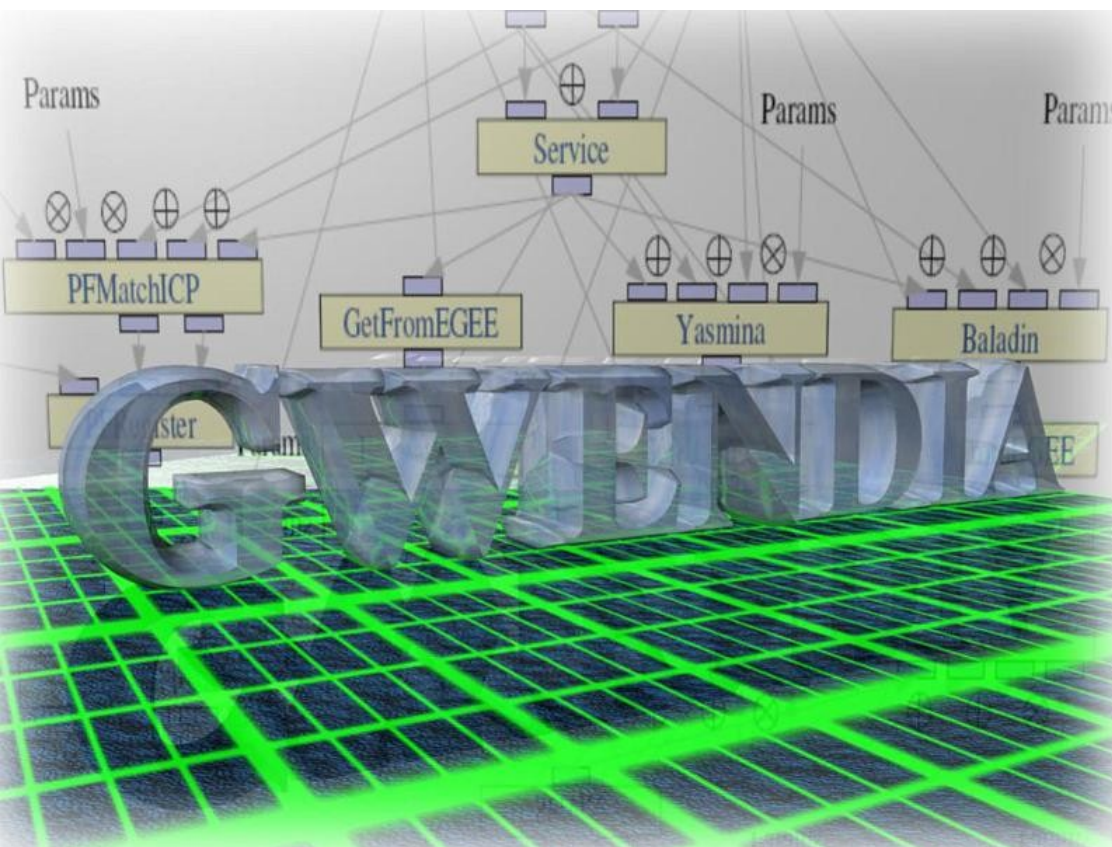




Grid Workflow Efficient Enactment for Data Intensive Applications

Software composition for scientific workflows



Johan Montagnat
CNRS / UNS, I3S, MODALIS
April 17, 2009

modalis



Financé par
ANR

- **Scientific production**

- Europe: EGEE (www.eu-egee.org), USA: OSG (www.opensciencegrid.org), NAREGI (www.naregi.org)...



EGEE European infrastructure:

- > 250 computing centers in 45 countries
- > 80 000 CPU cores
- > 13 000 users
- > 250 000 jobs / day

- **behavior**

- Very large scale (WAN)
- Coarse grain parallelism, loose coupling
- High heterogeneity, low reliability, large latencies



MR Images

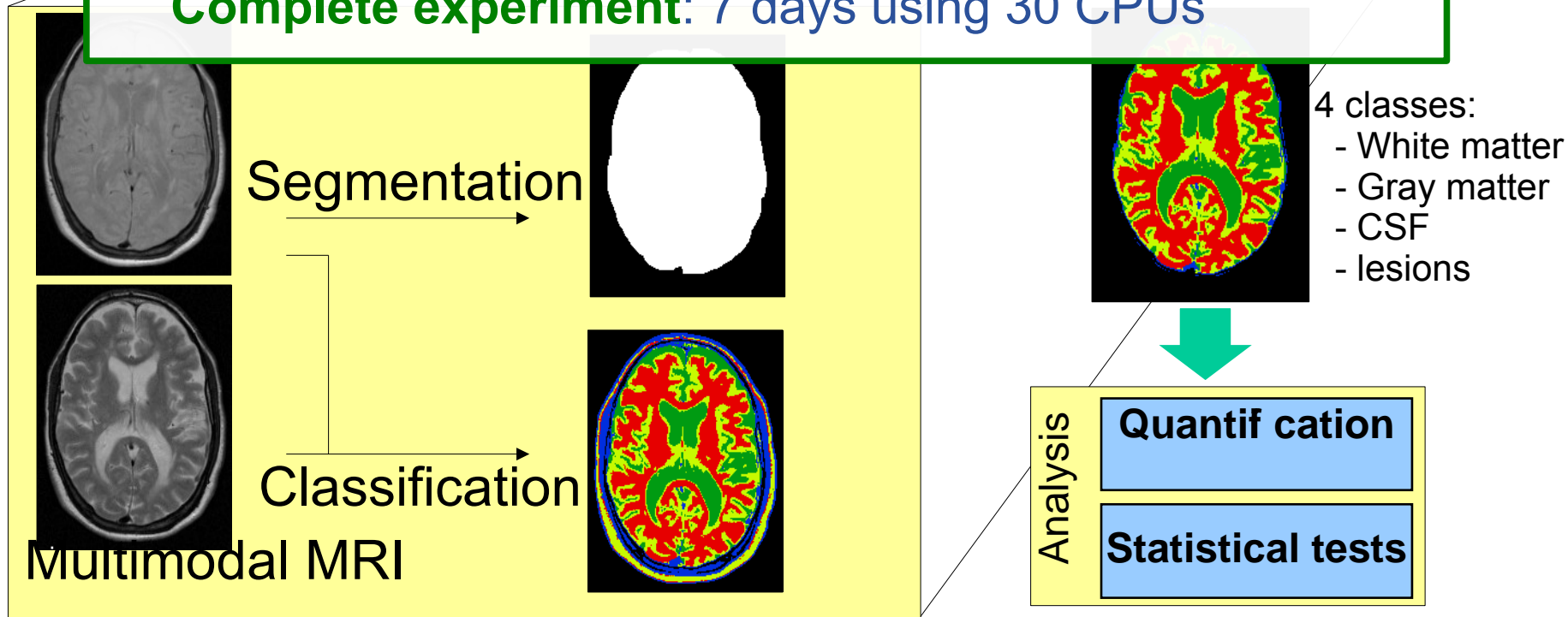
Clinical experiment: interferon beta drugs trials for multiple sclerosis

Population: 200 subjects (patients and normal control)

Images: 2400 longitudinal acquisition (1.3 TB)

Processing pipeline: 0h45 to 2h30 per image

Complete experiment: 7 days using 30 CPUs



Images analysis pipelines

MR Images

Clinical experiment: interferon beta drugs trials for multiple sclerosis

Population: 200 subjects (patients and normal control)

Images: 2400 longitudinal acquisition (1.3 TB)

Processing pipeline: 0h45 to 2h30 per image

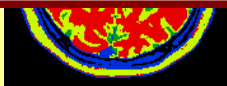
Complete experiment: 7 days using 30 CPUs

Other compute intensive health applications:

- . Clinical and statistical studies
- . Epidemiology
- . Anatomical and physiological models design
- . Validating medical image analysis procedures
- . Image-based medical treatment validation

....

Multimodal MRI



Analysis

Statistical tests

4 classes:

- White matter
- Gray matter
- CSF
- lesions

- ## Filtering, initialization



Quantification

Visualization, Decision taking

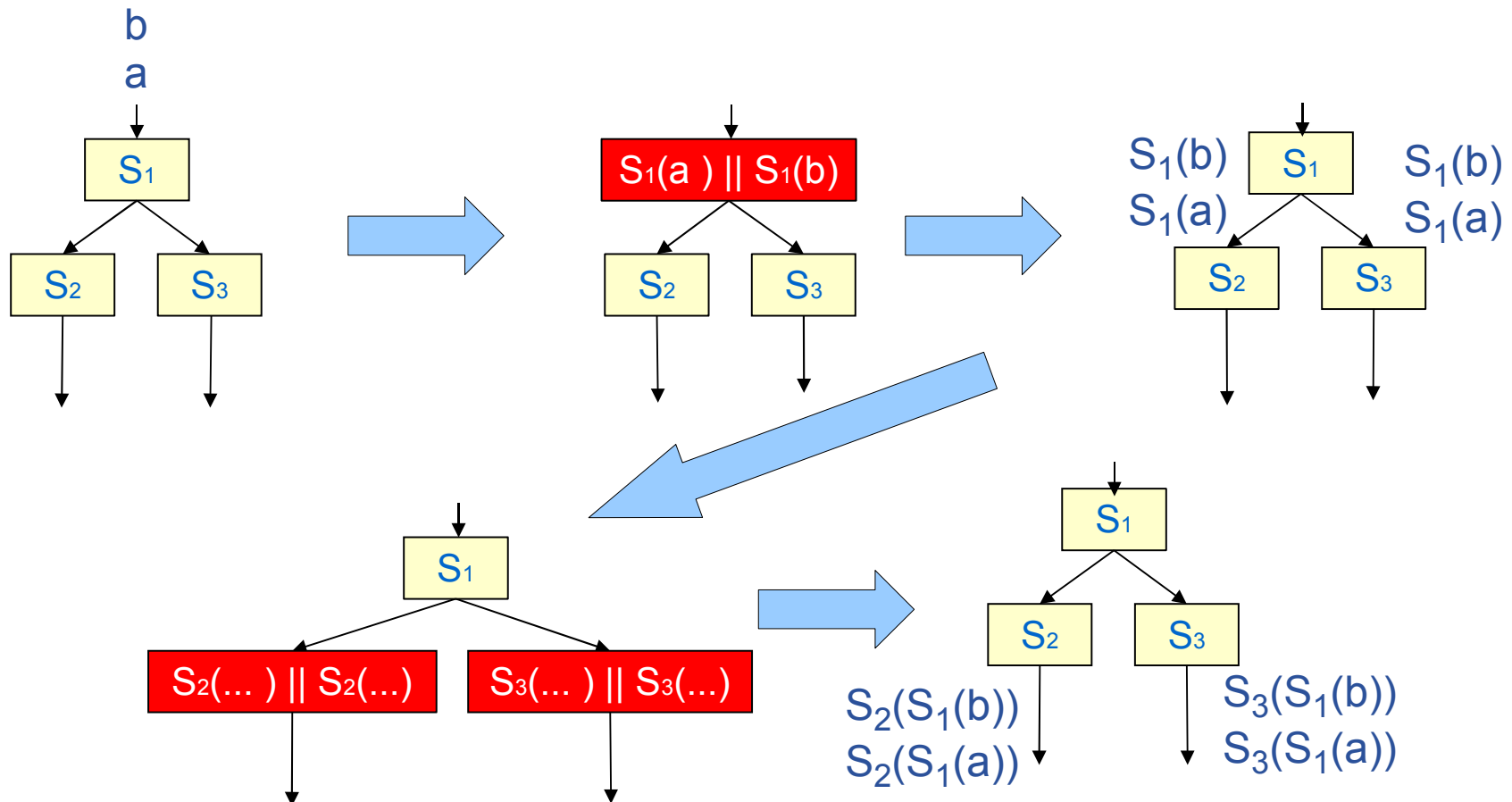


Data intensive medical imaging

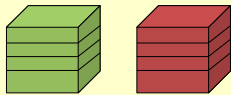
Grid Workflow Efficient Enactment for Data Intensive Applications

- **Application community**
 - Compute and data intensive applications
 - Non-expert end users
 - Distributed (medical centers)
- **Coarse grain parallelism**
 - Grid computing
 - Massive data parallelism
- **Platform independence**
 - Common representation / submission interface to
 - Different grids
 - Multiple grids
- **Pipelines are pure data flows**
 - Successive image processing filters
 - Data intensive and **data driven**

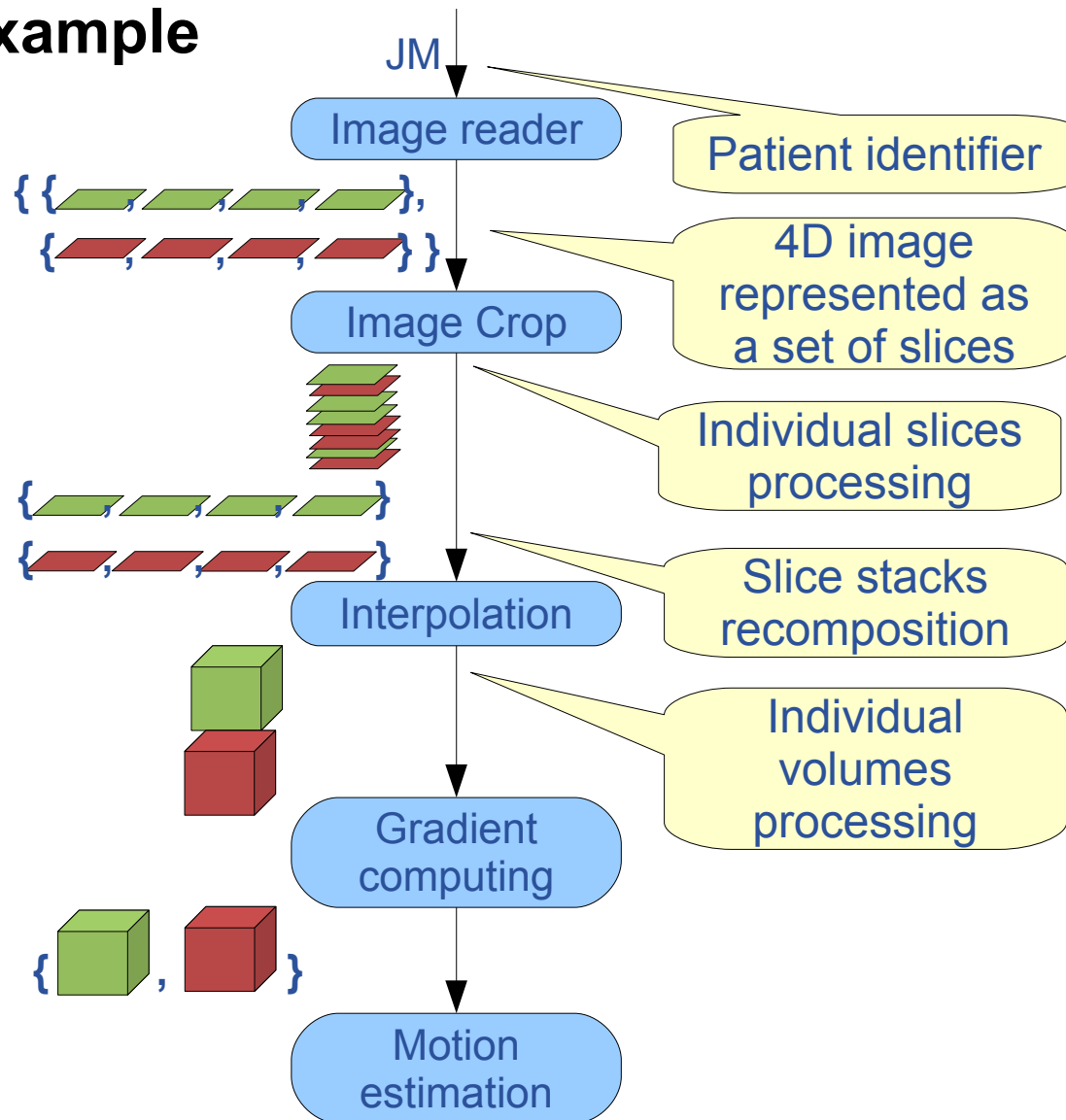
- Data flow resolved dynamically
- Services invoked multiple times

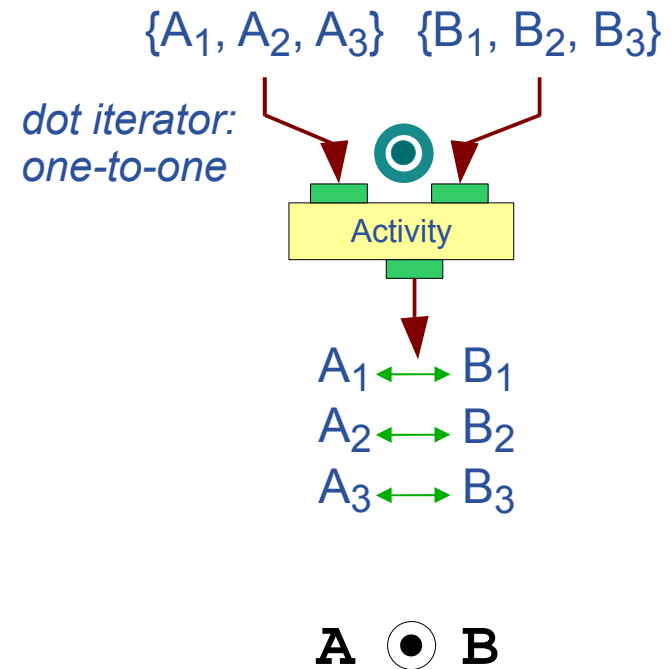
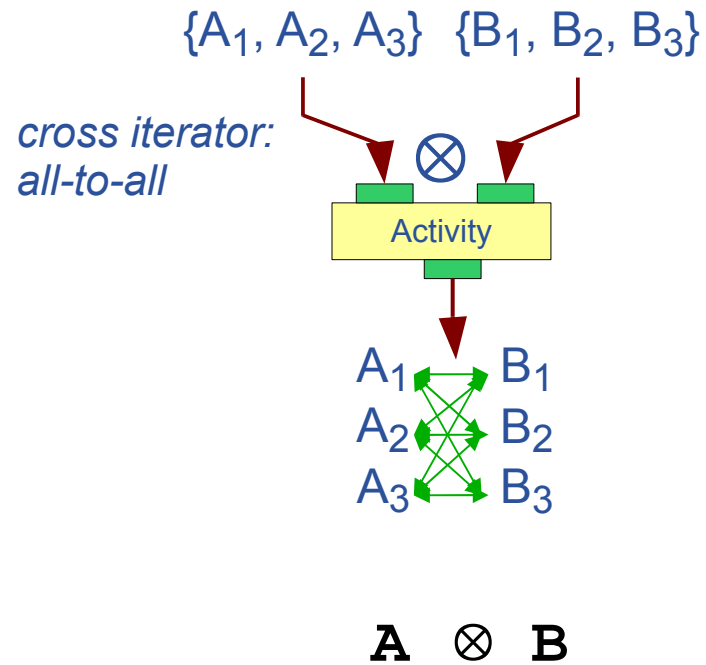


Cardiac sequences example

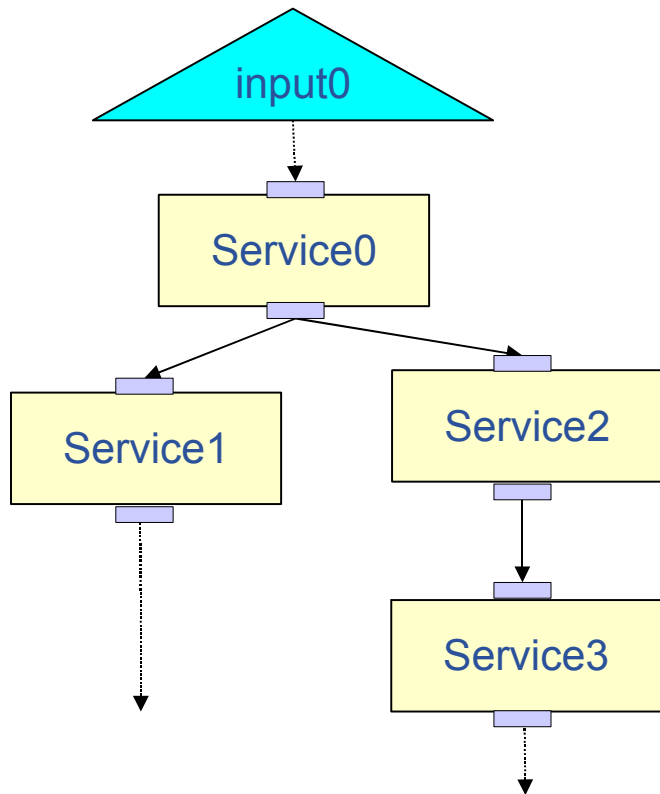


The input data set is a 4D image composed of 2 volumes (labelled green and red). Each volume is composed of 4 slices.

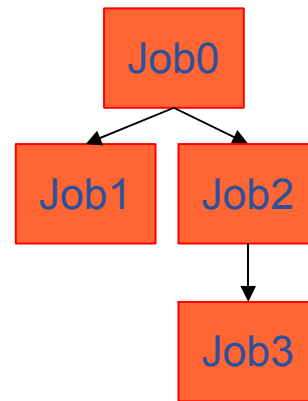




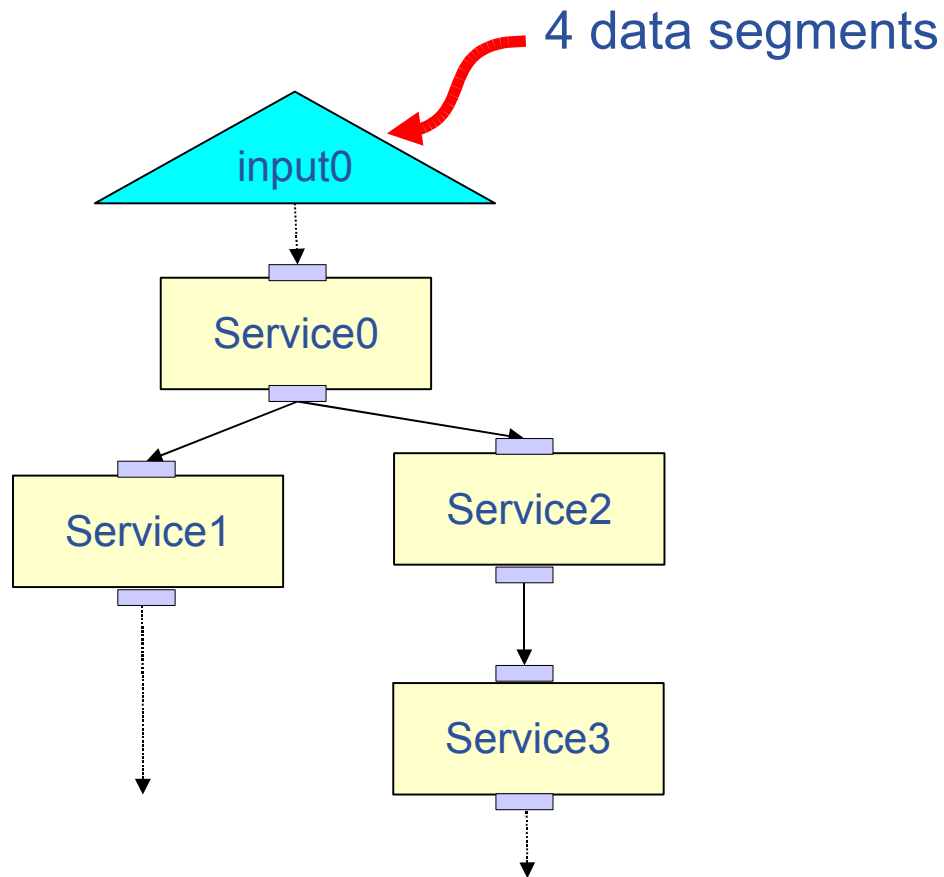
- Graph of services (+ data)**



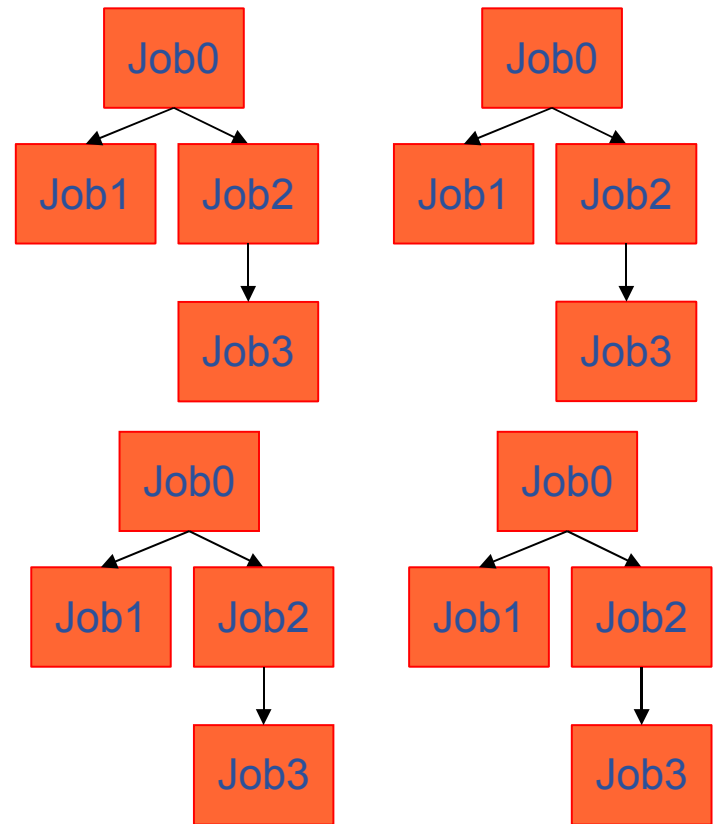
DAG of tasks



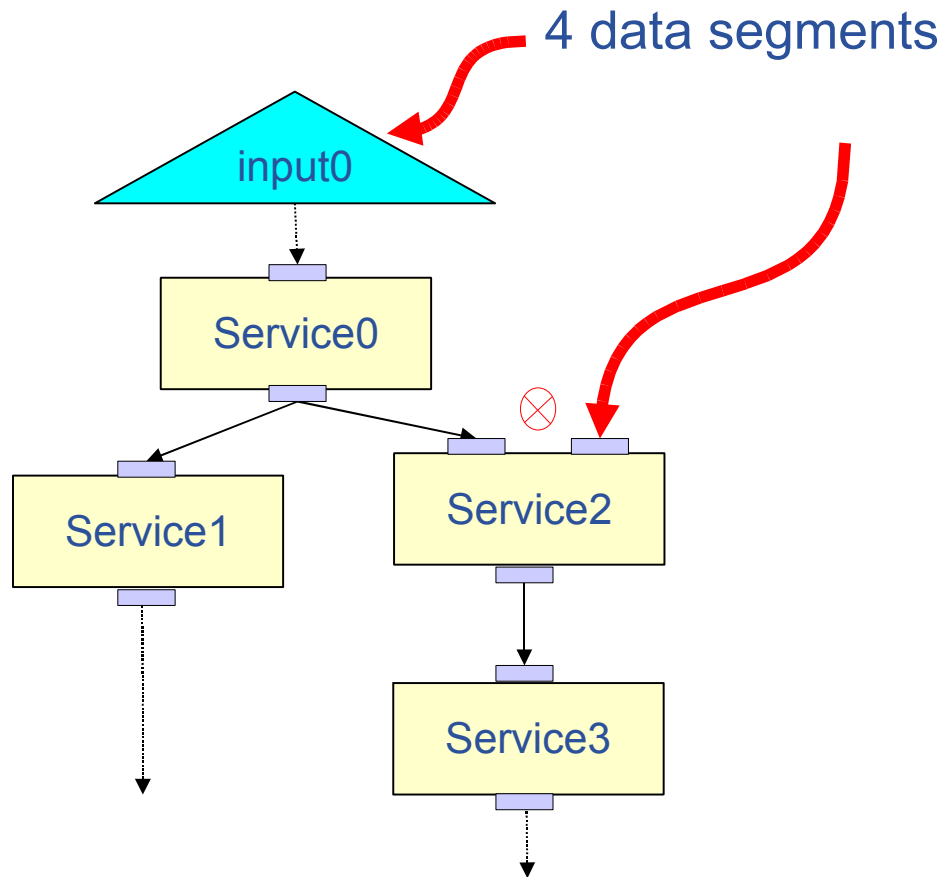
- Graph of services (+ data)**



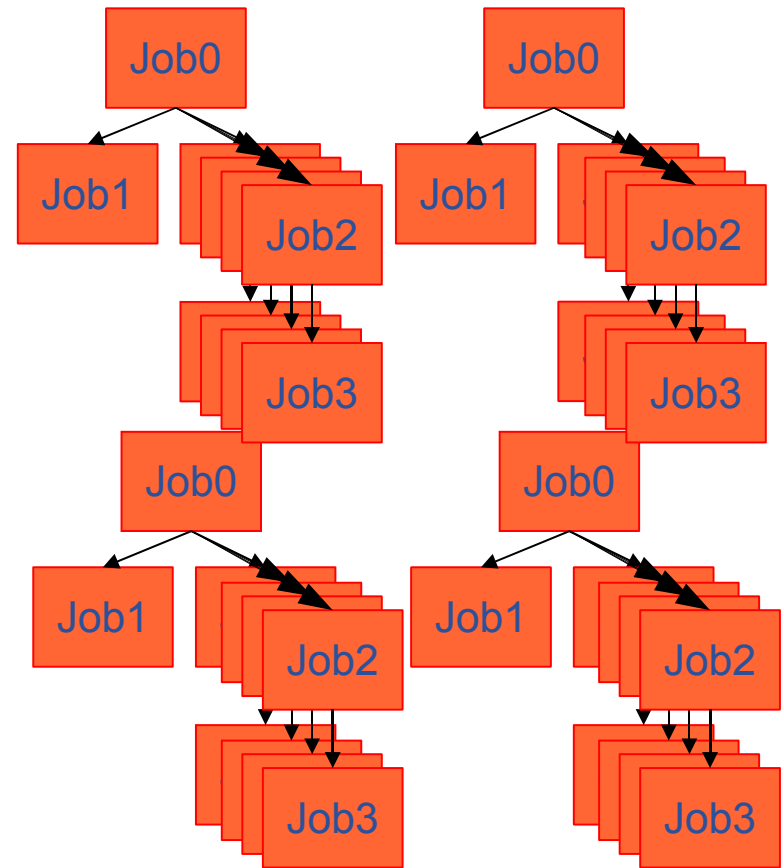
DAG of tasks



- Graph of services (+ data)**



DAG of tasks



-
- The diagram illustrates the dot product operation for string similarity. It starts with two input strings, "String 0" and "String 1", represented by cyan triangles. These strings are processed by "Concat 0" (a yellow rectangle) to produce two vectors, represented by purple rectangles. These vectors are then combined using a dot product operation, indicated by a circle with a plus sign (\oplus), to produce a single vector. This vector is then processed by "Concat 1" (another yellow rectangle) to produce two more vectors, represented by purple rectangles. Finally, these two vectors are used to calculate the "right" and "reverse" similarity scores, represented by cyan diamonds.

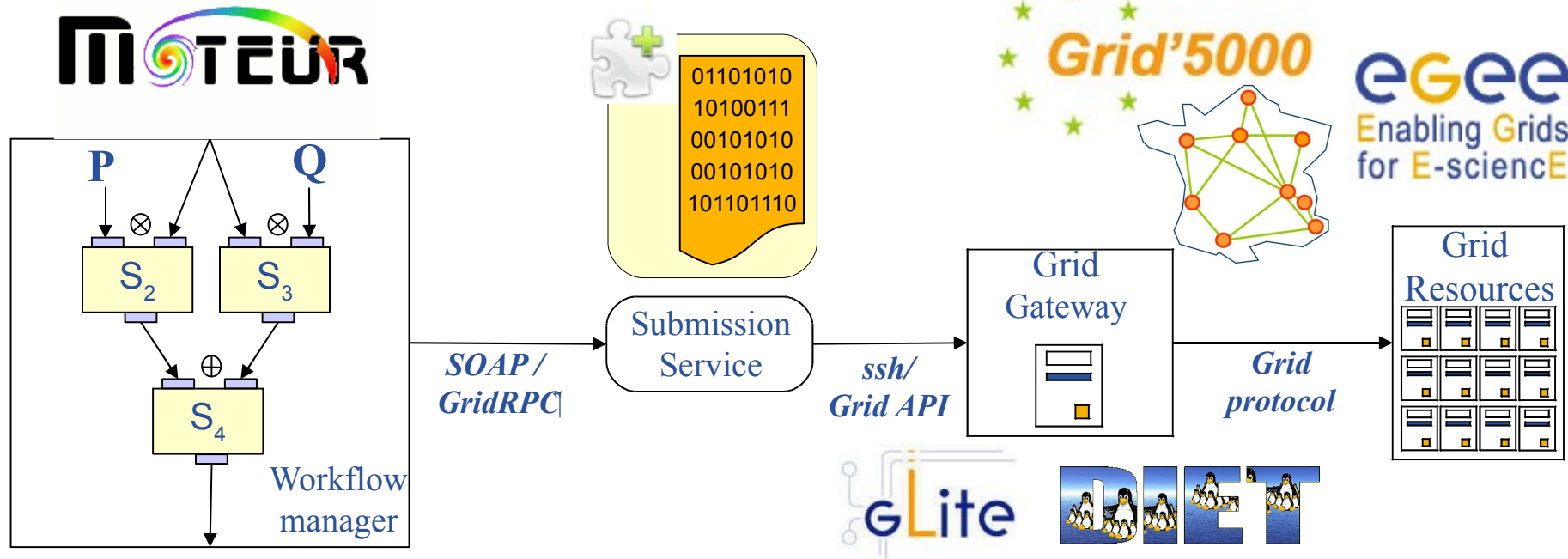


MOTEUR workflow manager

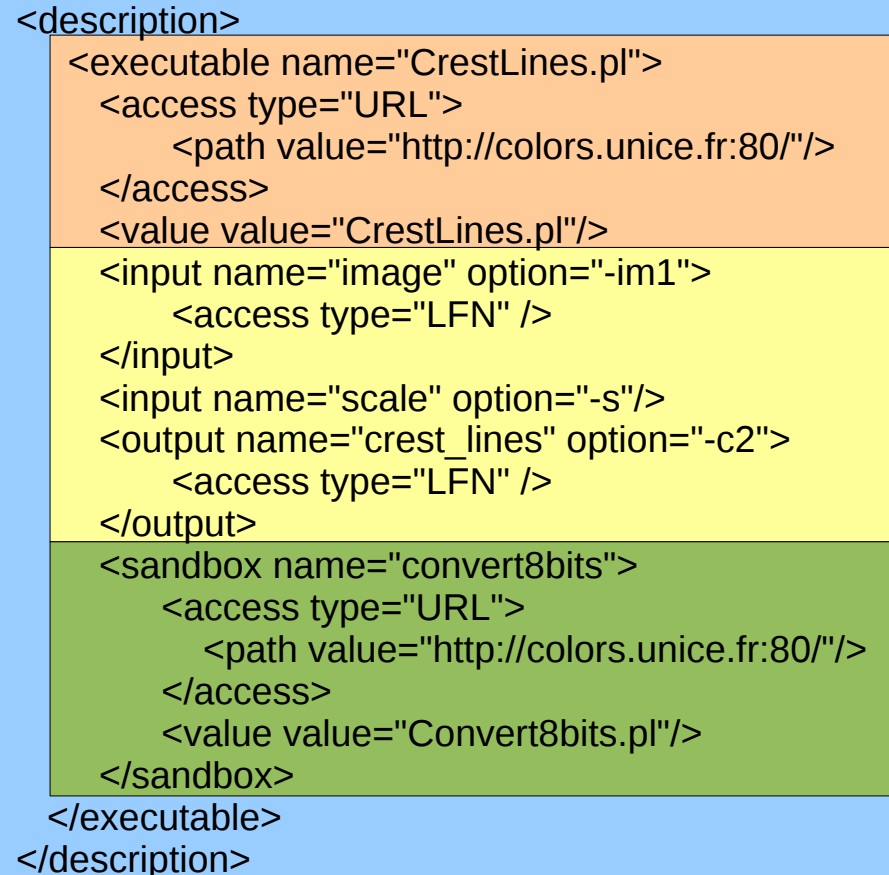
Grid Workflow Efficient Enactment for Data Intensive Applications

- **Open source workflow enactor**
 - Code + docs + tutorial: <http://egee1.unice.fr/MOTEUR>
 - Developed at the I3S CNRS laboratory
 - With the support of French national projects
 - AGIR: <http://www.agir.org>
 - GWENDIA: <http://gwendia.polytech.unice.fr>
- **Targets**
 - Ease of use, flexibility, service-oriented approach
 - Performance, transparent exploitation of application parallelism
- **Supports**
 - Scufi language (from myGrid/Taverna)
 - Service based invocation (WS)
 - Grid middlewares (EGEE / Grid'5000)

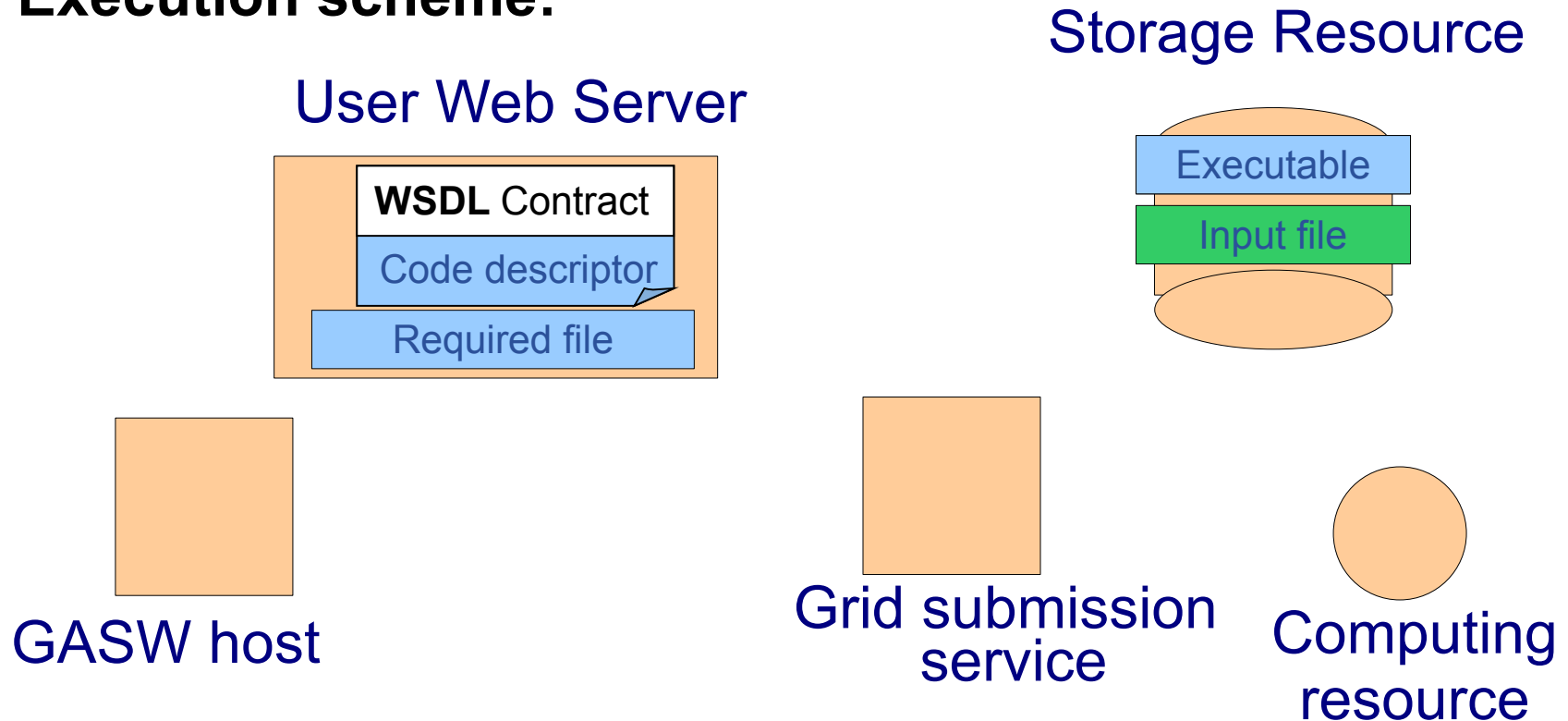




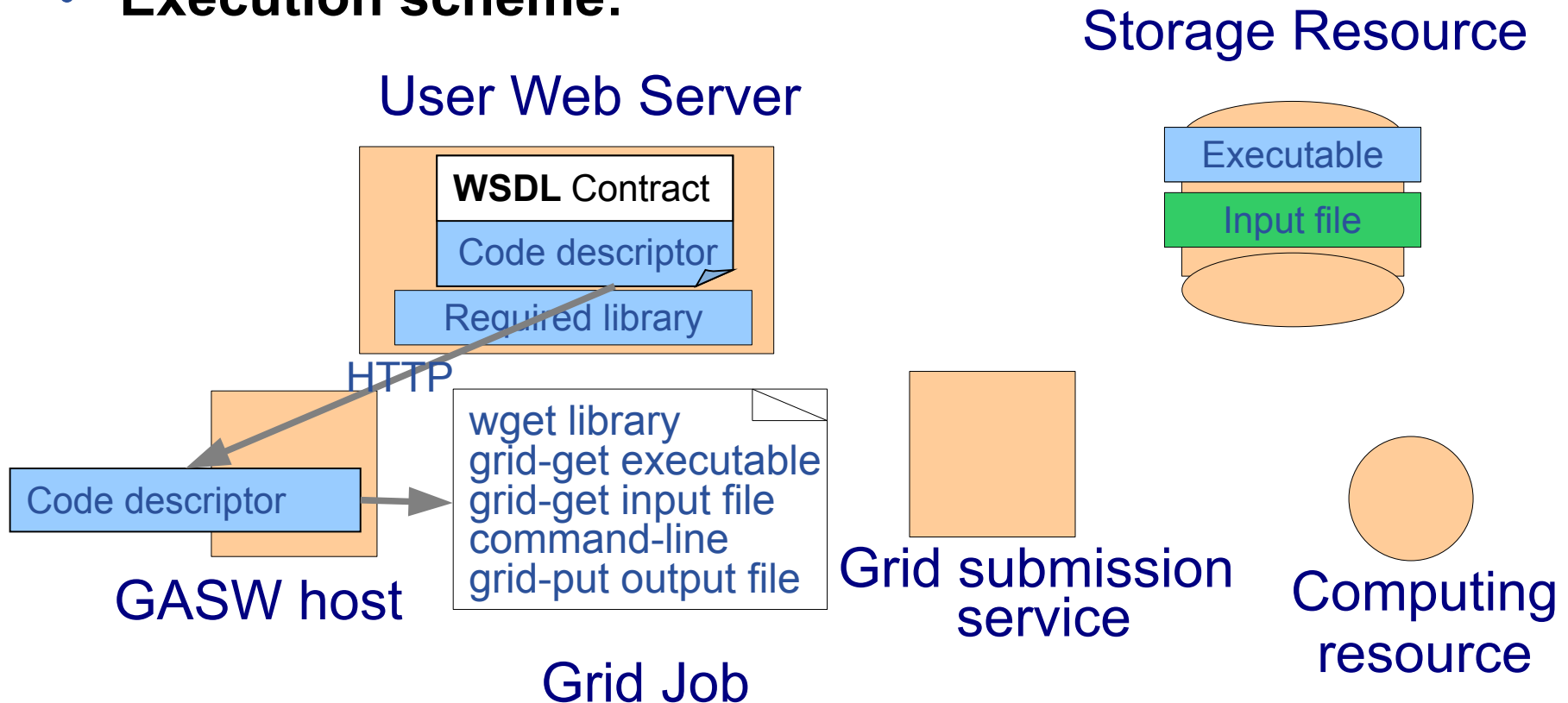
- **Batch-oriented infrastructure: no deployment**
 - EGEE / gLite middleware
 - Grid'5000 / OAR, using idle resources
- **Service oriented infrastructure: pre-deployment**
 - Web Services
 - DIET Services



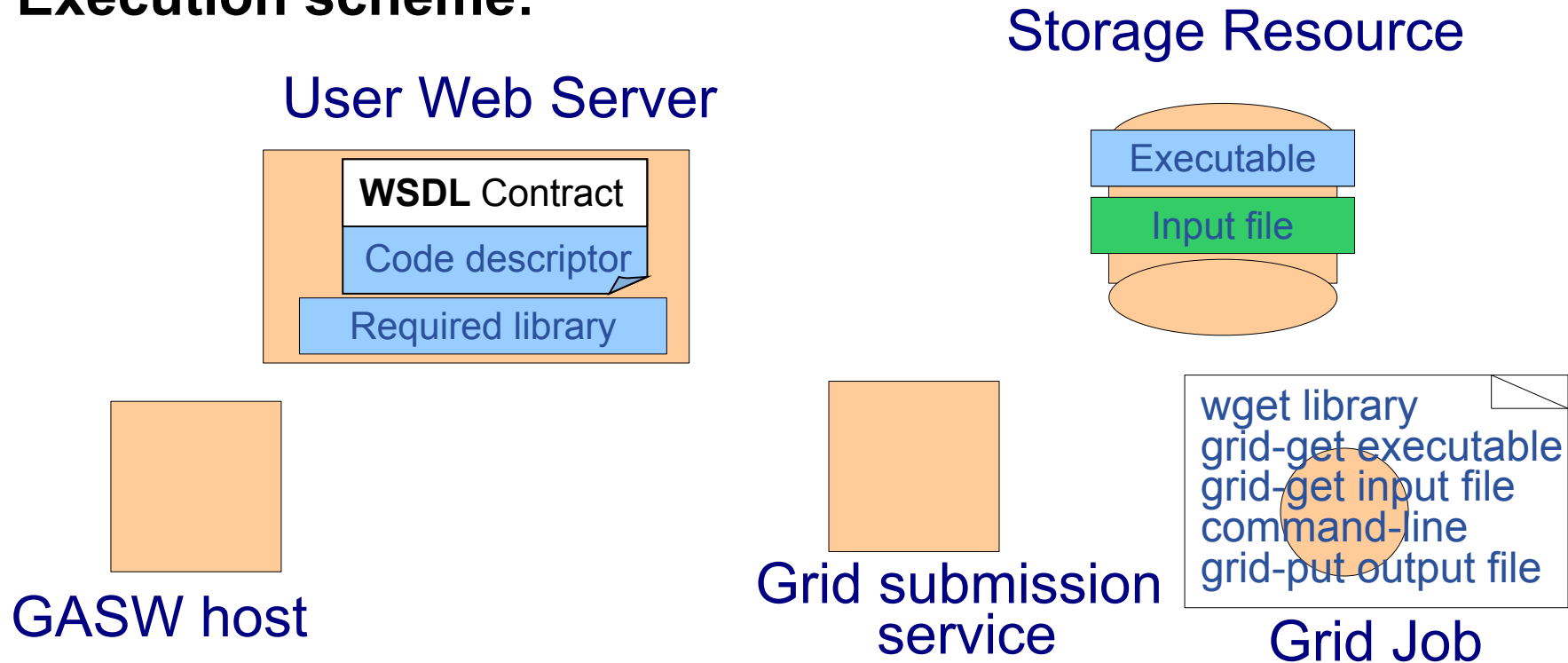
- **Generic Application Service Wrapper**
 - Provide service wrapper to non instrumented code
 - Handle data transfer (references to grid data)
- **Execution scheme:**



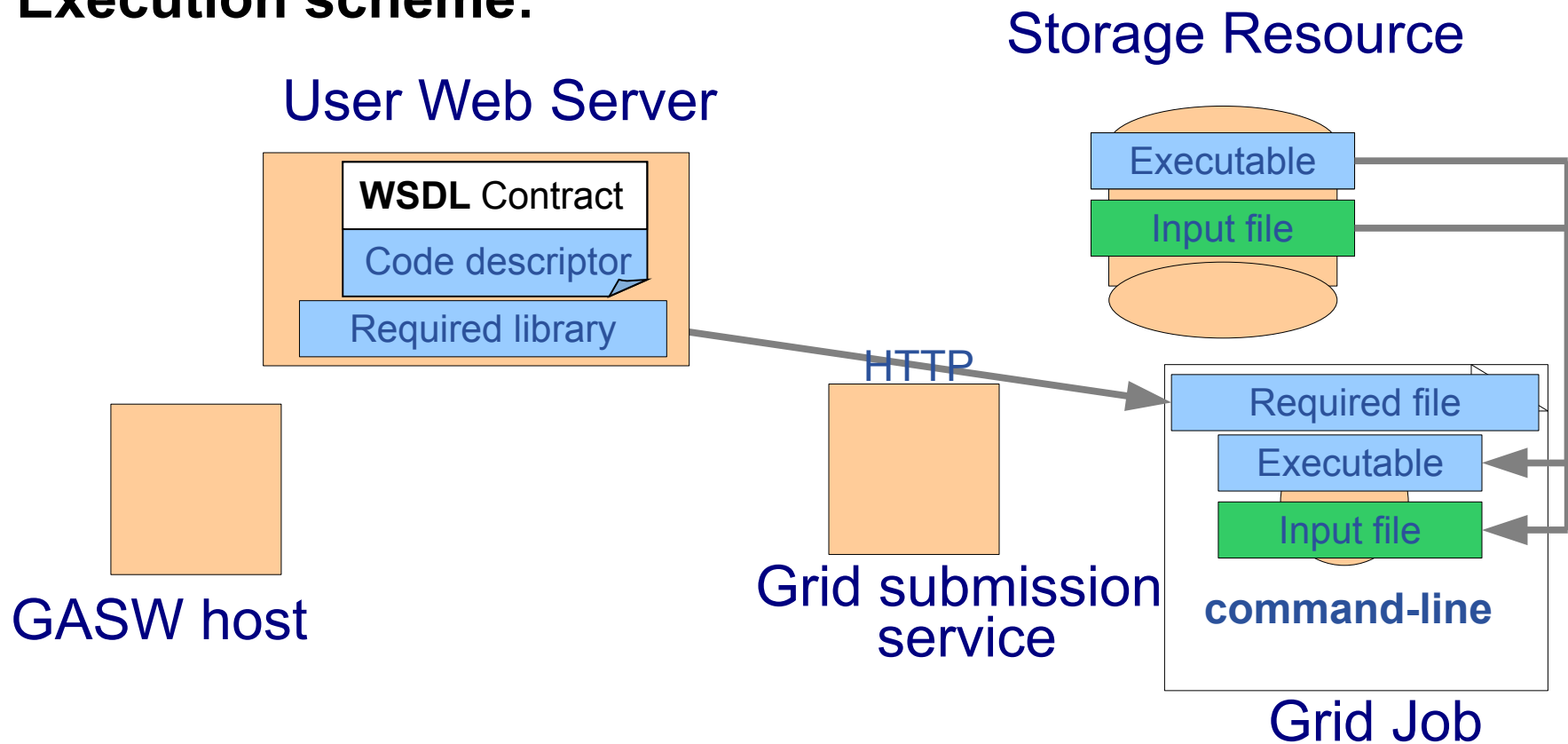
- **Generic Application Service Wrapper**
 - Provide service wrapper to non instrumented code
 - Handle data transfer (references to grid data)
- **Execution scheme:**



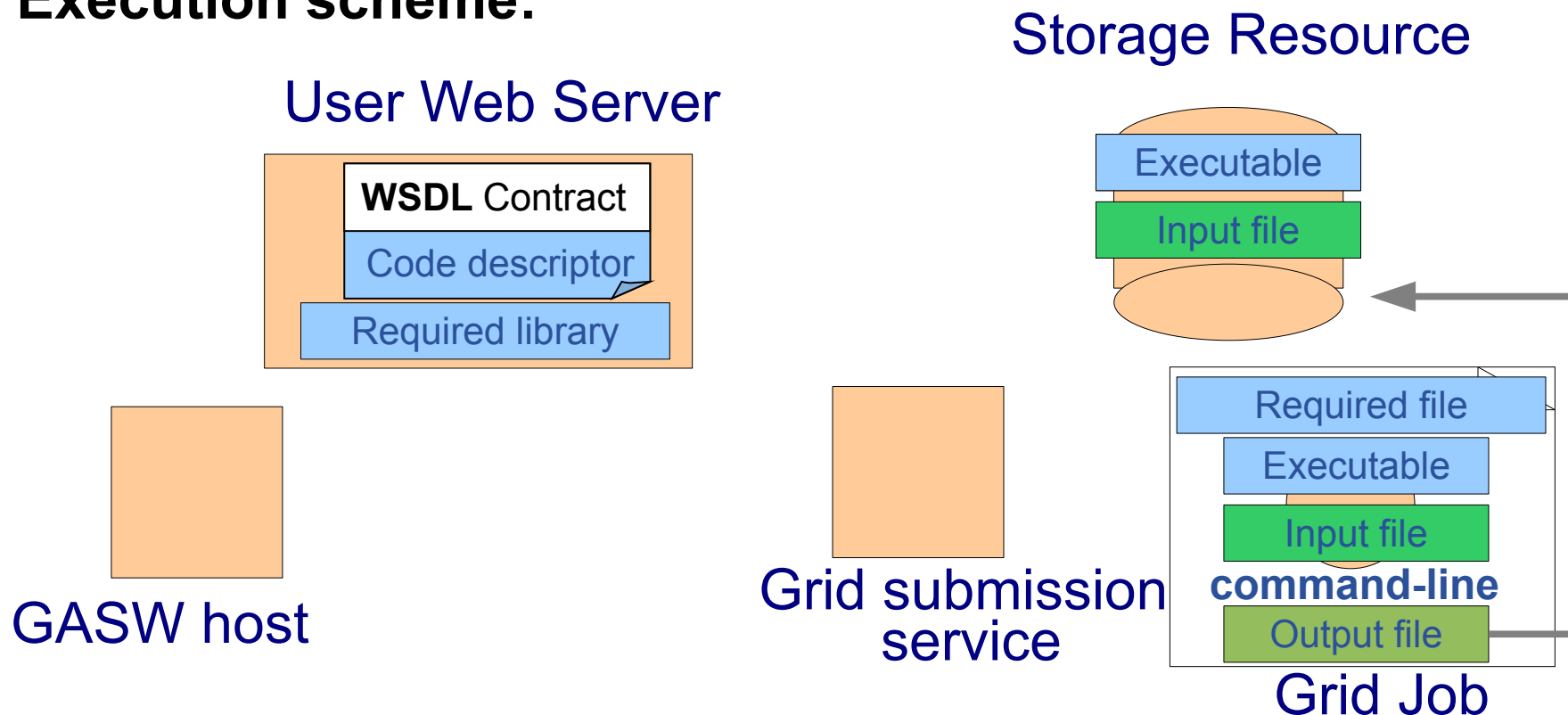
- **Generic Application Service Wrapper**
 - Provide service wrapper to non instrumented code
 - Handle data transfer (references to grid data)
- **Execution scheme:**



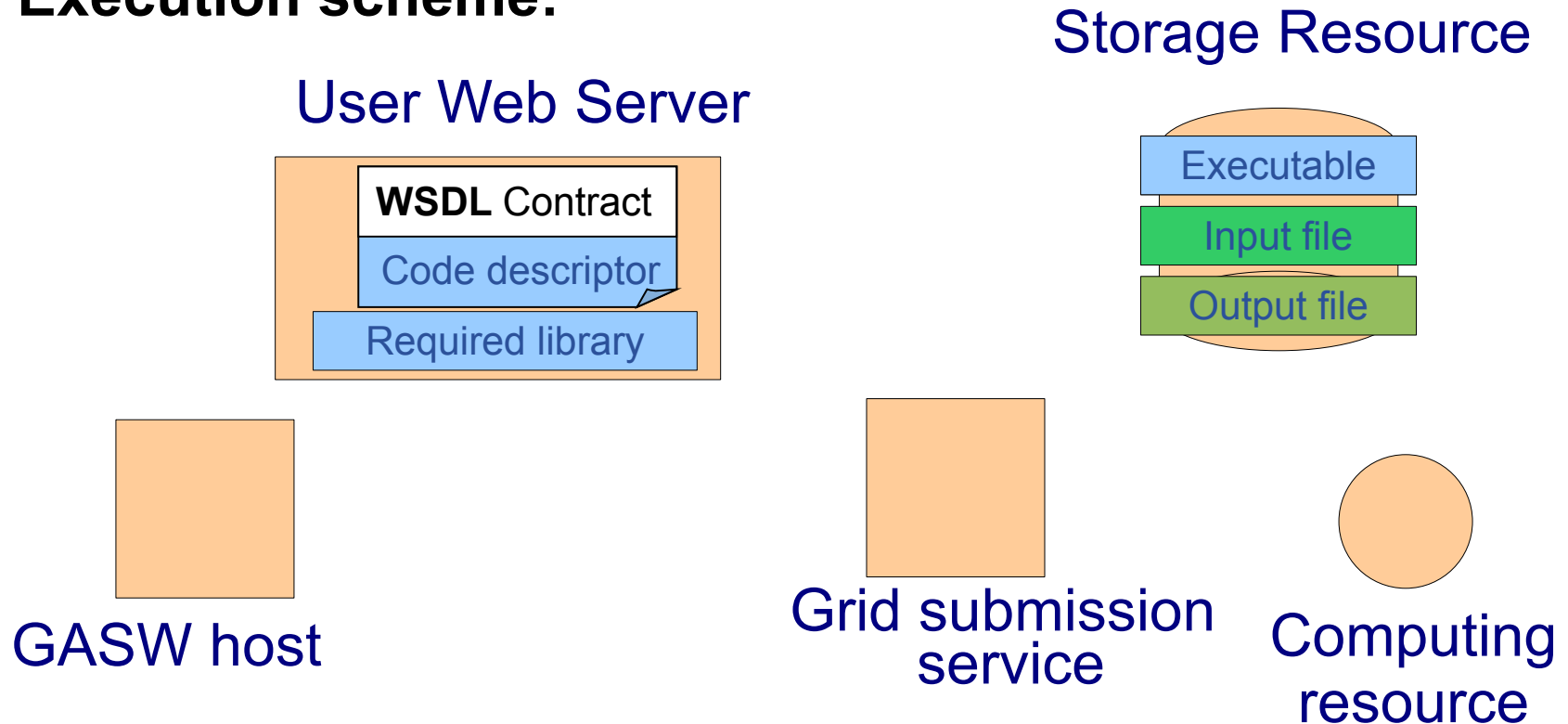
- **Generic Application Service Wrapper**
 - Provide service wrapper to non instrumented code
 - Handle data transfer (references to grid data)
- **Execution scheme:**



- **Generic Application Service Wrapper**
 - Provide service wrapper to non instrumented code
 - Handle data transfer (references to grid data)
- **Execution scheme:**

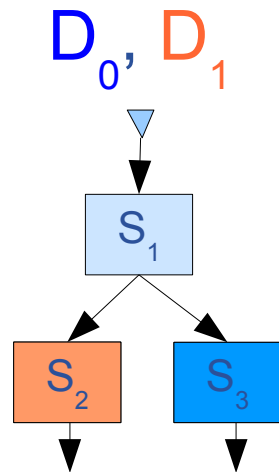


- **Generic Application Service Wrapper**
 - Provide service wrapper to non instrumented code
 - Handle data transfer (references to grid data)
- **Execution scheme:**

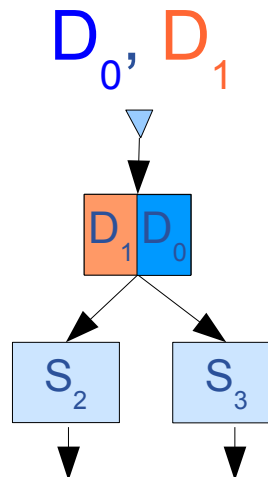


- A workflow naturally provides application parallelization
- 3 kinds of parallelism exploited

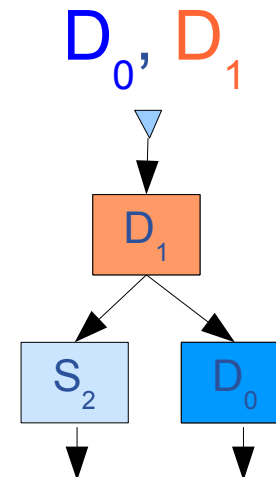
Workflow parallelism



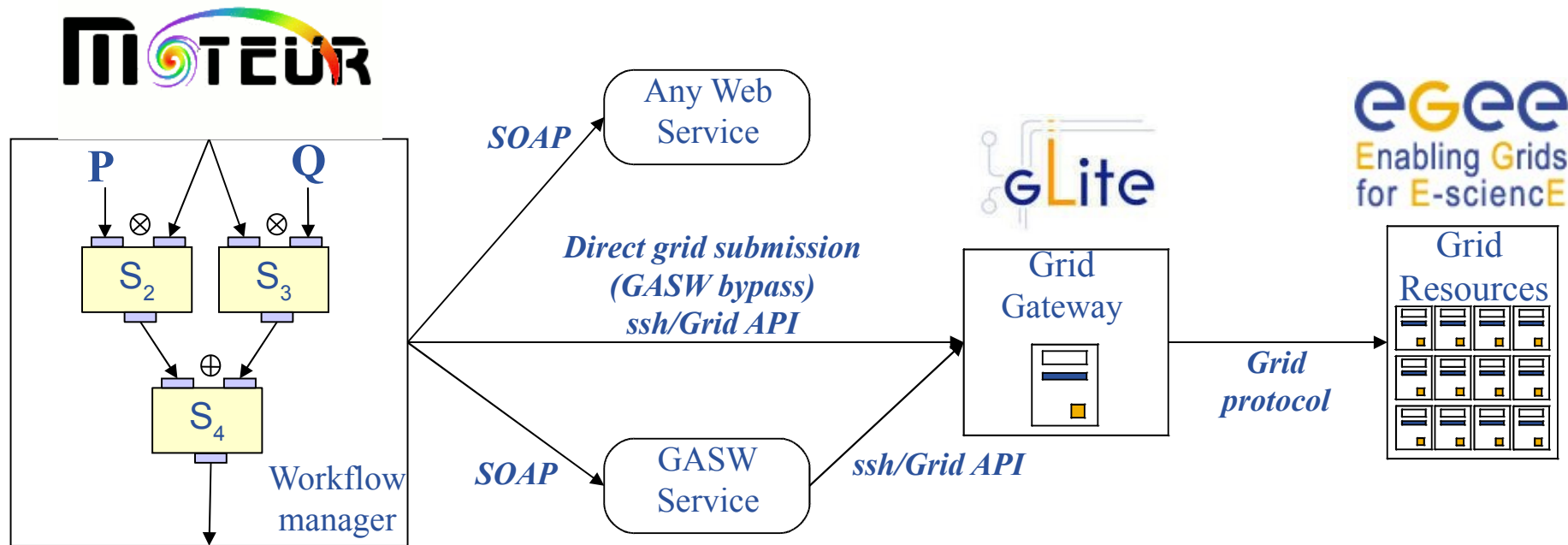
Data parallelism



Service parallelism



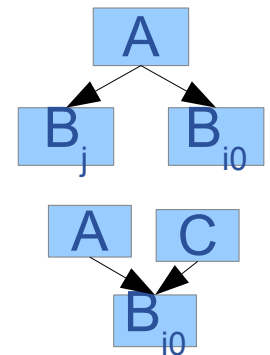
- Data sets are composed dynamically



- **GASW service is recognized by MOTEUR**
 - Not a real blackbox
- **Submission can be directly made to the grid**
 - Grid mode: Web Service call bypassed
 - No difference from a user point of view (no change in ScufI)

-
- The diagram illustrates the interaction between a Workflow manager and Application services. It is divided into three main sections: Workflow manager, Application services, and Grouped services.
- Workflow manager:** Contains two process blocks, P_1 and P_2 . P_1 is connected to P_2 via a downward arrow labeled "Services invocation".
 - Application services:** Contains two service blocks. Each block has a "standard interface" on the left and a list of operations on the right:
 - Service 1: Input data transfer, Code 1 submission, Output data transfer.
 - Service 2: Input data transfer, Code 2 submission, Output data transfer.
 - Grouped services:** Contains a dashed box enclosing P_1 and P_2 . This group is connected to a single service block in the Application services section.
 - Grouped service block:** Has a "standard interface to generic wrapper service" on the left and a list of operations on the right:
 - Input data transfer
 - Code 1 + code 2 submission
 - Output data transfer
- A label "command lines generation" points to the "Code 1 + code 2 submission" operation in the grouped service block.

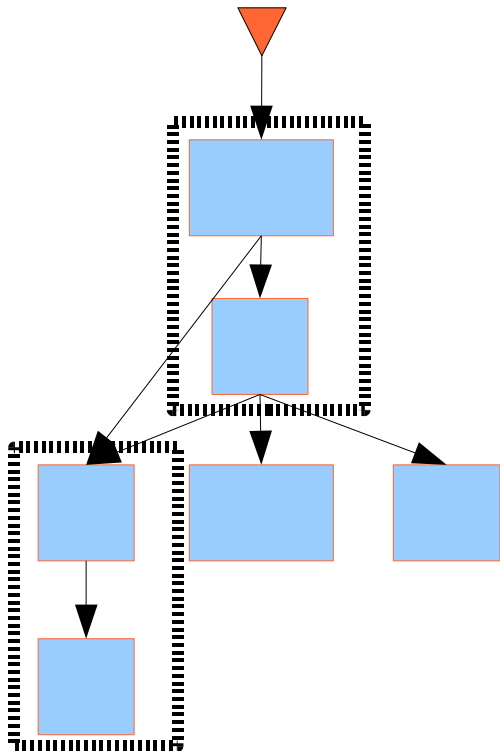
- **Goal: find a grouping rule that does not break parallelism**
- **Let A be a service of the workflow and $\{B_0, \dots, B_n\}$ its children**
- **For grouping A and B_{i0} : no parallelism loss \Leftrightarrow**
 - (1) B_{i0} is an ancestor of every B_j
 - (2) Every ancestor of B_{i0} is an ancestor of A (or A itself)
- **No parallelism loss \Rightarrow**
 - $\neg(1) \Rightarrow$ parallelism between B_j and B_{i0} is broken
 - $\neg(2) \Rightarrow$ parallelism between A and C is broken
- **(1) & (2) \Rightarrow no parallelism loss**
- **This rule is recursively applied on the workflow graph**



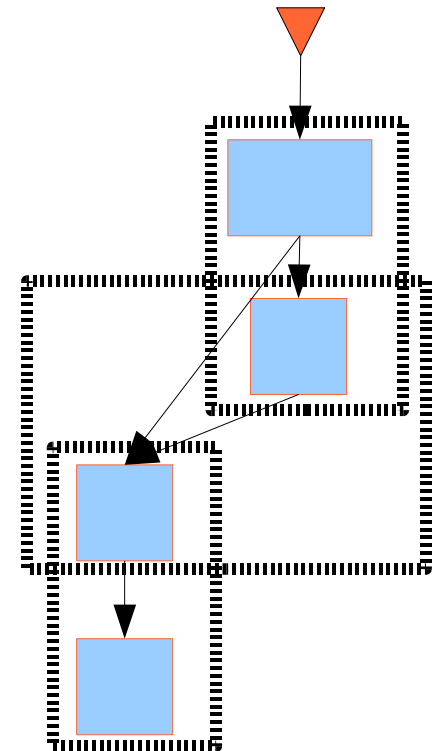
Job Grouping Experiments

Grid Workflow Efficient Enactment for Data Intensive Applications

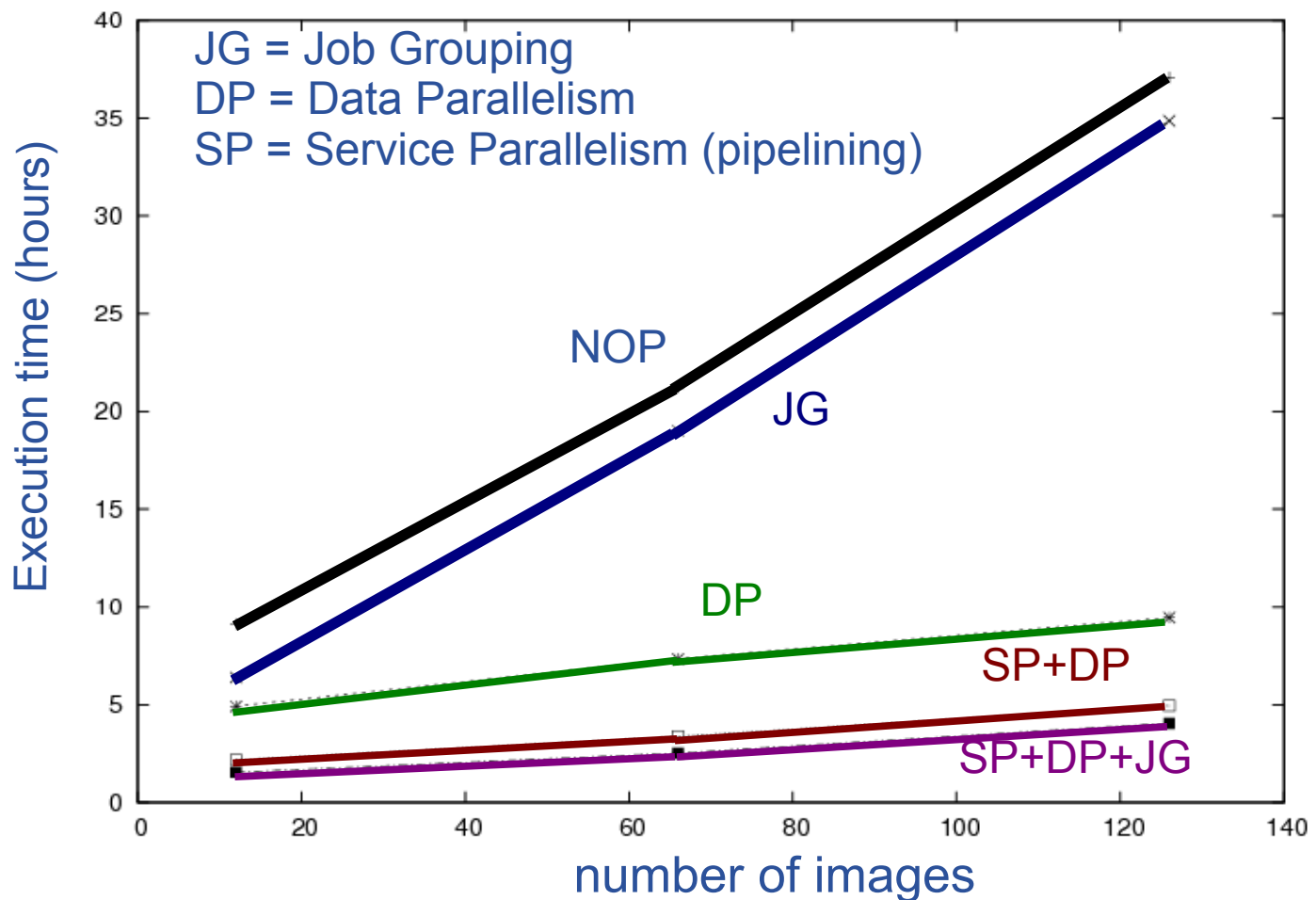
- 6 services – 2 groups
- 4 job submissions/input data set submission/input data set



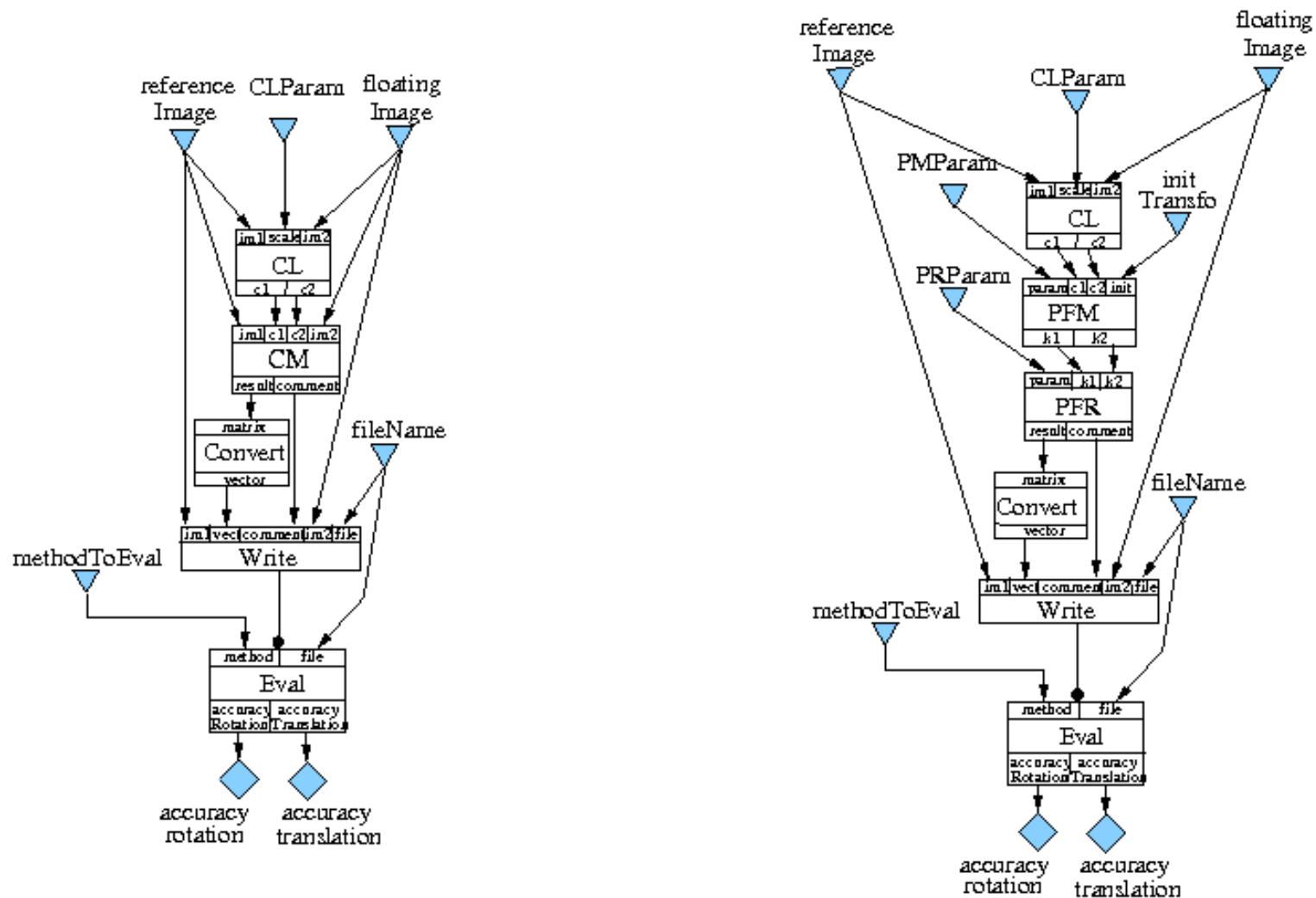
- 4 services - 3 groups
- Recursive application



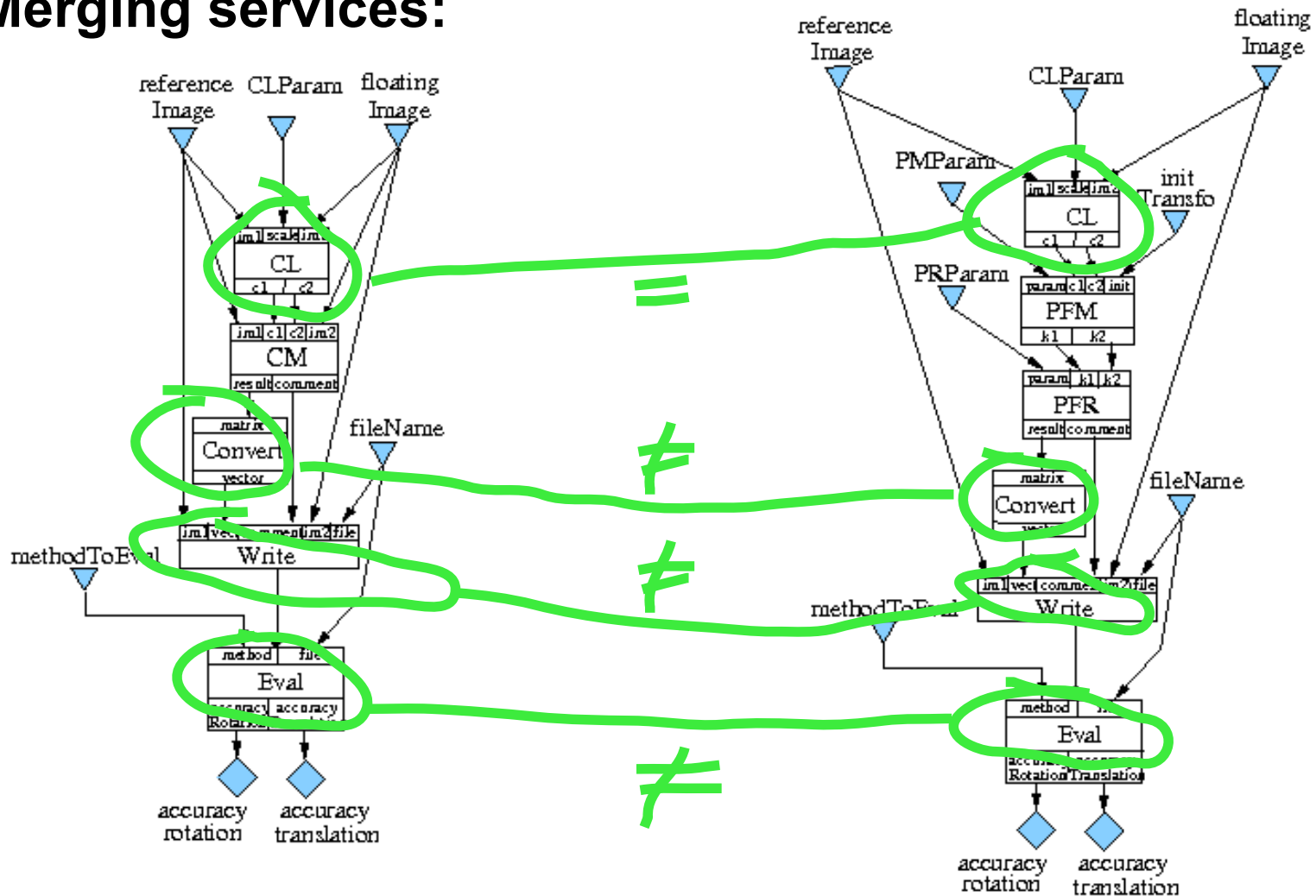
- On the EGEE infrastructure



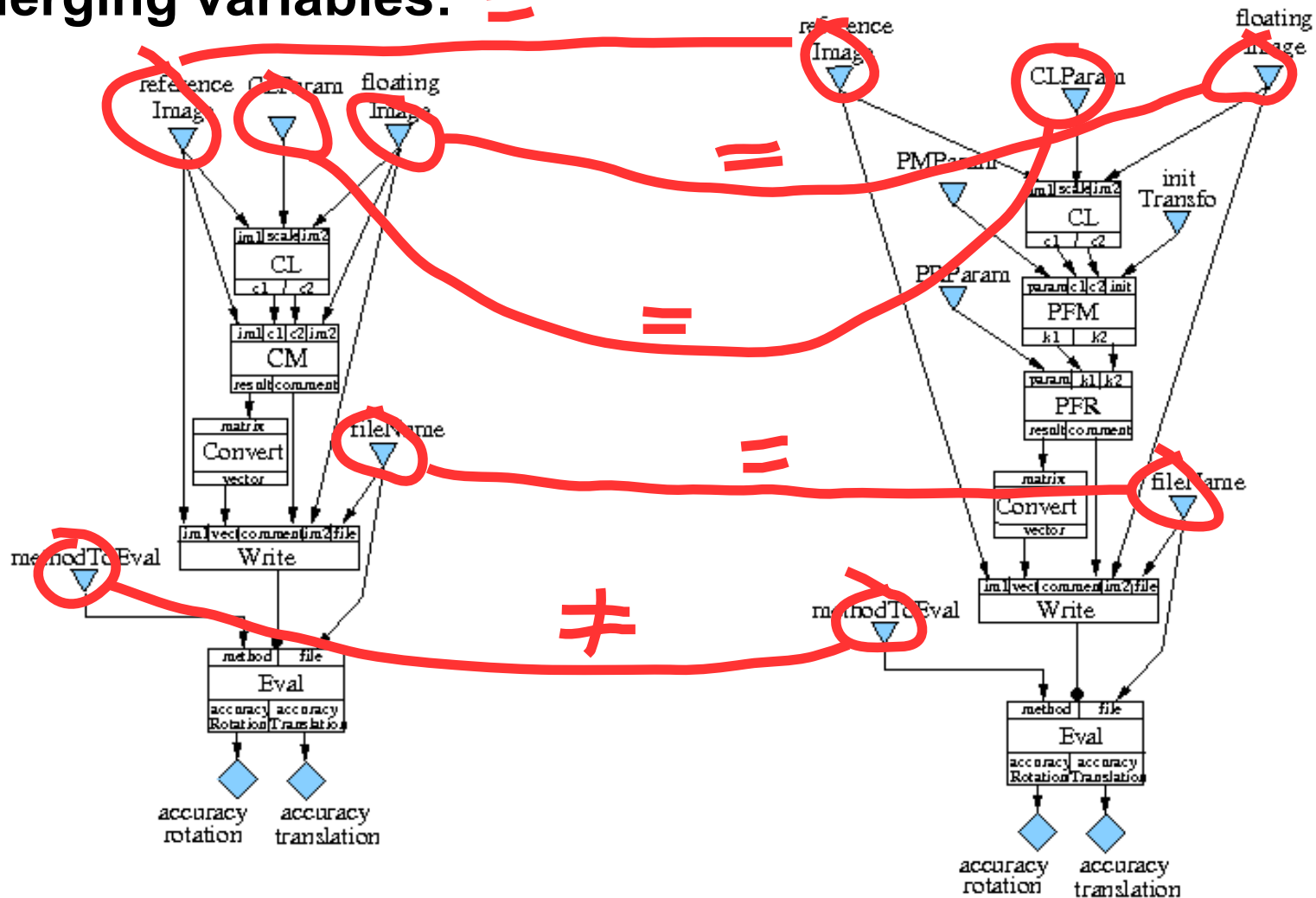
- Two basic registration orchestrations to be merged



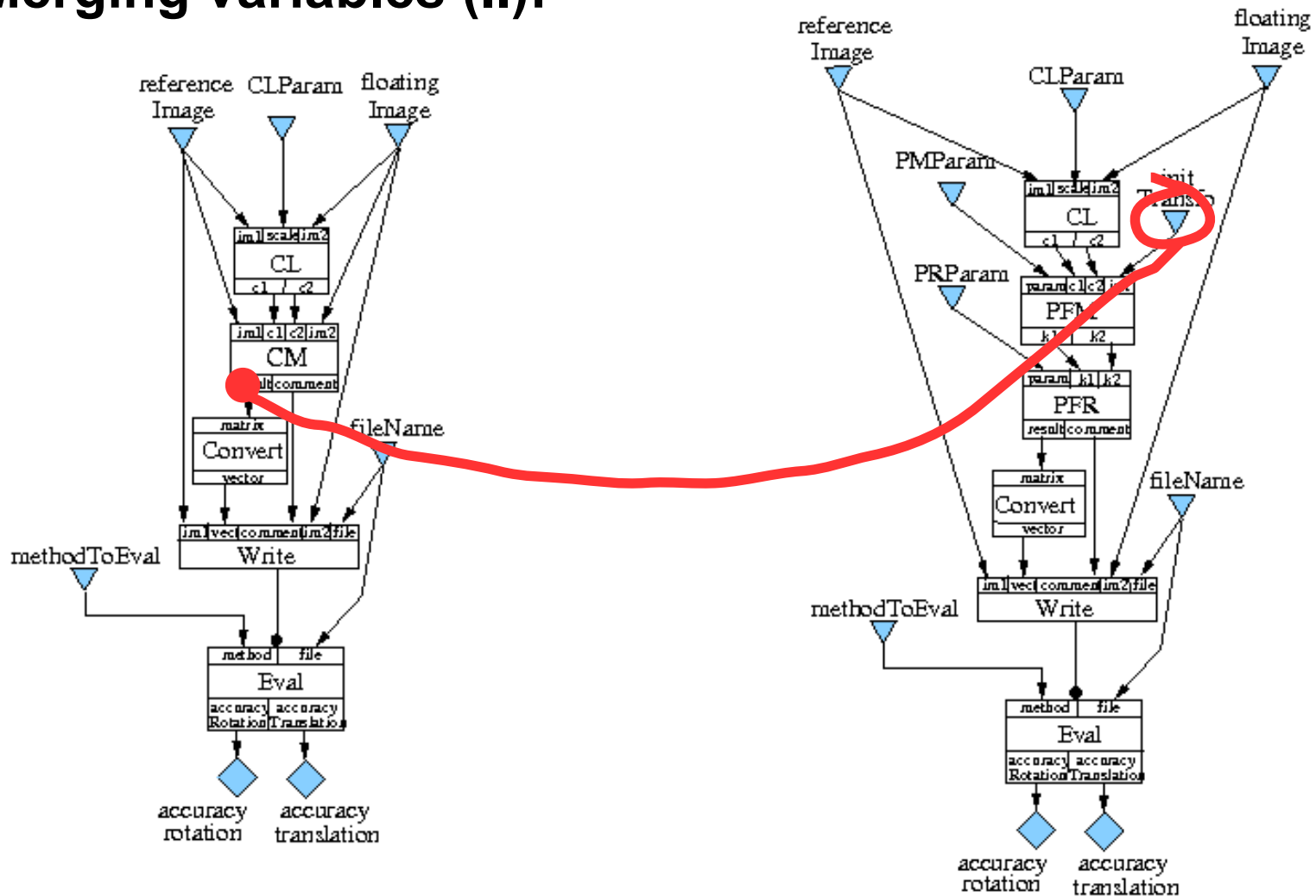
- Merging services:



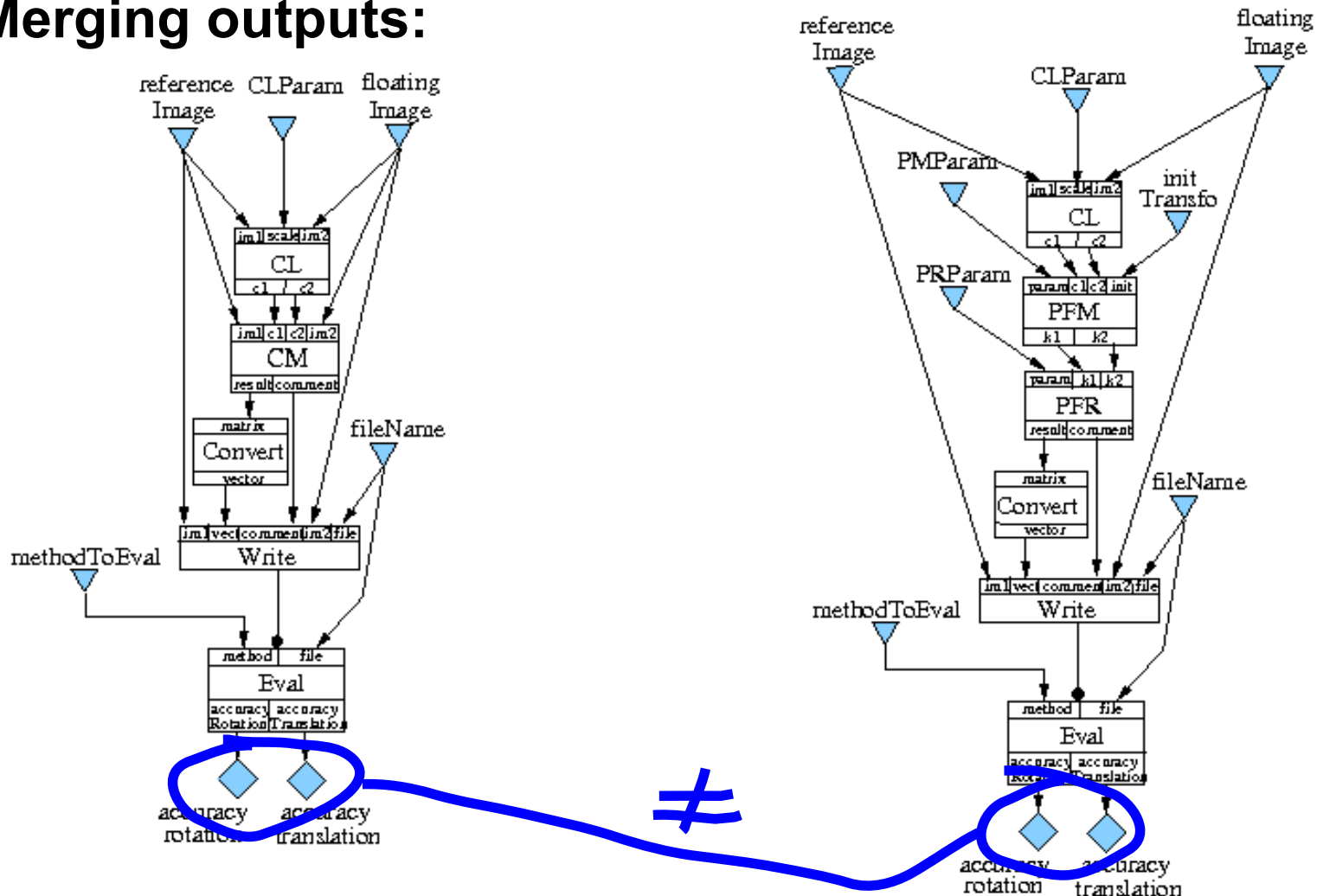
- Merging variables: =



- Merging variables (II):

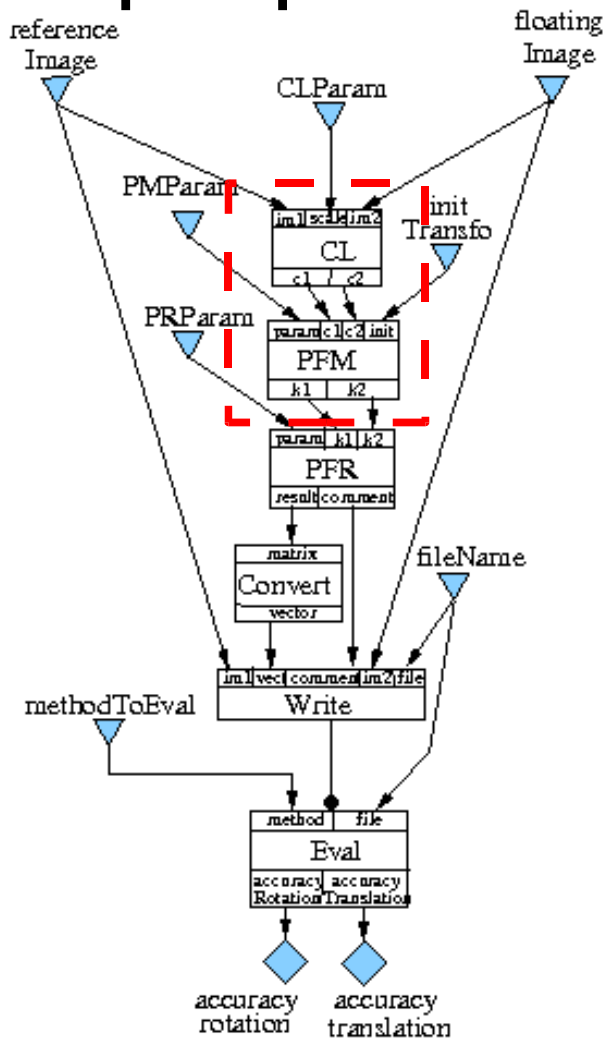


- Merging outputs:**



- A fully automatic merging procedure is not suitable**

- Graph representation**



- OMSM representation in Prolog**
[C. Nemo et al, SCC'07]

%Orchestration declaration

```
orch(pfm,[im0,im1,fN,method,transfo,cl,pm,pr],
    [i1,i2,i3,i4,...]).
```

% Instructions

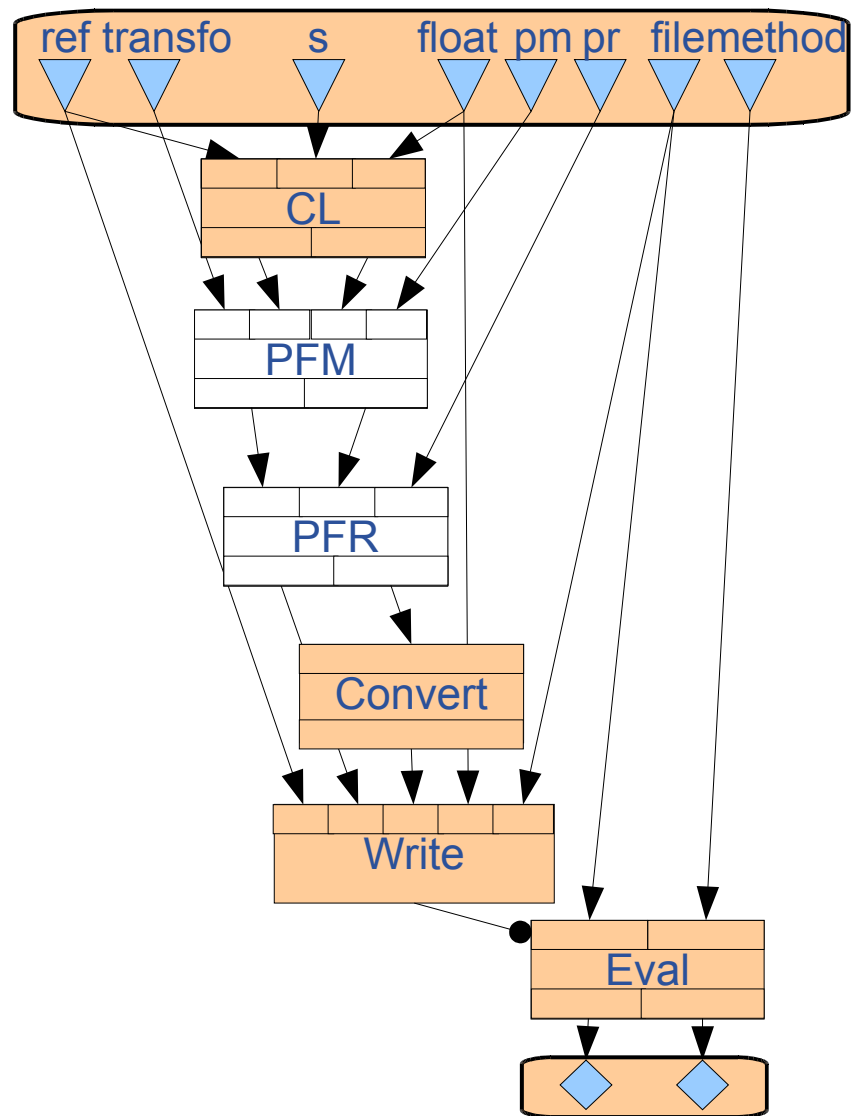
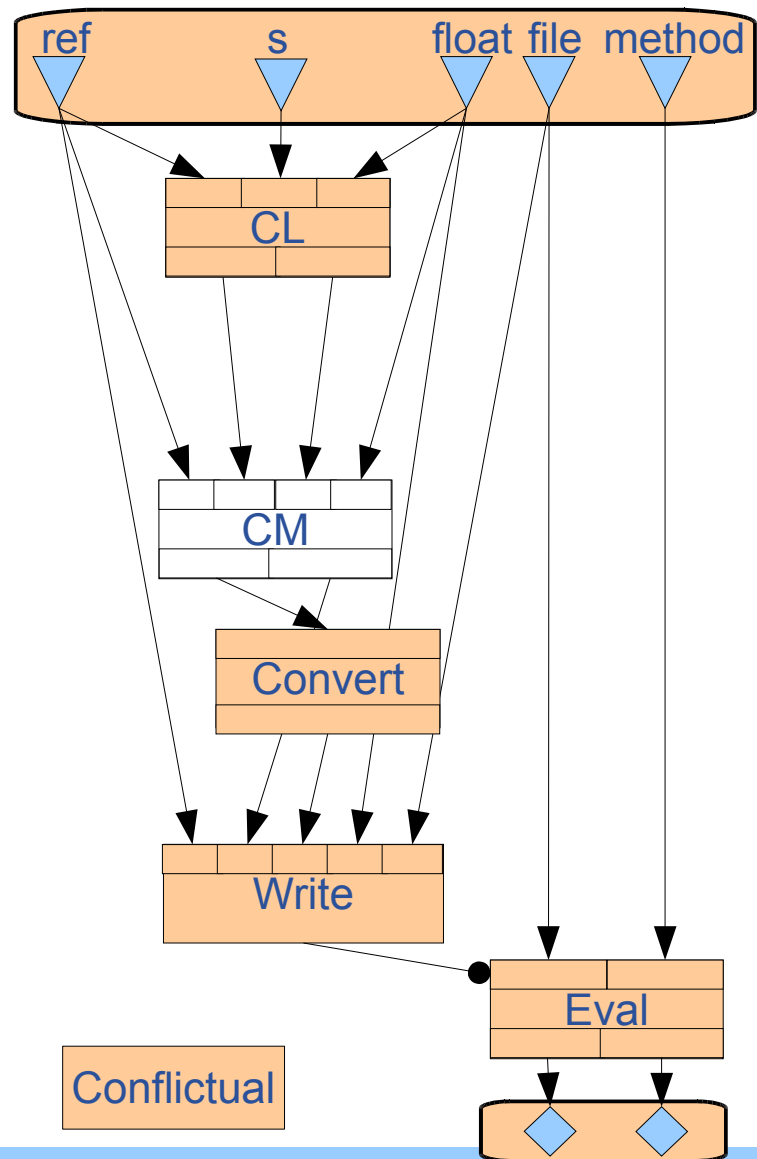
```
instr(i1,[im0,im1,cl],clOut,invoke(cl)).
instr(i2,[clOut],c1,assign(pfmC1)).
instr(i3,[clOut],c2,assign(pfmC2)).
instr(i4,
    [pm,pmC1,pmC2,transfo],pfmOut,invoke(pfm)).
...
```

% Constraints

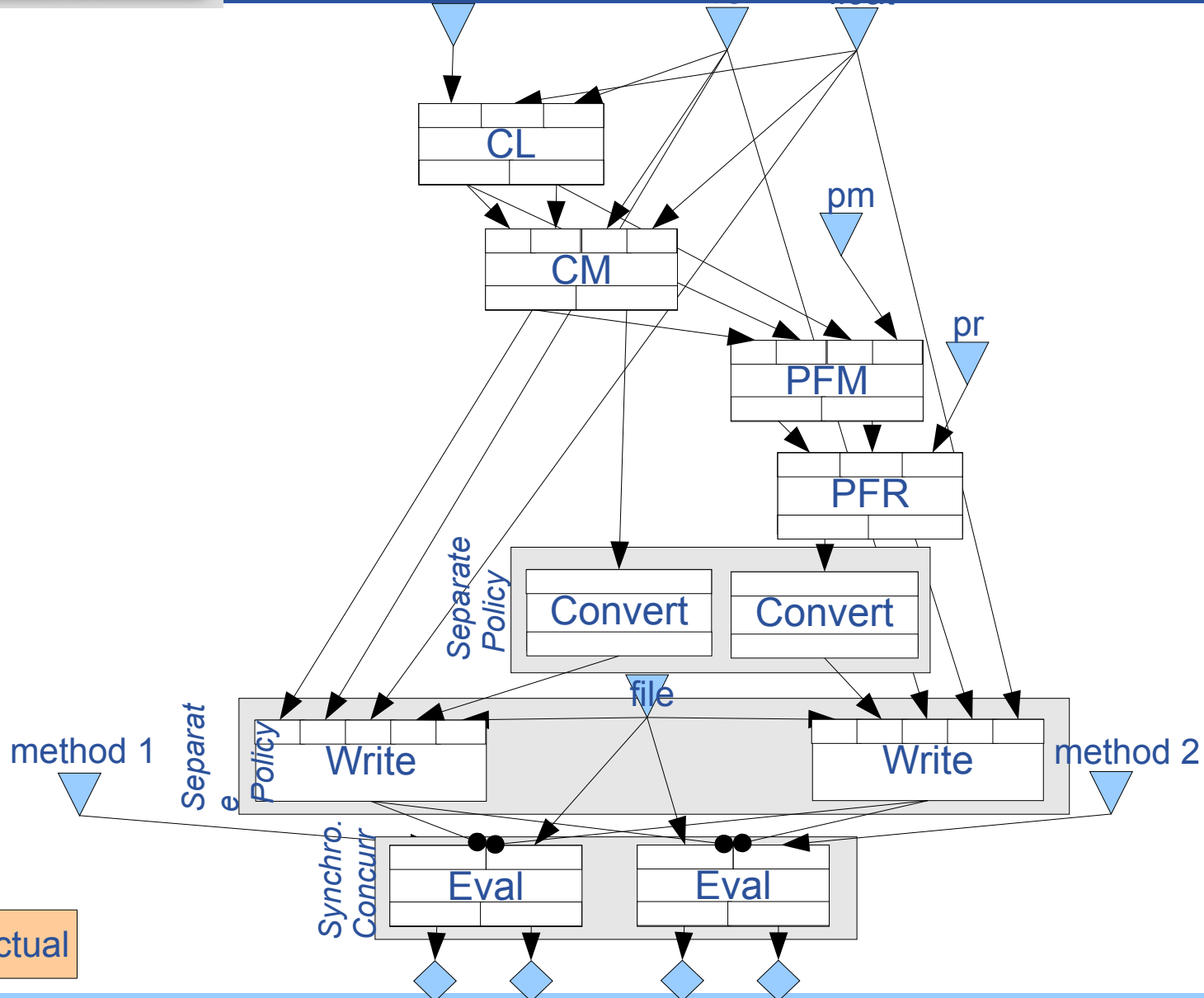
```
pred(i1,i2).
pred(i1, i3).
pred(i2, i4).
pred(i3, i4).
...
```

- **Consistency properties:**
 - P1: orchestrations have at most one invocation to a given basic service
 - Multiple invocations are wrapped in a complex invocation
 - P2: orchestrations have at most one input and one output
 - Multiple input / outputs are encapsulated in a structure
 - P3: orchestrations satisfy a set of constraints:
 - No concurrent write access to the variables
 - No cycle in instructions invocation order
- **Potential unification points are detected because they break either P1 or P2**
- **If an orchestration breaks P3, it is rejected from the process**

Conflicts detected







Conflictual



- **Scientific workflows**
 - Implicit parallelism
 - Complex data flows
- **Software composition**
 - Code reuse
 - Heterogeneous codes assembling
- **Preoccupations**
 - Performance
 - Ease of use
 - Separate application logic and data sets
 - Compact representation of complex processes