

Distilling structure in Taverna scientific workflows: a refactoring approach – Additional file

Sarah Cohen-Boulakia^{*1,2}, Jiuqiang Chen^{1,2,3}, Paolo Missier⁴, Carole Goble⁵, Alan R Williams⁵
and Christine Froidevaux^{1,2}

¹Laboratoire de Recherche en Informatique, CNRS UMR 8623, Université Paris Sud, France

²AMIB group, INRIA Saclay, France

³School of Information Science and Engineering, Lanzhou University, Lanzhou, Gansu, China

⁴University of Newcastle, United Kingdom

⁵University of Manchester, United Kingdom

Email: cohen@lri.fr; chenj@lri.fr; paolo.missier@ncl.ac.uk; carole.goble@manchester.ac.uk; alanrw@cs.man.ac.uk;
christine.froidevaux@lri.fr;

*Corresponding author

Abstract

This document provides the additional file for the paper “Distilling structure in Taverna scientific workflows: A refactoring approach”. It describes in details the algorithms designed in the refactoring approach introduced.

Pseudo code of the DownDistillation and Distill procedures

```
1 DownDistillation(IN  $GG[q, v]$ , IN/OUT  $DSGG$ : graphs, IN  $q$ : node,  
2                 IN/OUT  $SetAU$ : set of graphs, IN/OUT  $ListRed$ : set of nodes)  
3 Distill( $GG[q, v]$ ,  $DSGG$ ,  $q$ );  
4  $ListRed \leftarrow ListRed \cup \{\text{new reduction nodes of } GG[q, v]\}$ ;  
5  $SetAU \leftarrow SetAU \cup \{\text{new autonomous subgraphs of } GG[q, v]\}$ ;  
6 foreach autonomous subgraph  $GG[a, b]$  in  $SetAU$  do  
7   | Distill( $GG[a, b]$ ,  $DSGG$ ,  $a$ )  
8 end  
9 End DownDistillation
```

Figure 1: Pseudo-code of the DownDistillation procedure

Additional functions used are introduced here after.

- The following procedure makes use of the function **OKTransformation**(p, q, GG) which specifies the conditions for p and q to be merged. It is true iff the following conditions are satisfied:
 - p and q are copies of each other;
 - p and q are involved in some anti-pattern (A) or (B) in GG ;
 - for any autonomous subgraph G' of GG , every time p appears in G' , q appears in G' too. This last condition ensures us that we do not remove an anti-pattern by a transformation that would make an SP-graph becoming non-SP.
- The function **SameOrientedPath**(p, q, GG) is true iff there is at least a directed path dp in GG such that p and q belong to dp .
- **Visited** is a function allowing to mark nodes as visited or unvisited.

Note that, for example, in Figure 5 of the main paper, $OKTransformation(9, 11, GG)$ is false.

The next page provides the complete pseudo-code of the **Distill** procedure.

```

1 Distill(IN GG: graph; IN/OUT DSGG: graph; IN x: node)
2 v ← sink(GG);
3 ListRed ← set of reduction nodes of GG;
4 SetAU ← set of autonomous subgraphs of GG;
5 Visited(GG[q, v]) ← false /* set all the nodes of GG[q, v] unvisited */
6 foreach successor p of x in GG do
  /* search for copies of p */
7   if Visited(p) ← false then
8     p1 ← p; flagp ← true;
9     while flagp do
10      /* flagp allows to consider all the unvisited descendant of p1 if necessary */
11      Distilled ← false /* Distilled says if some transformation on p1 has been done */
12      foreach successor q of x in GG, such that q ≠ p1 do
13        /* successors of x different from p1 are potentially copies of p1 */
14        q1 ← q; flagq ← true;
15        while flagq do
16          /* flagq allows to consider all the unvisited descendant of q1 if necessary */
17          if Visited(q1)=false and SameOrientedPath(p1,q1, GG)=false then
18            if OKTransformation(p1,q1,GG)=true then
19              /* q1 is a copy of p1 in some anti-pattern and transformation can be performed */
20              transformation on DSGG, replacing q1 by mergeq;
21              flagq ← false; distilled ← true; /* loop on q is stopped */
22            else
23              /* no transformation has been done on p1 and q1 */
24              if outDegree(q1) ≠ 1 then
25                if there exists a single autonomous subgraph GG[q1,y] in SetAU then
26                  q1 ← y; /* the loop on q1 is continued with the sink of the unique autonomous
27                  subgraph */
28                else
29                  /* there is no autonomous subgraph GG[q1,y] in SetAU or more than one */
30                  if q1 is a reduction node in Listred then
31                    /* search for anti-patterns from reduction node q1 */
32                    DownDistillation(GG[q1,v], DSGG, q1, SetAU, ListRed);
33                    Visited(GG[q1,v]) ← false;
34                    if outdegree(q1) > 1 then flagq ← false;
35                  else
36                    flagq ← false;
37                    /* q1 is not a reduction node or there is no autonomous subgraph GG[q1,y] in
38                    SetAU the loop on q is stopped */
39                  end
40                end
41                q1 ← the successor of q1 /* outDegree(q1) = 1 */
42            end
43          end
44          flagq ← false;
45        end
46      end
47    end
48  end
49  /* while loop on flagp continues*/
50  if distilled then
51    /* if p1 has been merged with some other node then search for anti-patterns from p1 */
52    DownDistillation(GG[p1,v], DSGG, p1, SetAU, ListRed);
53    Visited(x,mergeq) ← true; /* set all the nodes on all paths from x to mergeq as visited */
54  else
55    /* p1 has not been merged */
56    if outDegree(p1) ≠ 1 then
57      if there exists a single autonomous GG[p1,y] in SetAU then
58        p1 ← y;
59      else
60        flagp ← false;
61      end
62    else
63      let p1 ← the successor of p1 /* outDegree(p1) = 1 */
64    end
65  end
66 end
67 end
68 EndDistill

```

Figure 2: Pseudo-code of the Distill procedure