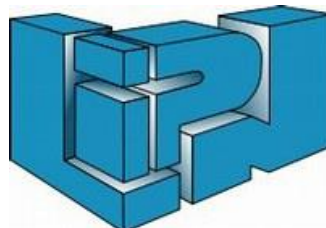    - Privacy aware provenance

# Square Predict

**Prédiction financière:**

**Gestion de données assurances et open data**

# Square Predict

- Collecte de données sur Teralab

- Fusion de données sémantiques (RDFs):

    -Evaluation de requêtes par réecriture distribuée sur Spark

    -Requêtes virtuelles

    -réutilisation des données pour d'autres requêtes

- Qualité de données:inconsistance, l'incertitude

- Prendre les Données en streaming

- Clustering en streaming sur Spark

- Visualisation

# SP3.1 – Architecture fonctionnelle
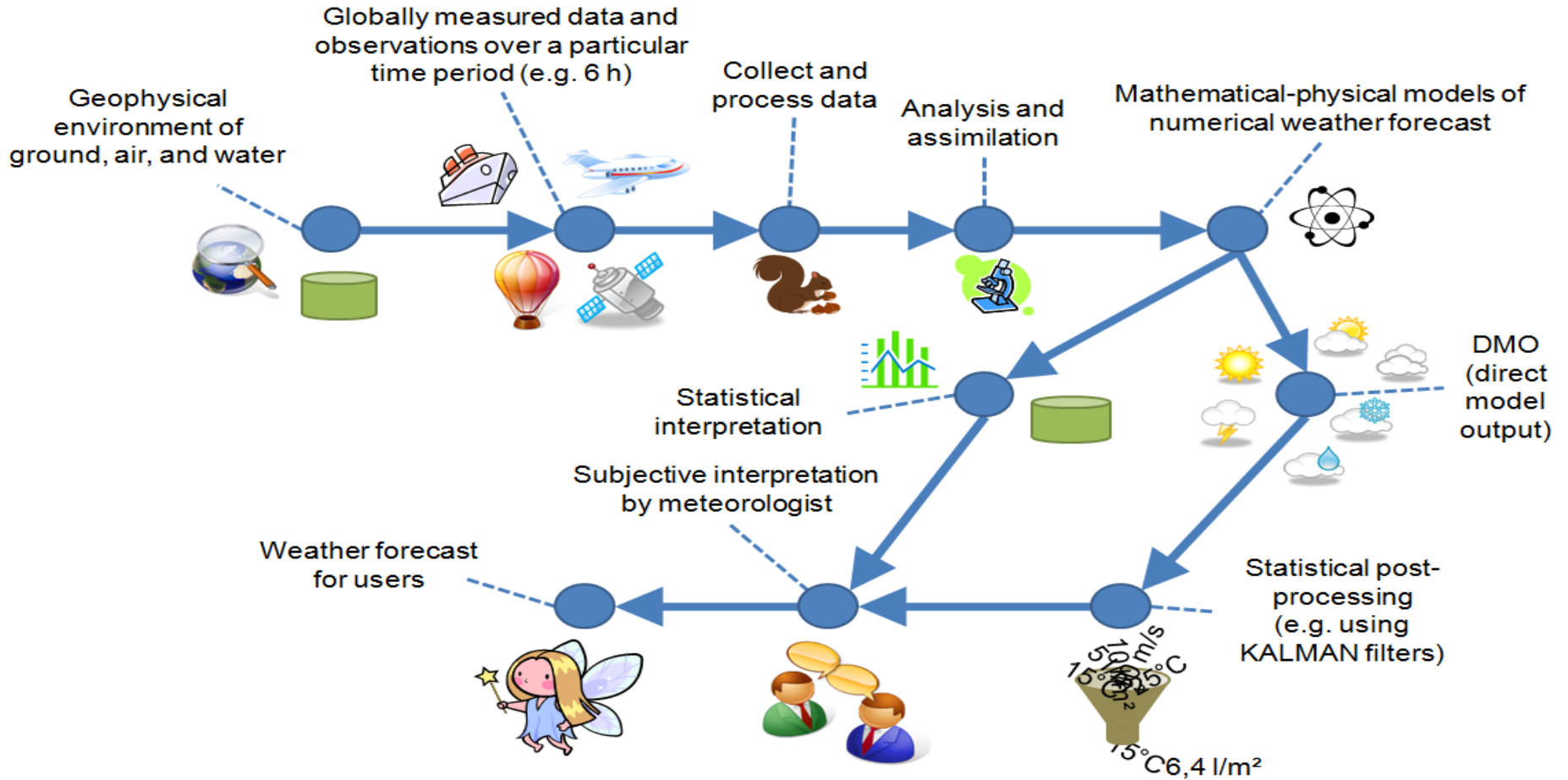
# SP3.1 – Architecture technique

# Scientific workflow

# What are scientific workflows?

- Scientific experiments/computations/simulation modeled and executed as workflows  called scientific workflow (SWf).

- Deal with intensive data, are long running, data driven, can integrate multiple data sources (i.e. sensors)

# SWF Examples



Geophysical environment of ground, air, and water

Globally measured data and observations over a particular time period (e.g. 6 h)

Collect and process data

Analysis and assimilation

Mathematical-physical models of numerical weather forecast

DMO (direct model output)

Statistical interpretation

Subjective interpretation by meteorologist

Weather forecast for users

Statistical post-processing (e.g. using KALMAN filters)

# SWF Examples

Functional MRI Analysis Workflow

# SWF Examples (cont)



ATLAS experiment (simplified)

31

# SWF Examples (cont)



A snapshot of the Taverna Workbench.

# Scientific workflow systems

- Workflow are already used in e-science
- This is not always the conventional workflow technology
- Some workflow systems in e-science: Kepler, Taverna, Pegasus, Trident, Simulink, Karajan …

- To be improved
  - Robustness, fault handling
  - Flexibility and adaptability
  - Reusability
  - Scalability
  - Interaction with users, user-friendliness of tools
  - science skills required from scientist
  - No generic approach
  - Domain specific solutions (in term of modeling and execution)

# Scientific workflow systems

- Data-driven applications are more and more developed in science to exploit the large amount of digital data today available

- Adequate workflow composition mechanisms are needed to support the complex workflow management process including workflow creation, workflow reuse, and modifications made to the workflow over time.

- Use conventional technologies **(Business processes)**

# Business workflows (i.e, BPEL)

- independent of the application domain, can be used for every type of scenario

- The concept of workflow models and instances is inherently capable of enabling parameter sweeps.

- Asynchronous messaging features are predestinated for non-blocking invocation of long running scientific computations

- Business workflows are usually based on agreed-upon standards for workflow modeling and execution as well as for integration technologies.

- facilitate collaboration between scientists (e.g., with the help of Web services).

- Services computing technology enables scientists to expose data and computational resources wrapped as publicly accessible Web services

# Scientific workflows vs. Business workflows

- **Scientific "Workflows"**
  - Dataflow and data transformations
  - Data problems: volume, complexity, heterogeneity
  - Grid-aspects
    - Distributed computation
    - Distributed data
  - User-interactions/WF steering
  - Data, tool, and analysis integration
  - ➔ Dataflow and control-flow are *married!*

- **Business Workflows**
  - Process composition
  - Tasks, documents, etc. undergo modifications (e.g., flight reservation from *reserved* to *ticketed*), but modified WF objects still identifiable throughout
  - Complex control flow, task-oriented: *travel reservations; credit approval*
  - ➔ Dataflow and control-flow are *divorced!*

# Business and scientific Lifecycles

# Scientific workflow limitations

- Scientific workflow life cycle: scientits'perspective
  - Reflects how scientits actually work-trial and error fashion
  - Hidden technical details
  - « no » deployment phase
  - Operations to control workflow execution
  - Monitoring is the visualisation of the results only

# Scientific workflow and the scalability

- Service-Oriented Workflows on Cloud Infrastructures for **reusing**.

- The service-oriented paradigm will allow **large-scale distributed workflows** to be run on heterogeneous platforms and the integration of workflow elements developed by using different programming languages. **Web and Cloud services** are a paradigm that can help to handle workflow interoperability,

# SWF: Programmable and reproducible Scalability

- Access and query data
- Scale computational analysis
- Increase reuse
- Save time, energy and money
- Formalize and standardize

# Workflow, data science and big data

# Workflow, data science and big data (cont)

- many SWfMSs are not prepared to handle large data sets because of inadequate support for distributed computing.

- most SWfMSs that do support distributed computing only allow static task execution orders

# Workflow, data science and big data (cont)

- Develop <span style="color:red">new big data science technologies</span> and infrastructure

- Develop data science workflow application through combination of tools, technologies and best practices

- Hands on consulting on workflow technologies for big data and cloud systems, e.g., MapReduce, Hdoop, Spark, Yann, Cascading.Oosie, Nova .

- Technology briefings and applied classes on end-to-end support for data science.

# Some challenges

- Optimized execution on heterogeneous platforms.

- Representing and reasoning data: semantic, quality (unconsistent, uncertain)

- Increasing reuse within and across application domains

- Querying and integration of workflow provenance data

# BPaas vs SPaaS

- Business process vs swf outsourcing to take advantage of the Cloud computing model.

- Reusing process fragments to develop process-based service compositions and adapt the new swf according to the scientists  (reusing a partial differential equation program)

-  privacy risks aware.

- **Sharing scientific process fragments and hide provenance.**

- Many works studied data provenance but not hide provenance.

- Formal model of BPaaS vs Scientific Process as a Service (Icloud@vldb2012)

# Decomposition of business process vs scientific process



1. A business process is decomposed into a set of process fragments suitable for re-use.

**Hospital Business Process**

# Identification of fragments



1. A business process is decomposed into a set of process fragments suitable for re-use.

2. Each process fragment is identified and made available to be reused as a Web Service.

**Hospital Business Process**

**Insurance Business Process**

# Development of processes



3. A process-based service compositions is developed by reusing process fragments deployed in the Cloud.

**Hospital Business Process**

**Employer Business Process**

**Insurance Business Process**

# Provenance and privacy in SPaaS

An adversary (a curious) **can discover the provenance** of the reused process fragments.

**Can infer connections** between end-users and scientists that outsource fragments to the Cloud.

✓**No related work!**

# Formal model



Start ◎   End ○   Data element name ●   Activity name ⊕   Data attribute ○   Data value ●

- **Business Process:** business graph. [Beeri,VLDB'06]

- ✓ **Process Fragment:** business subgraph.

- ✓ **BPaaS:** a finite set of business processes.

- ✓ **Reusing Function.**

# Anonymous Views



- **View on SPaaS:**

  A set of process fragments having the same objective (called clones).

- **Anonymous View on SPaaS:**

  View on SPaaS having at most K clones.

- **Objective : Make it hard for an adversary to know the provenance of a reused process fragment. (**Anonyfrag)

```
┌─────────────────────────┐
│  Let P a business process │
│  to be developed in the   │
│         SPaaS             │
└─────────────────────────┘
              │
              ▼
┌─────────────────────────┐
│    For each Process       │ ◄──────────────────────────┐
│    Fragment F in P        │                             │
└─────────────────────────┘                             │
     │           │                                        │
     │           └──────────┐                             │
     │                      ▼                             │
     │            ┌──────────────────┐    ┌◇◇◇◇◇◇◇◇◇◇┐   │
     │            │ Verify if F exists│───▶│ not exists│──┘
     │            │   in the SPaaS    │    │           │
     │            └──────────────────┘    │   exists  │──┐
     │                                     └◇◇◇◇◇◇◇◇◇◇┘  │
     ▼                                                    ▼
┌─────────────────────────┐                  ┌─────────────────┐
│   P' a process-based      │                  │   Reuse F to     │
│  service compositions     │                  │   develop P      │
│  developed in the SPaaS   │                  └─────────────────┘
└─────────────────────────┘
```
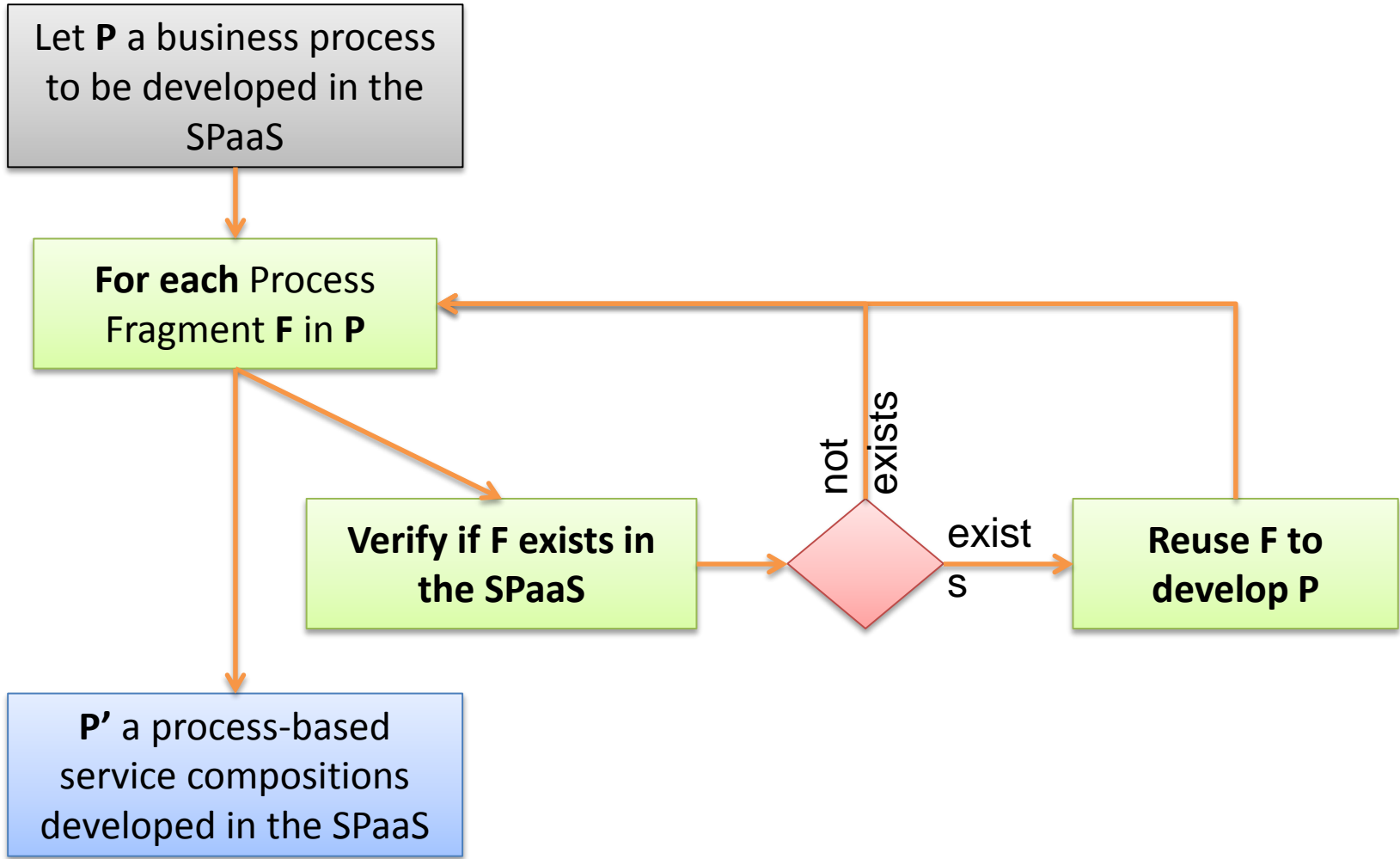
Let **P** a business process to be developed in the SPaaS

**For each** Process Fragment **F** in **P**

**Verify if F exists in the SPaaS**

not exists

exists

**Reuse F to develop P**

**P'** a process-based service compositions developed in the SPaaS

# Workshop organisation

1st Workshop on LOng term Preservation for big Scientific data (LOPS) to be held in conjunction with ICDE 2014, Mach 31-April 4, Chicago , IL, USA

Lipade.math-info.univ-paris5.fr/lops/