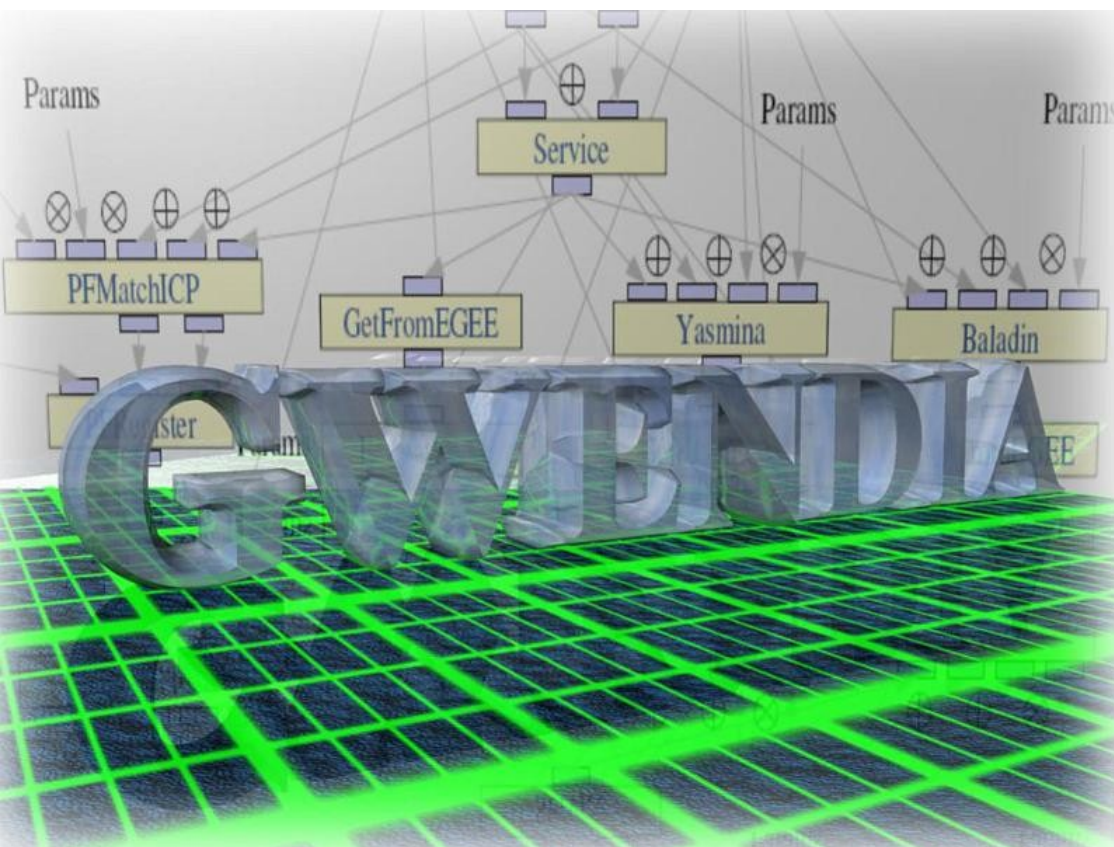




Grid Workflow Efficient Enactment for Data Intensive Applications

Grid Workflows

Sophia Antipolis, March 23rd 2010



University of Winsconsin

*University of Southern
California*

*Ecole Normale Supérieure de
Lyon*

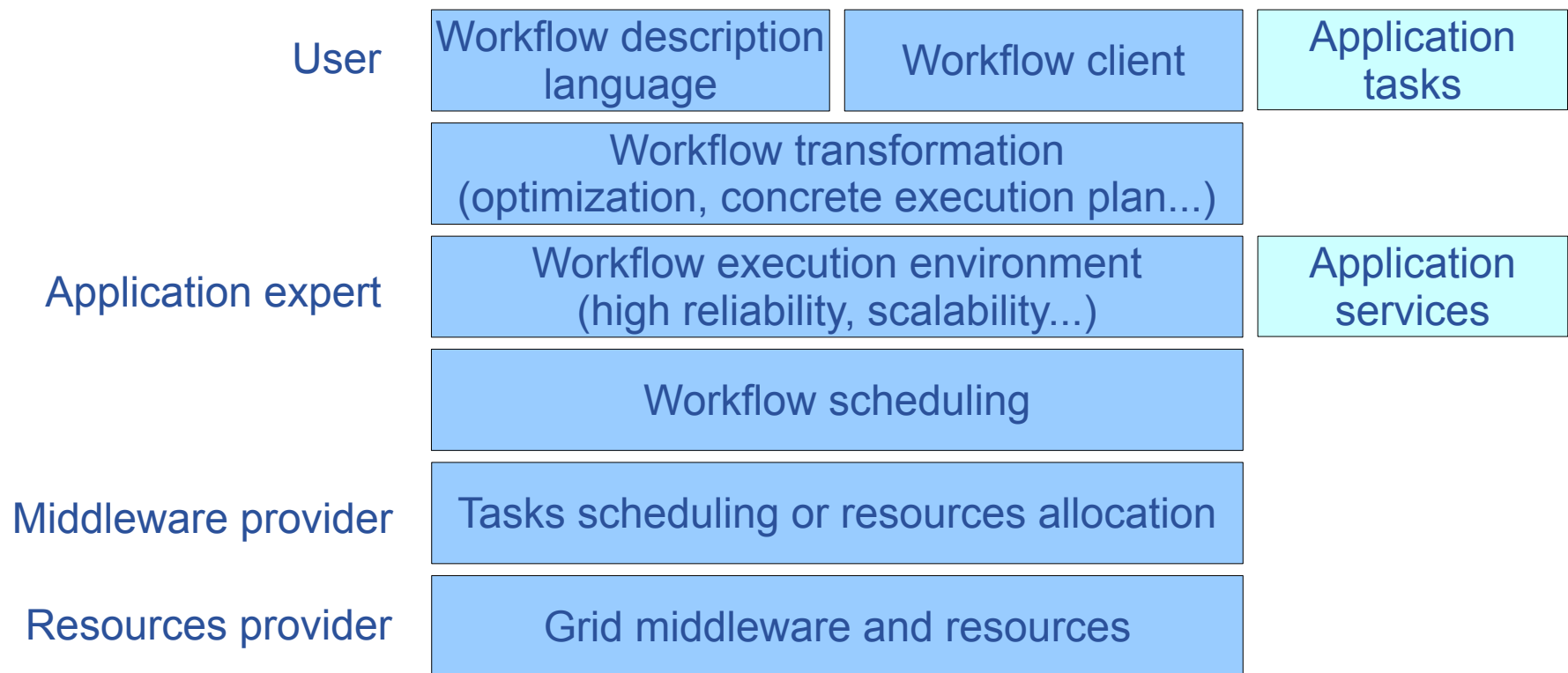
CNRS, I3S, MODALIS



Financé par
ANR



- **Grid, HTC, Data Parallelism, Workflows**
 - Grid infrastructure support for HTC
 - Embarassingly parallel (data parallel) applications
 - Coarse-grained workflow representation
- **Grid workflow enactment stack**

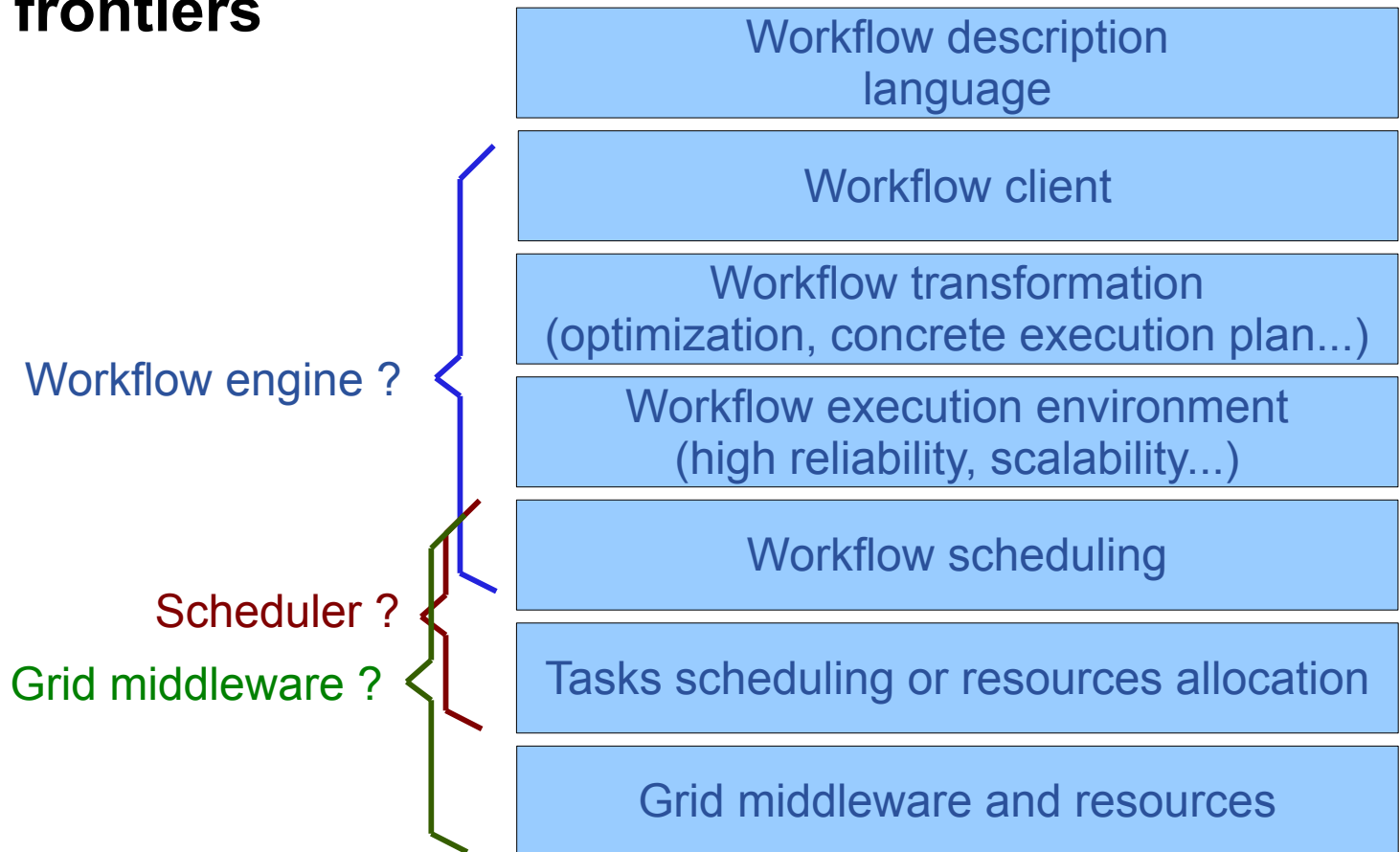




Grid workflow enactment stack

Grid Workflow Efficient Enactment for Data Intensive Applications

- **Blurred frontiers**



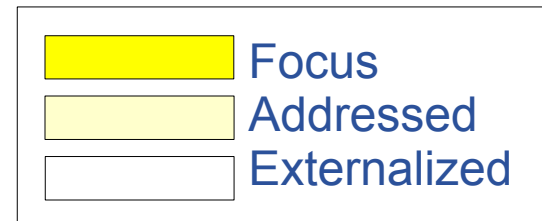
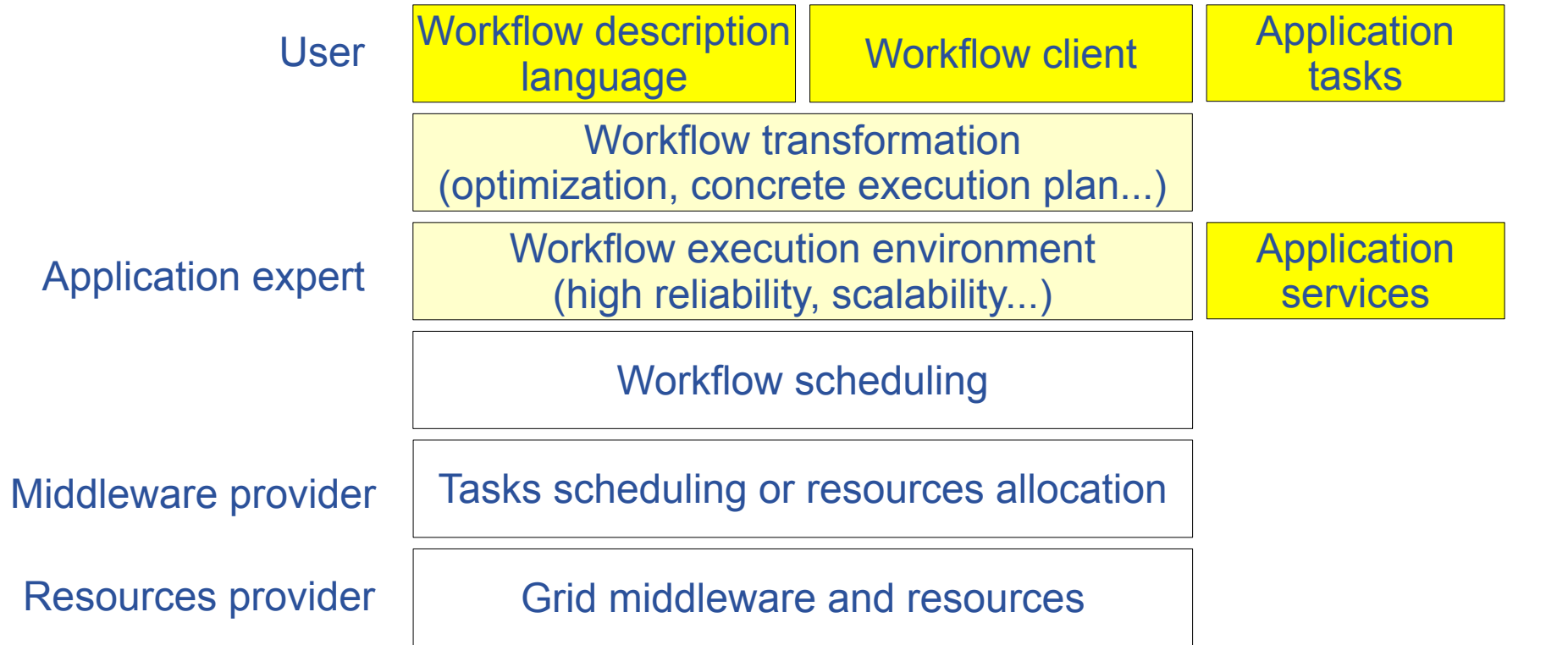
- **Tightly coupled stages**



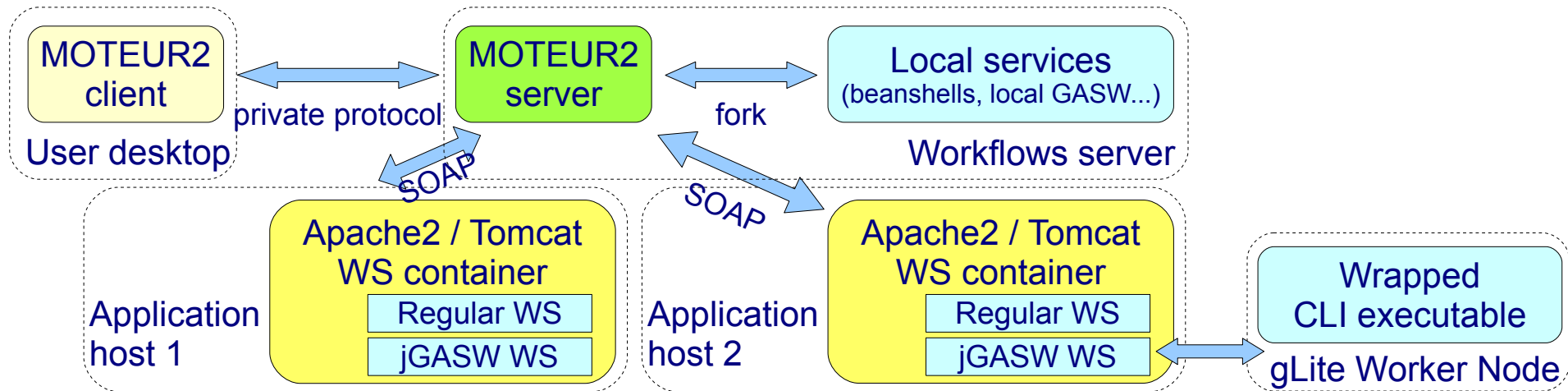
MOTEUR Workflow Engine

Grid Workflow Efficient Enactment for Data Intensive Applications

- Focus on application description**



- **Typical deployment**

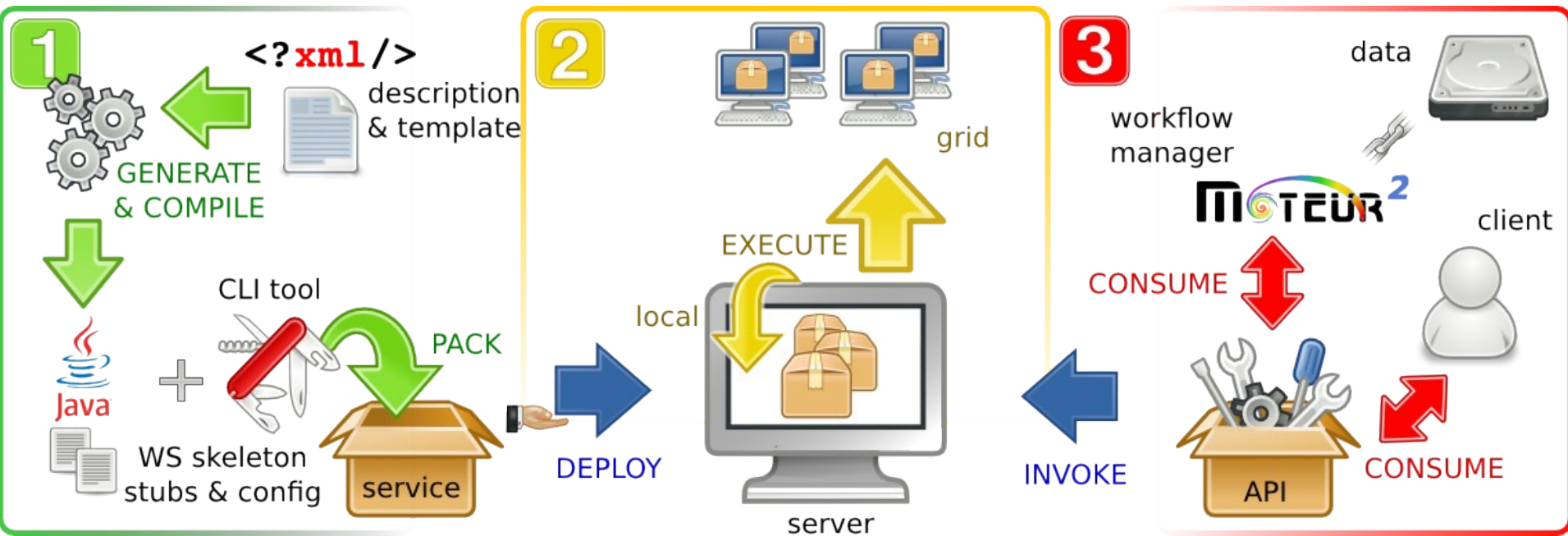


- **gLite middleware interface**

- CLI wrapped as Web Services
- WS invocation cause task to be relocated on a grid worker node

- **jGASW**

- Code packaging + invocation interface
- Based on a CLI interface descriptor

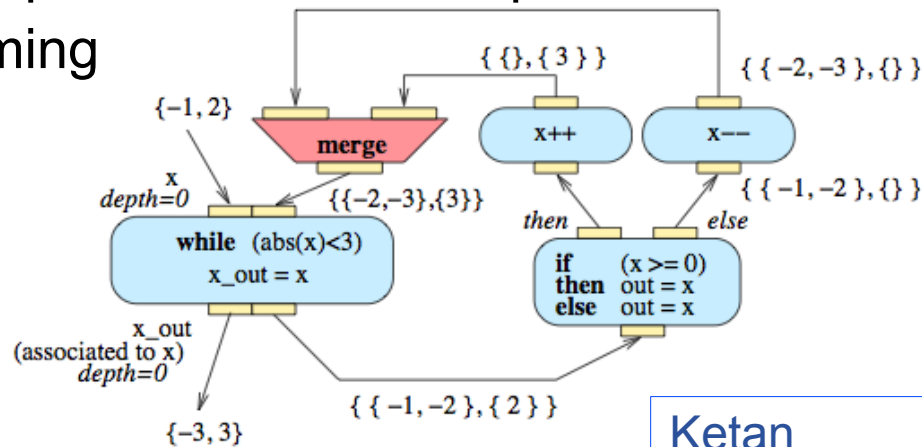


- **Tasks in a Service-Oriented Architecture**

- Service-Oriented framework; standard Web Service interface
- Execution unit relocation; embedded legacy CLIs

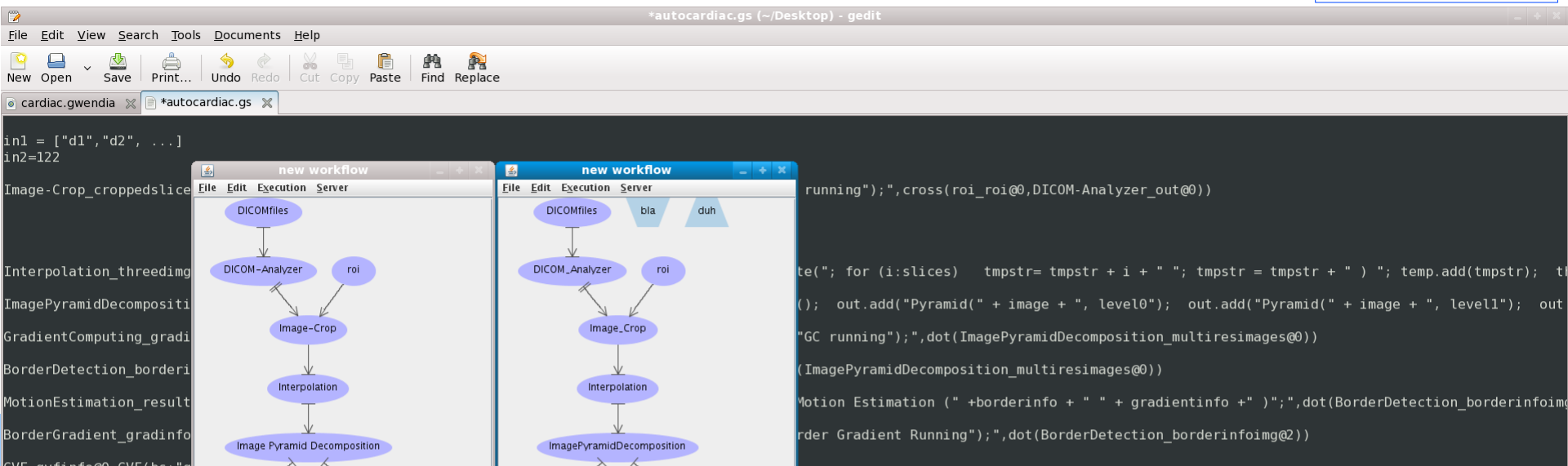
- **GWENDIA language**

- Use array-programming concepts to describe complex data flows
- Graph-based visual programming



Ketan Maheshwari

- **Dual visual-programming / scripting approach**





- **Objectives: highly-expressive, user-friendly workflow language for grids**
- **Motivations**
 - Little emphasis on expressiveness of workflow languages nor equivalence between different languages
 - Study how “compatible” different programming models are
- **We strongly believe there is no ideal language anyway**
 - Many different users, use cases and levels of expertise
 - User-friendliness is not an optimizable criterion
- **In practice, many other consideration will prevail**
 - Interface to grid infrastructure
 - Reliability
 - Performance
 - GUI...

These are environment and engines-related considerations. Different implementations could exist.

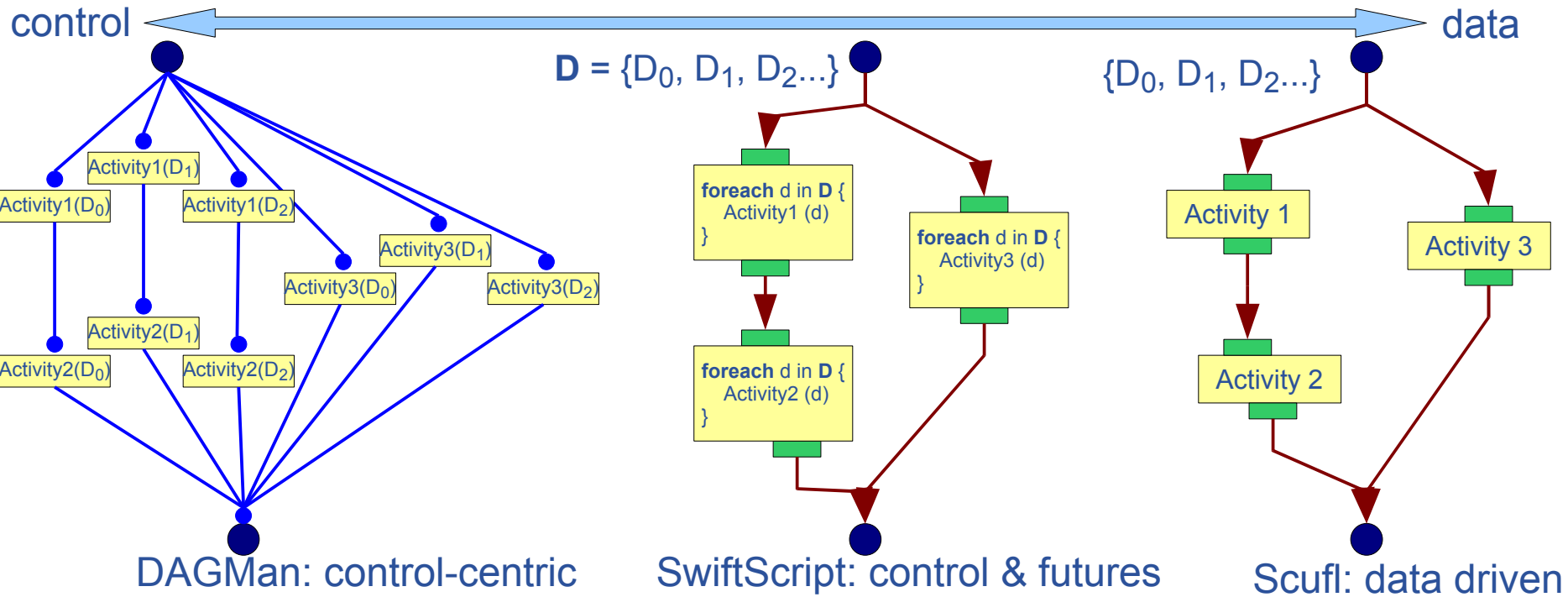


- **Any programming language could be adopted**
 - From C to Makefiles
- **But all do not fit well a distributed and data-intensive environment**
 - Grid workflows target human-scale activities
 - Prototyping is the rule
 - Scientific applications with heavy user communities and large data sets require high performance and/or high throughput
- **Workflow language for grids**
 - Coarse grain
 - Data intensive
 - Heavy legacy code
 - Interfaced to external Job / Workload manager(s)



- **Language ↔ Expressiveness**
 - Separate design and enactment phase
 - High expressiveness is not enough: easy access to high expressiveness is a real challenge
- **Not so much emphasis on languages**
 - Workflow environments often use visual programming
 - Internal representations not always clearly specified and documented
 - This impacts the semantics of applications: semantic might change over time or implemented using undocumented, specific assumptions that make understanding of programs difficult in non-trivial cases.

- **Many dependencies are data dependencies**



- **Most languages are data-driven in the end**

- Scuff: by construction
- MA DAG: by generation
- SwiftScript: using futures
- Askalon: mixed data/control flows with coherent, fixed patterns

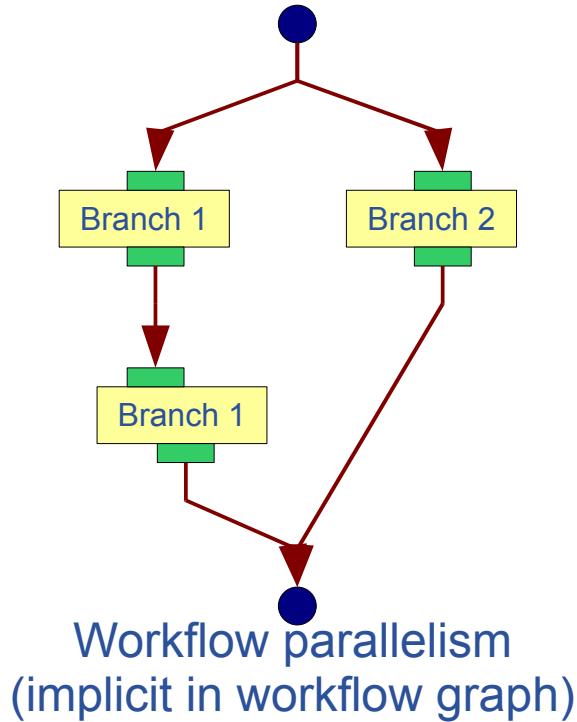


Data-driven fits data-intensive well

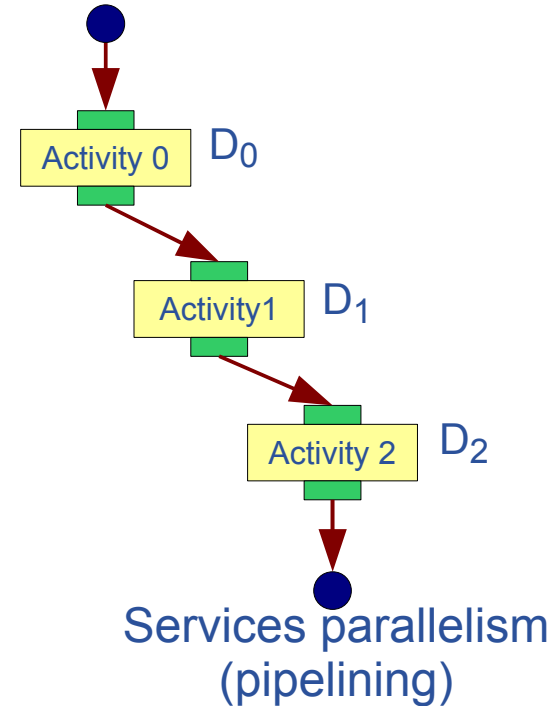
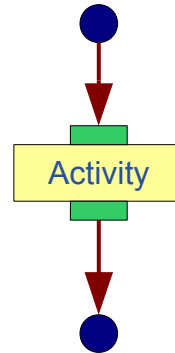
Grid Workflow Efficient Enactment for Data Intensive Applications

- **Main source of parallelism comes from data sets**
 - Embarrassingly parallel
 - Parameter sweep
 - SPMD
- **Clear separation of application logic and experiment**
 - Experiment = raw data and parameter data
 - This is a clear difference with BPEL or other imperative programming inspired approaches
- **Implicit description of parallelism**

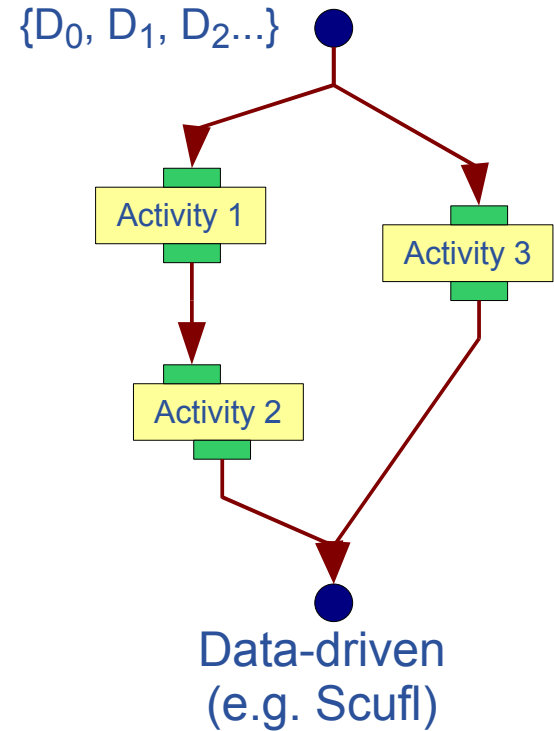
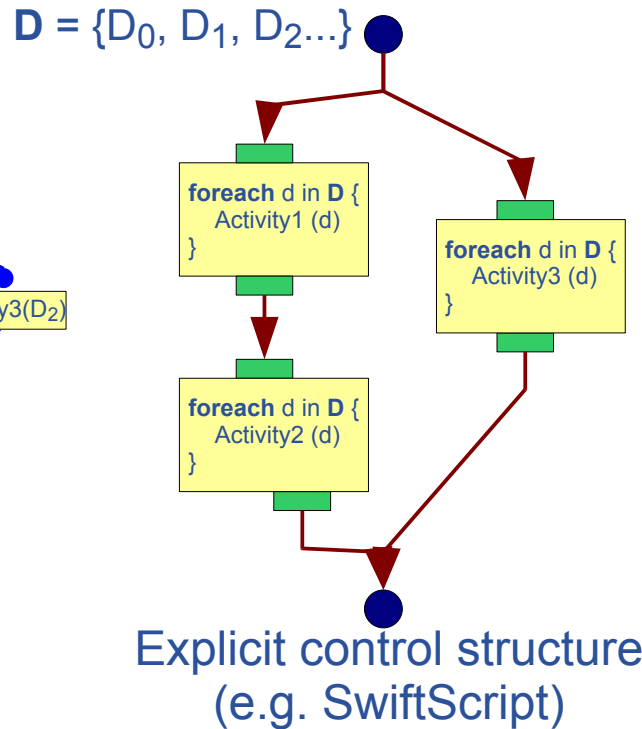
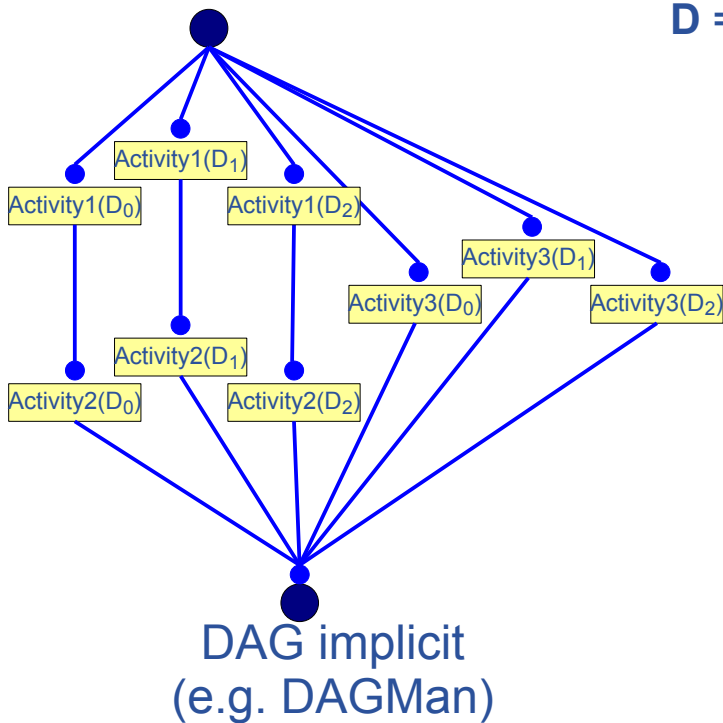
- **Three levels of parallelism**



$\{D_0, D_1, D_2 \dots\}$



- **The expression of data and service parallelisms depend on the workflow language considered**
- **Data parallelism is implicit in data-driven approaches**



- **DAGs: implicit in the (large scale) workflow graph**
- **Parallel languages: explicit control structures**
- **Data-driven: transparent**



Parallel languages vs data-driven flows

Grid Workflow Efficient Enactment for Data Intensive Applications

- **Parallel languages**

- Bounded: `foreach d in {D0, D1, D2}...`
- Unbounded: `foreach d in D...`
- Involves data synchronization at each data parallel service (no pipelining)

- **Data-driven flows (e.g. ScufI)**

- Independent definition of data flows and data sets
- Enable completely asynchronous enactment (data parallelism + pipelining): no data synchronization
- May require explicit data synchronization barriers when this is needed

- **Futures (e.g. SwiftScript)**

- Non-blocking assignment operations, blocking read
- Data-centric approach, asynchronous execution



Control structures or not?

Grid Workflow Efficient Enactment for Data Intensive Applications

- **Many data-driven / visual programming approaches do not define control structures**
 - Hardly any, yet some (fail-if-false, fail-if-true)
- **DAGs do not include any loop**
 - Because complex application logic is described inside the workflow activities
- **Yet control on data flows is sometimes needed**
 - Different execution conditions
 - Exceptions / Retry on errors at the application level
 - ...



Gwendia data-driven workflow languages

Grid Workflow Efficient Enactment for Data Intensive Applications

- **Search for a language integrating**
 - Data-driven approach (transparent parallelism)
 - Data arrays manipulation
 - Control structures
 - Maximum asynchronous execution
 - Compact representation
- **Applying array programming principles**
 - Arrays are first class entities
 - Operators apply both to scalar and arrays
- **Indexing scheme to preserve computation coherence in an asynchronous environment**
- **Arrays nesting to represent data-sets decomposition / recomposition**



- **Array programming principles**

- Lightweight syntax for handling arrays:

```
X+Y ≡ foreach i in indices(X) do
    Xi + Yi
done
```

- Convenient representation for vectorial processors
- Extendable to any operation on arrays of values
- Use nested arrays $x = \{\{1, 2\}, \{-1, -2\}\}$ is a 2-nesting levels array

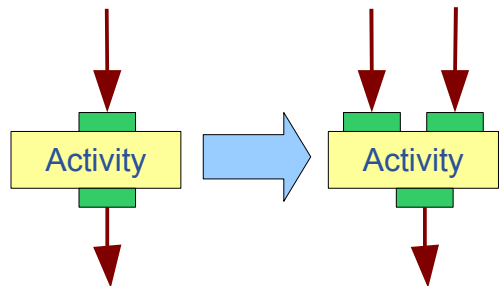
- **Activities with array parameters**

- $A(X) \equiv \{ A(X_0), \dots, A(X_n) \}$
- $\{\{1, 2\}, \{-1, -2\}\} \equiv \{ \{A(1), A(2)\}, \{A(-1), A(-2)\} \}$

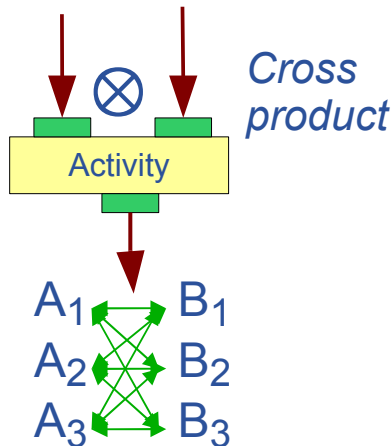
- **Ports depth**

- Activity ports have a defined depth that corresponds to the number of array nesting levels being synchronized:
- $\text{Mean}(\{1, 2, 3\}) = 2$ if `Mean` has input port depth 1

Iteration strategies

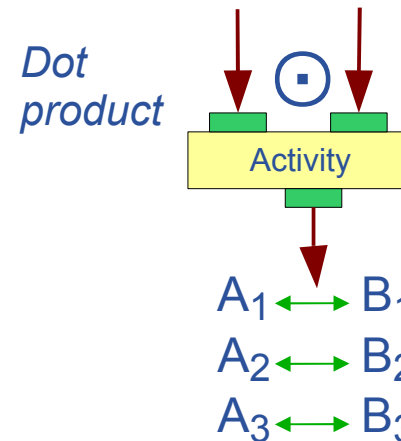


$\{A_1, A_2, A_3\}$ $\{B_1, B_2, B_3\}$



$$\mathbf{A} \otimes \mathbf{B}$$

$\{A_1, A_2, A_3\}$ $\{B_1, B_2, B_3\}$



$$\mathbf{A} \odot \mathbf{B}$$

- In Scufi

- Parallel language

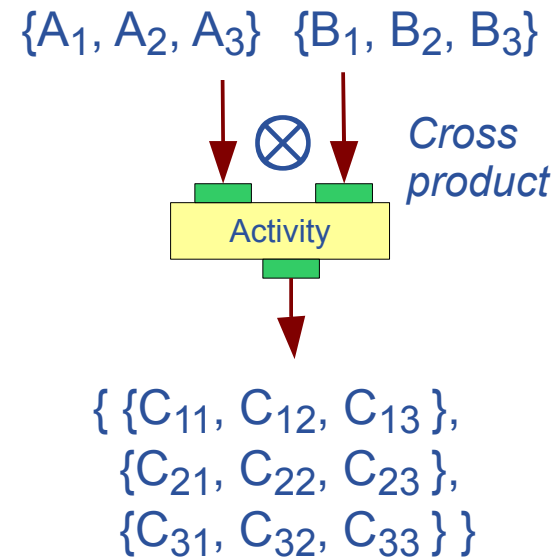
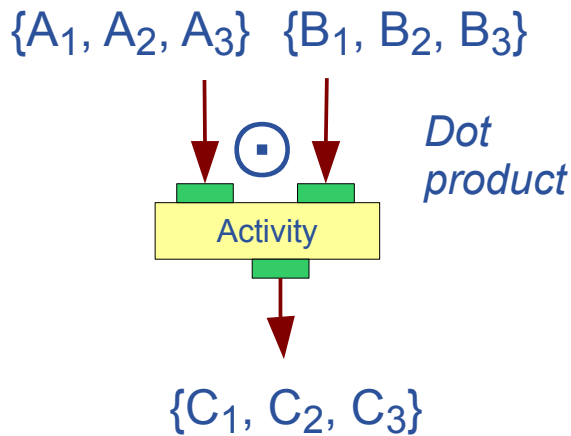
```

foreach a in A
  foreach b in B
    fire Activity(a,b)
  
```

```

foreach i in indices(A)
  fire Activity(Ai, Bi)
  
```

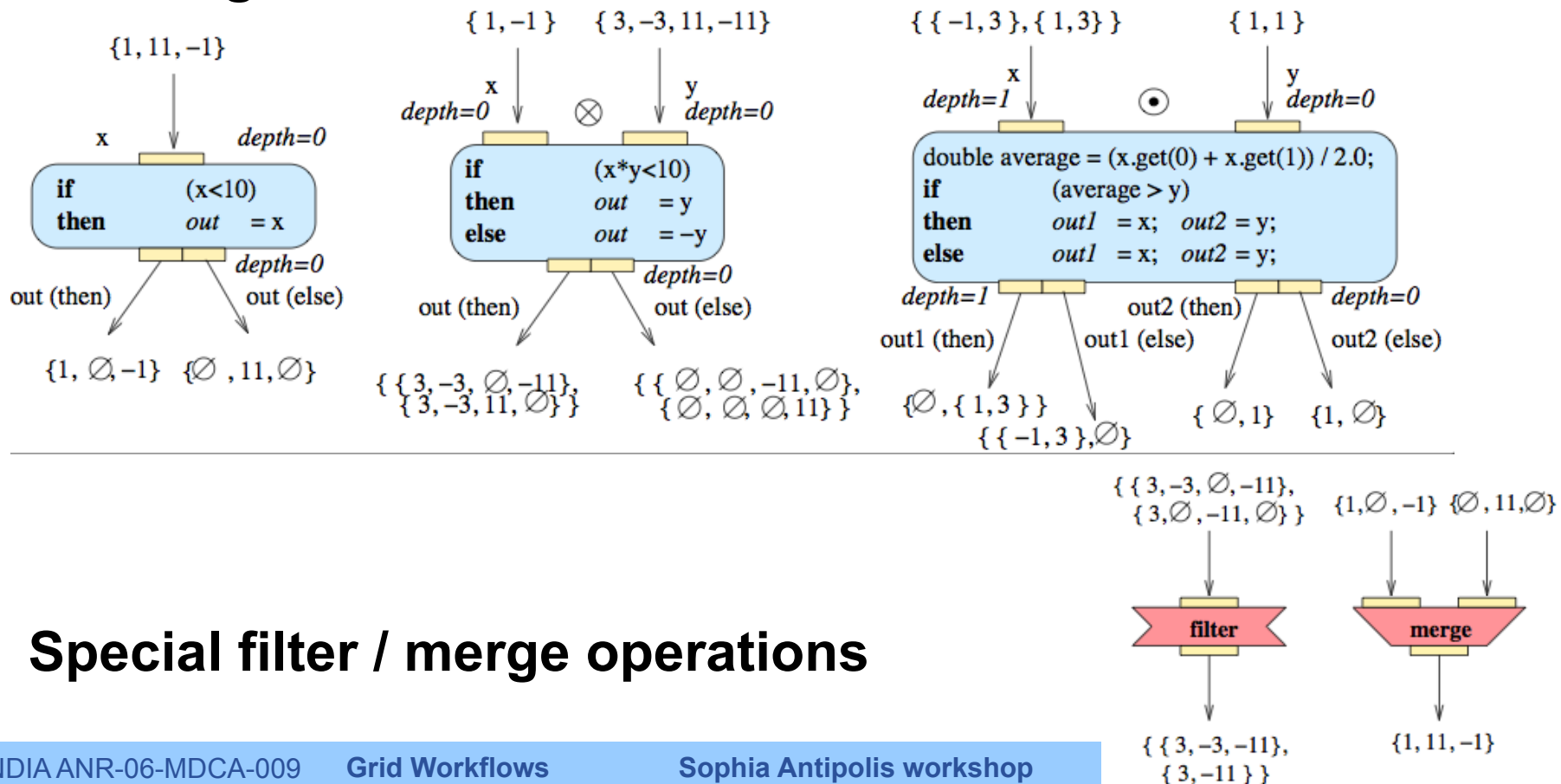
- Iterations strategies define output indexing strategies
- They may change data nesting levels





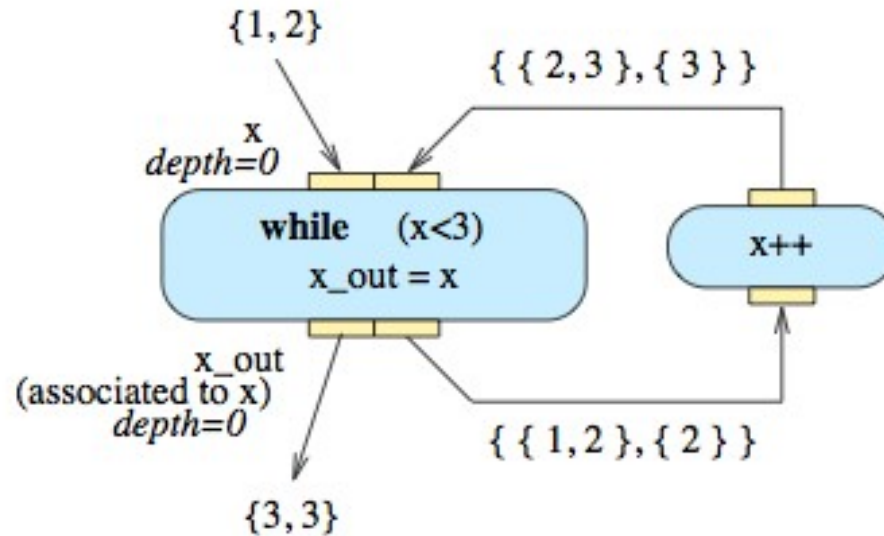
- **Special *void* value \emptyset**
- **$\mathbf{A}(\emptyset) \equiv \emptyset$**
 - No evaluation of the activity (\emptyset is only known from the workflow engine)
- **Iteration strategies semantics**
 - $\mathbf{x} \odot \emptyset \equiv \emptyset \odot \mathbf{x} \equiv \emptyset$
 - $\mathbf{x} \otimes \emptyset \equiv \emptyset \otimes \mathbf{x} \equiv \emptyset$
- **\emptyset has an index in the array it belongs to**

- Test expression variables mapped to input ports
- Dual “then / else” output ports
- Use of the special “void” data element to preserve indexing scheme coherence

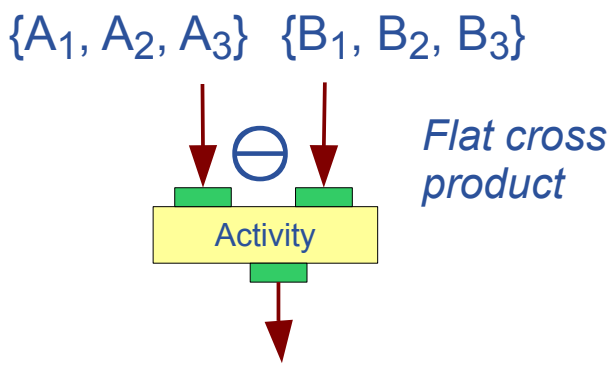


- Special filter / merge operations

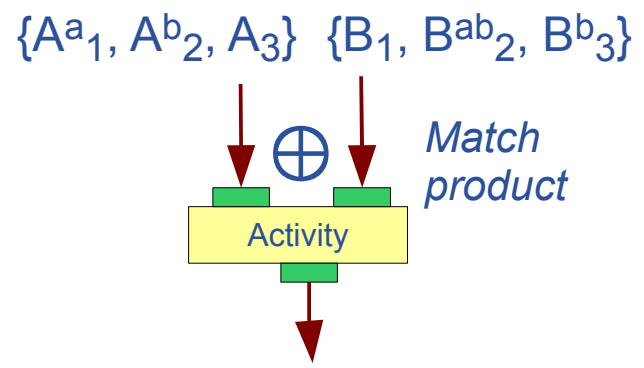
- **Test expression variables mapped to input ports**
- **Dual “inner / outer” ports**
 - Input ports: initialization / iterated values
 - Output ports: iterated values / end of loop



- Flat-cross product
- Match product



$\{C_{11}, C_{12}, C_{13}, C_{21}, C_{22}, C_{23}, C_{31}, C_{32}, C_{33}\}$



$\{ \{\emptyset, C_{12}, \emptyset\}, \{\emptyset, C_{22}, C_{23}\}, \{\emptyset, \emptyset, \emptyset\} \}$



- **Control structures**
 - Forseeable: data flow generated statically from inputs (dot, cross, flat-cross and for loops)
 - Unforseeable: dynamic generation (match, conditionals, while loop, output array with unknown size)
- **Complete asynchronous execution**
 - With best effort to deal with unknown indexes (e.g. flat cross)
- **Moteur2 workflow engine**
<http://modalis.polytech.unice.fr/moteur/start>
 - On-the-fly interpreter
 - All structures processed dynamically
- **MA-DAG engine** <http://www.ens-lyon.fr/graal/~diet>
 - Partial DAGs generation, greedy approach
 - Multi-DAGs scheduling heuristics



- **Coherent integration of:**
 - Data-driven approach: Scufi-based
 - Array programming principles: Nested arrays as first-class entities
 - Asynchronous execution: coherence through data items indexing
 - Control structures
- **Transparent parallelism**
 - No *foreach* nor *fork/join* constructs
- **Compact**
 - Only application logic described, no redundancy
 - No *for* kind of loop over data sets
 - Iteration strategies
- **Reusability**
 - Decoupled data and application logic