*Janus*:
Fine-grained and efficient
provenance querying for Taverna
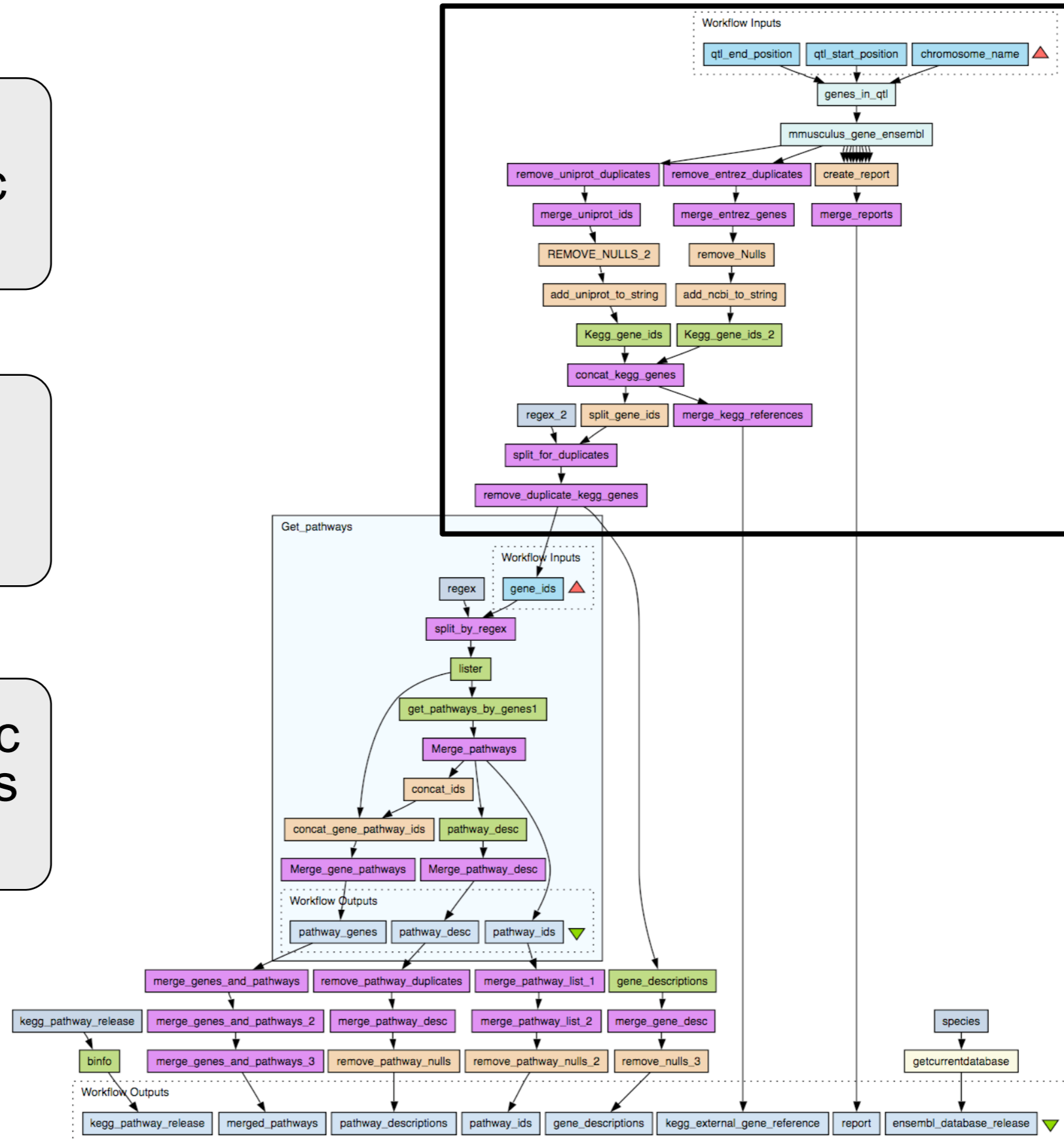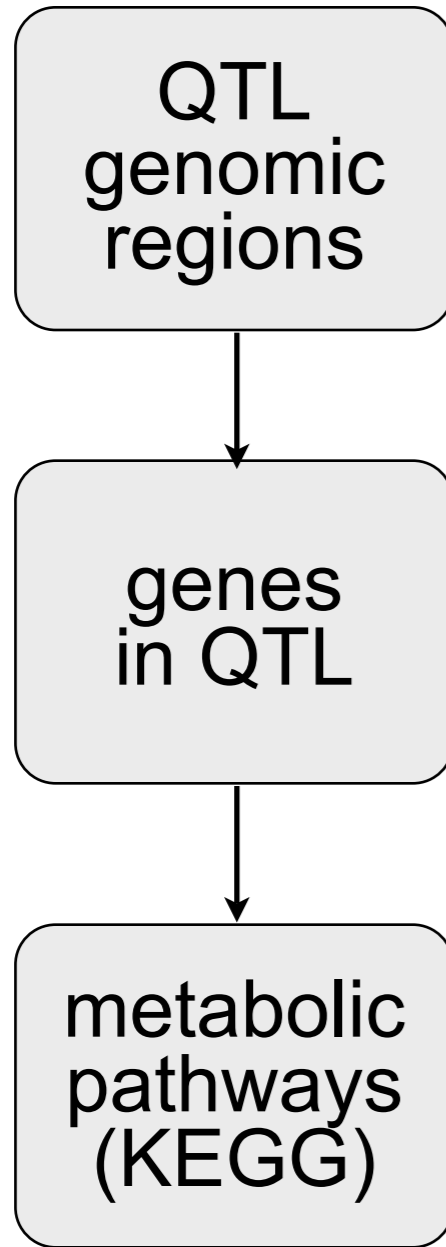
Paolo Missier

Information Management Group

School of Computer Science, University of Manchester, UK
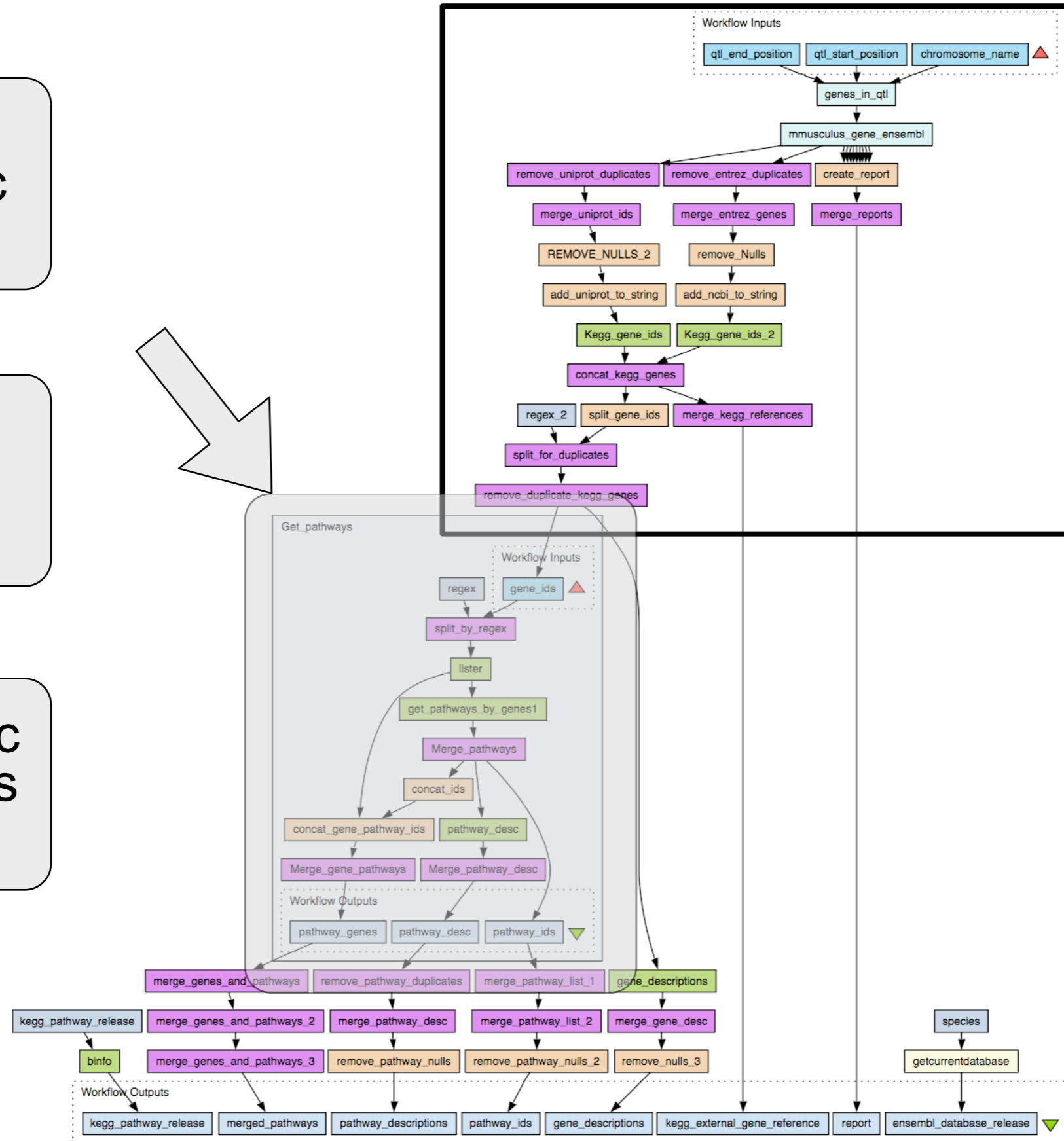

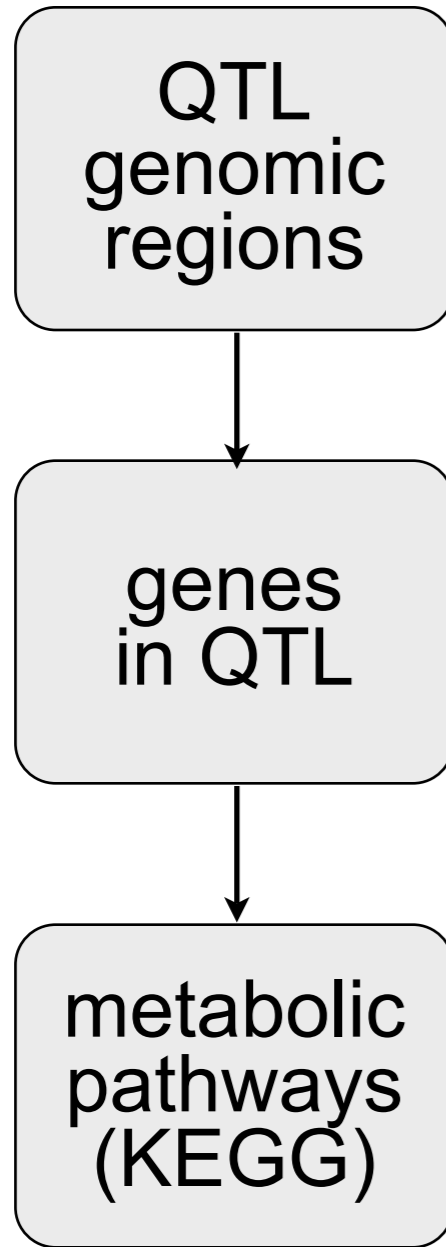in collaboration with:

Norman Paton, Khalid Belhajjame

1

- Motivation: Taverna for rapid information integration

- Fine-grained process provenance and its role in Taverna

- Context and scope: forms and uses for process provenance

- Technical challenges in querying provenance traces

- A solution, and experimental results

# Workflow as data integrator

QTL genomic regions

genes in QTL

metabolic pathways (KEGG)

June 2009 - P. Missier

# Workflow as data integrator

June 2009 - P. Missier

KEGG gene ids:
"mmu:20816   (g1)
mmu:26416   (g2)
mmu:328788 (g3)"

[ p1, p2, p3, ...]

"p1 MAPK signaling pathway
 p2 VEGF signaling pathway
…  "

4

KEGG gene ids:
"mmu:20816   (g1)
mmu:26416   (g2)
mmu:328788 (g3)"

[ p1, p2, p3, ...]

KEGG DB lookups

"p1 MAPK signaling pathway
 p2 VEGF signaling pathway
…  "

4

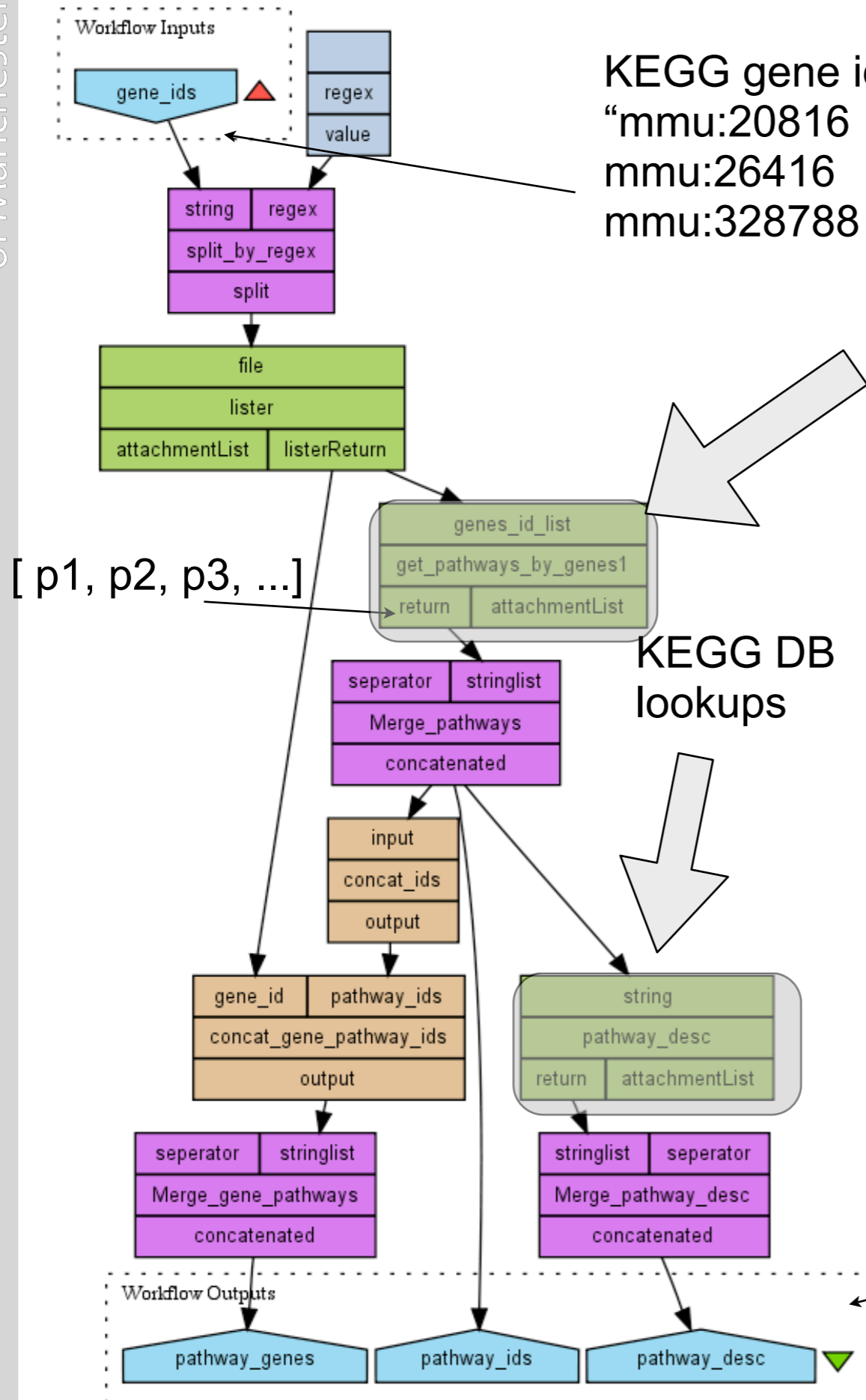KEGG gene ids:
"mmu:20816  (g1)
mmu:26416  (g2)
mmu:328788 (g3)"

[ p1, p2, p3, ...]

KEGG DB
lookups

"p1 MAPK signaling pathway
 p2 VEGF signaling pathway
… "

4

KEGG gene ids:
"mmu:20816   (g1)
mmu:26416   (g2)
mmu:328788 (g3)"

[ p1, p2, p3, ...]

KEGG DB lookups

goal:
- list all pathways that are mapped from both sets of genes

- substantial list manipulation involved in achieving this
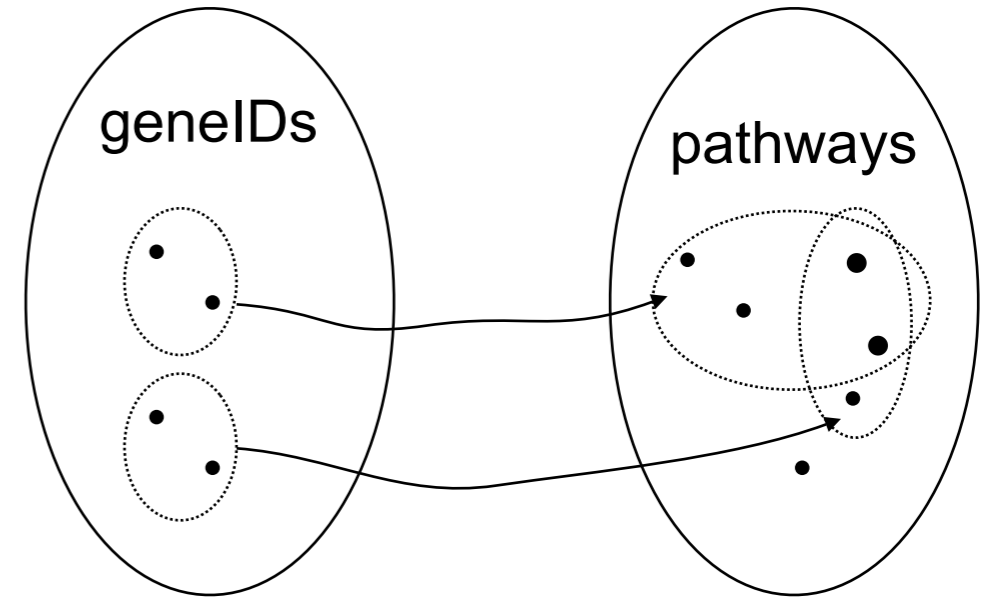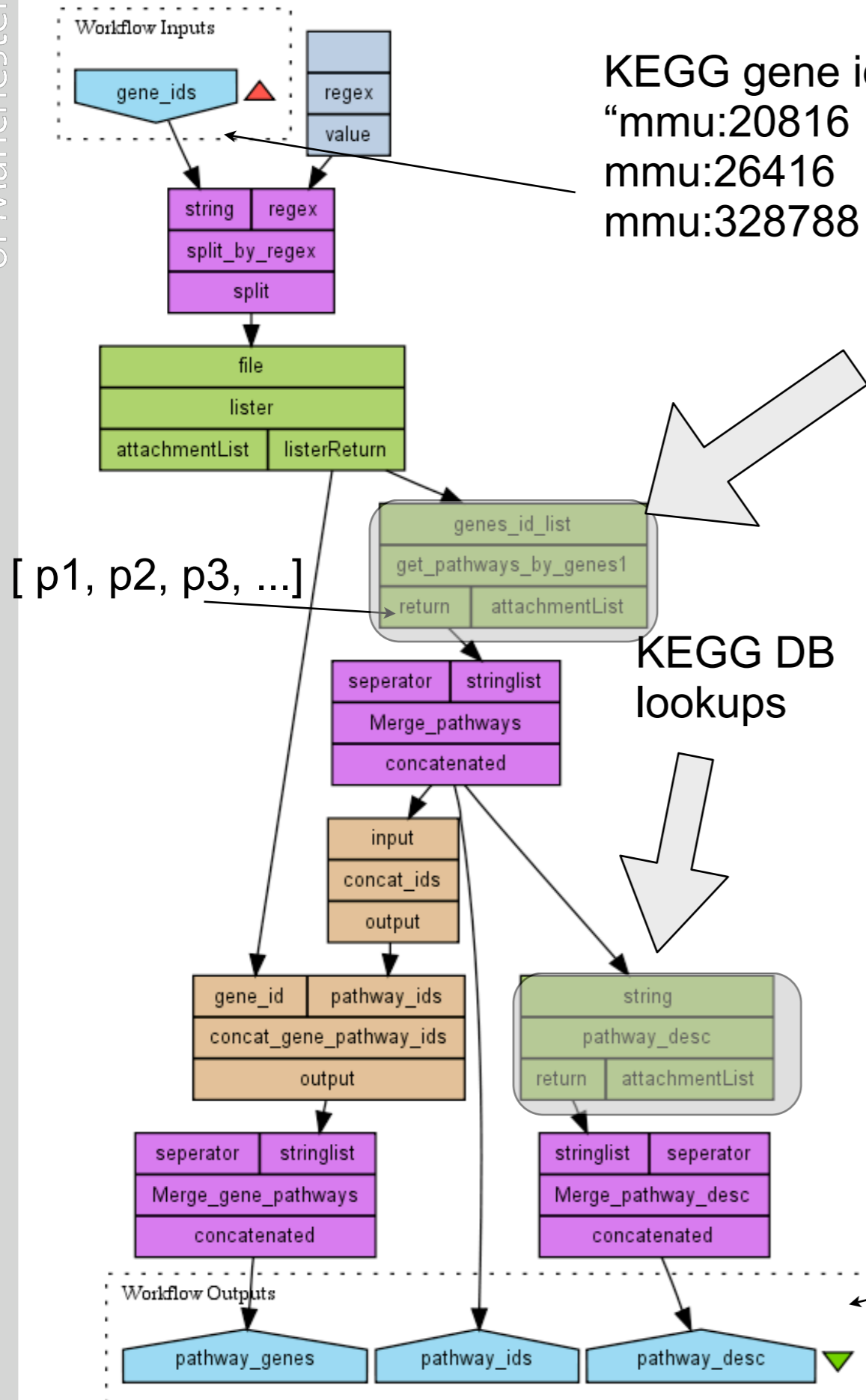
"p1 MAPK signaling pathway
 p2 VEGF signaling pathway
 ... "

4

KEGG gene ids:
"mmu:20816   (g1)
mmu:26416   (g2)
mmu:328788 (g3)"

[ p1, p2, p3, ...]

KEGG DB
lookups

goal:
- list all pathways that are mapped from both
  sets of genes

- substantial list manipulation involved in
achieving this

[ [ g1, g2, g3],
  [ "g1 p1 p2 ...", "g2 p1 p2 ..." ,"g3 p1 p2 ..."]
  [ g1, g2, g3] ]

"p1 MAPK signaling pathway
 p2 VEGF signaling pathway
 ... "

4

List-structured
KEGG gene ids:

[ [ mmu:26416 ],
  [ mmu:328788 ] ]

[ path:mmu04010 MAPK signaling,
  path:mmu04370 VEGF signaling ]

[ [ path:mmu04210 Apoptosis, path:mmu04010 MAPK signaling, ...],
  [ path:mmu04010 MAPK signaling , path:mmu04620 Toll-like receptor, ...] ]

5

List-structured
KEGG gene ids:

[ [ mmu:26416 ],
  [ mmu:328788 ] ]

geneIDs

pathways

[ path:mmu04010 MAPK signaling,
  path:mmu04370 VEGF signaling ]

[ [ path:mmu04210 Apoptosis, path:mmu04010 MAPK signaling, ...],
  [ path:mmu04010 MAPK signaling , path:mmu04620 Toll-like receptor, ...] ]

5

List-structured
KEGG gene ids:

[ [ mmu:26416 ],
    [ mmu:328788 ] ]

geneIDs

pathways

geneIDs

pathways

[ path:mmu04010 MAPK signaling,
  path:mmu04370 VEGF signaling ]

[ [ path:mmu04210 Apoptosis, path:mmu04010 MAPK signaling, ...],
  [ path:mmu04010 MAPK signaling , path:mmu04620 Toll-like receptor, ...] ]

5

- Pros:
  - simpler to design and understand (hopefully)
  - (no shims!)
  - accepts multiple gene sets
    - returns list of pathways separately for each gene set
    - in addition to those shared by the union of all sets
- Cons:
  - no genes in output list:

    [ [ path:mmu04210 Apoptosis, path:mmu04010 MAPK signaling, ...],
      [ path:mmu04010 MAPK signaling , path:mmu04620 Toll-like receptor, ...] ]

  - so the relationship between the gene set and the pathway set is lost...

  ...Or is it?

6

List-structured
KEGG gene ids:

[ [ mmu:26416 ],
  [ mmu:328788 ] ]

[ path:mmu04010 MAPK signaling,
  path:mmu04370 VEGF signaling ]

[ [ path:mmu04210 Apoptosis, path:mmu04010 MAPK signaling, ...],
  [ path:mmu04010 MAPK signaling , path:mmu04620 Toll-like receptor, ...] ]
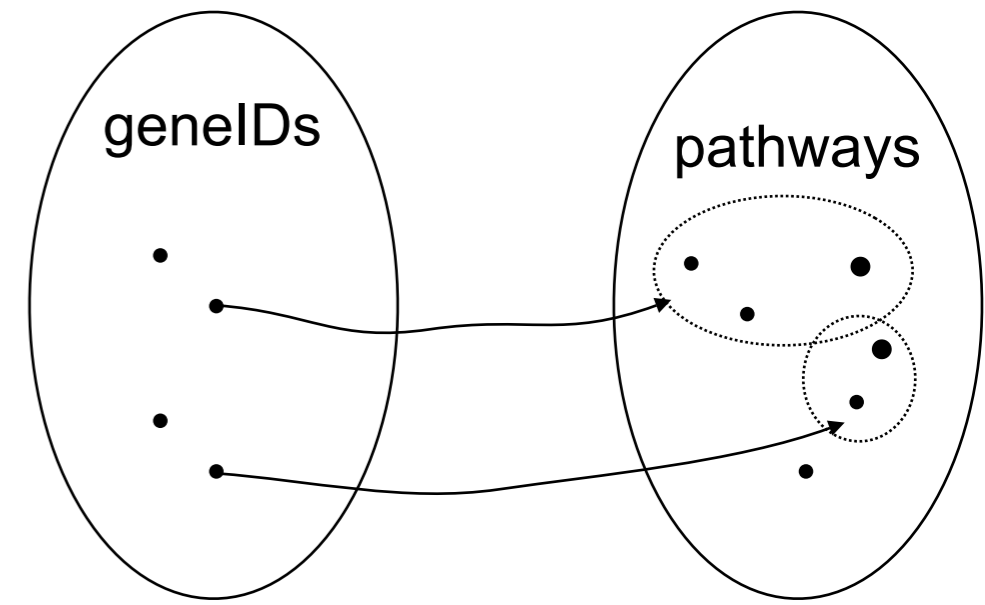
List-structured
KEGG gene ids:

[ [ mmu:26416 ],
  [ mmu:328788 ] ]

[ path:mmu04010 MAPK signaling,
path:mmu04370 VEGF signaling ]

[ [ path:mmu04210 Apoptosis, path:mmu04010 MAPK signaling, ...],
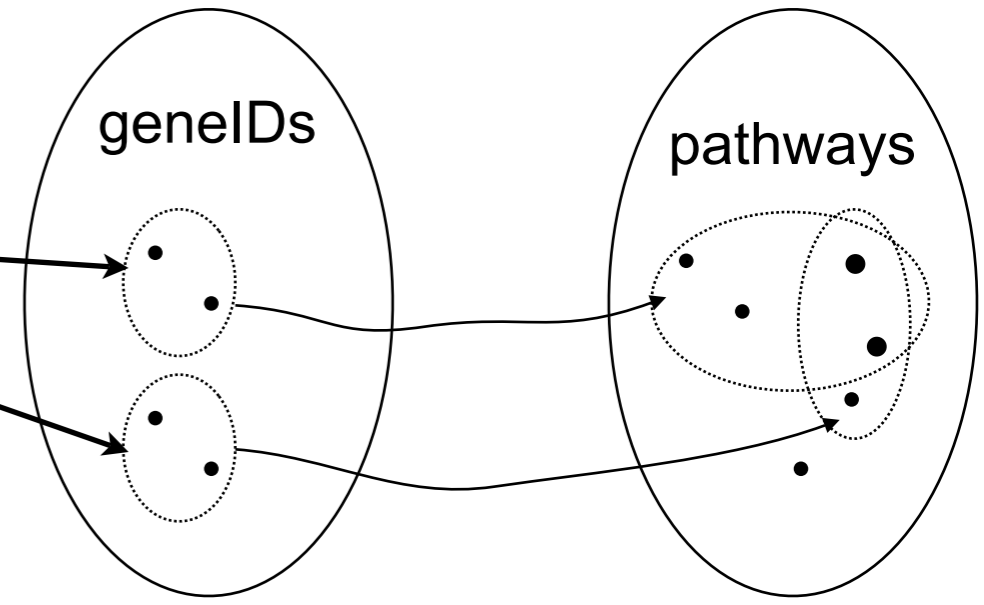  [ path:mmu04010 MAPK signaling , path:mmu04620 Toll-like receptor, ...] ]

List-structured
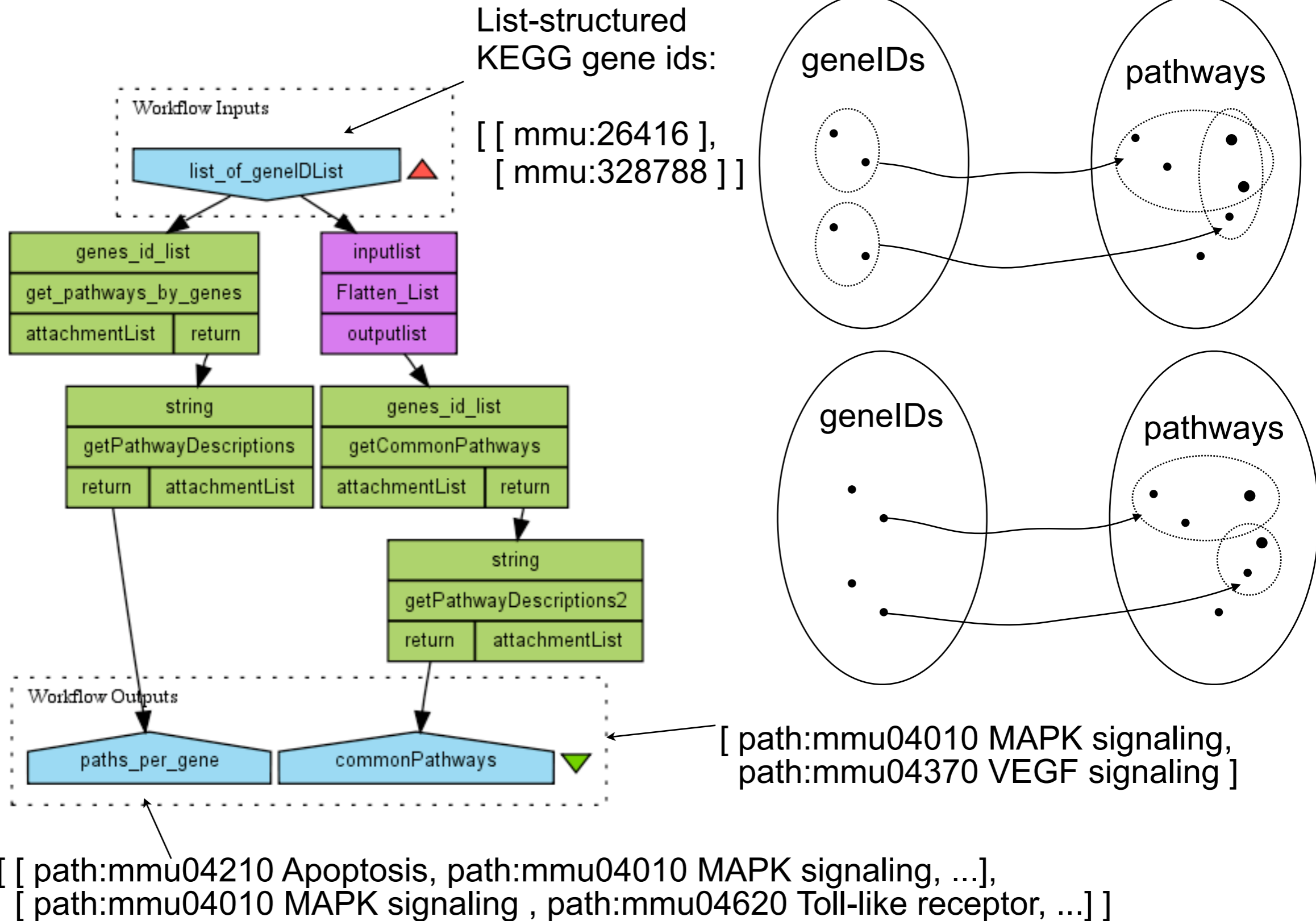KEGG gene ids:

[ [ mmu:26416 ],
[ mmu:328788 ] ]
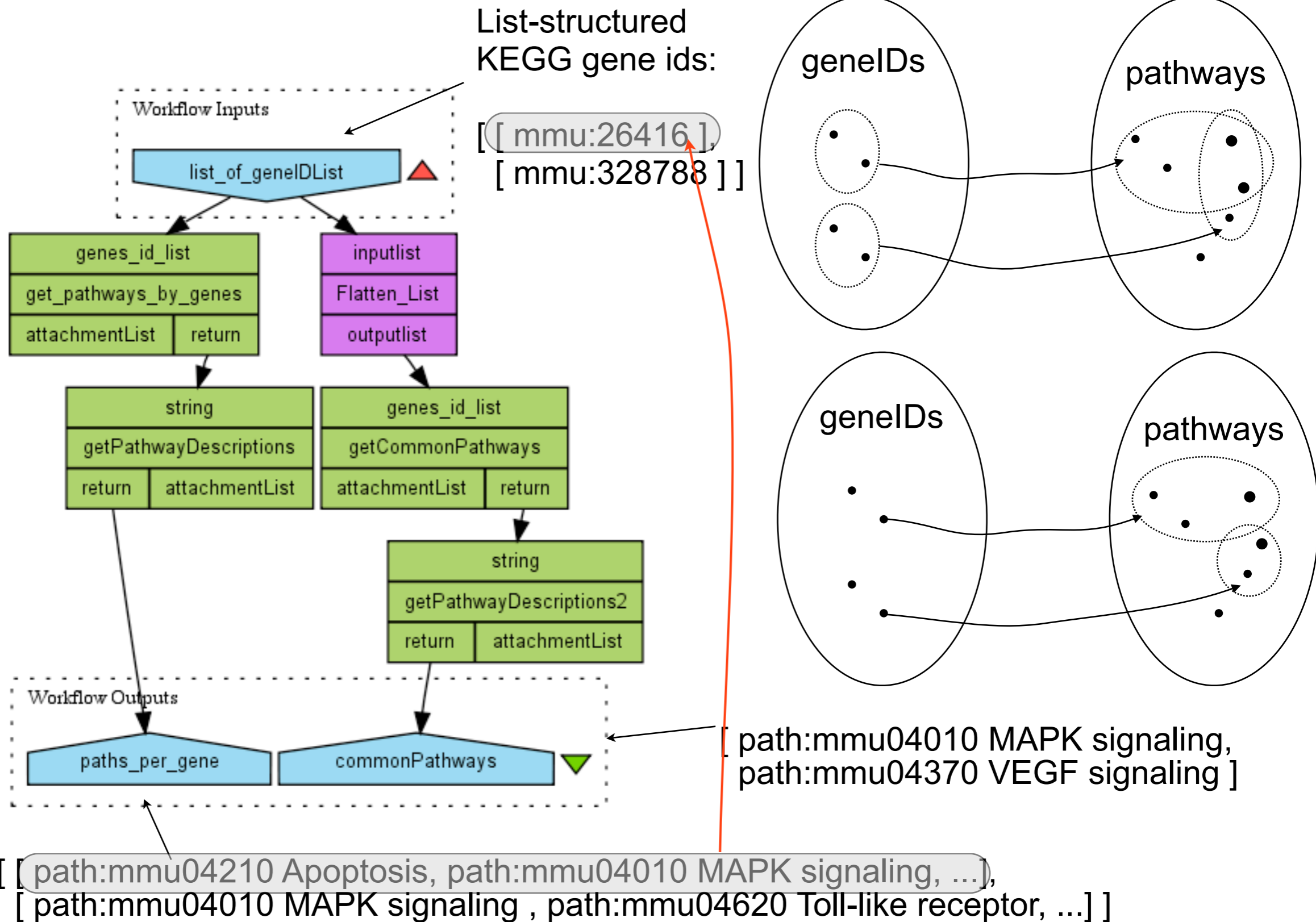
[ path:mmu04010 MAPK signaling,
path:mmu04370 VEGF signaling ]

[ [ path:mmu04210 Apoptosis, path:mmu04010 MAPK signaling, ...],
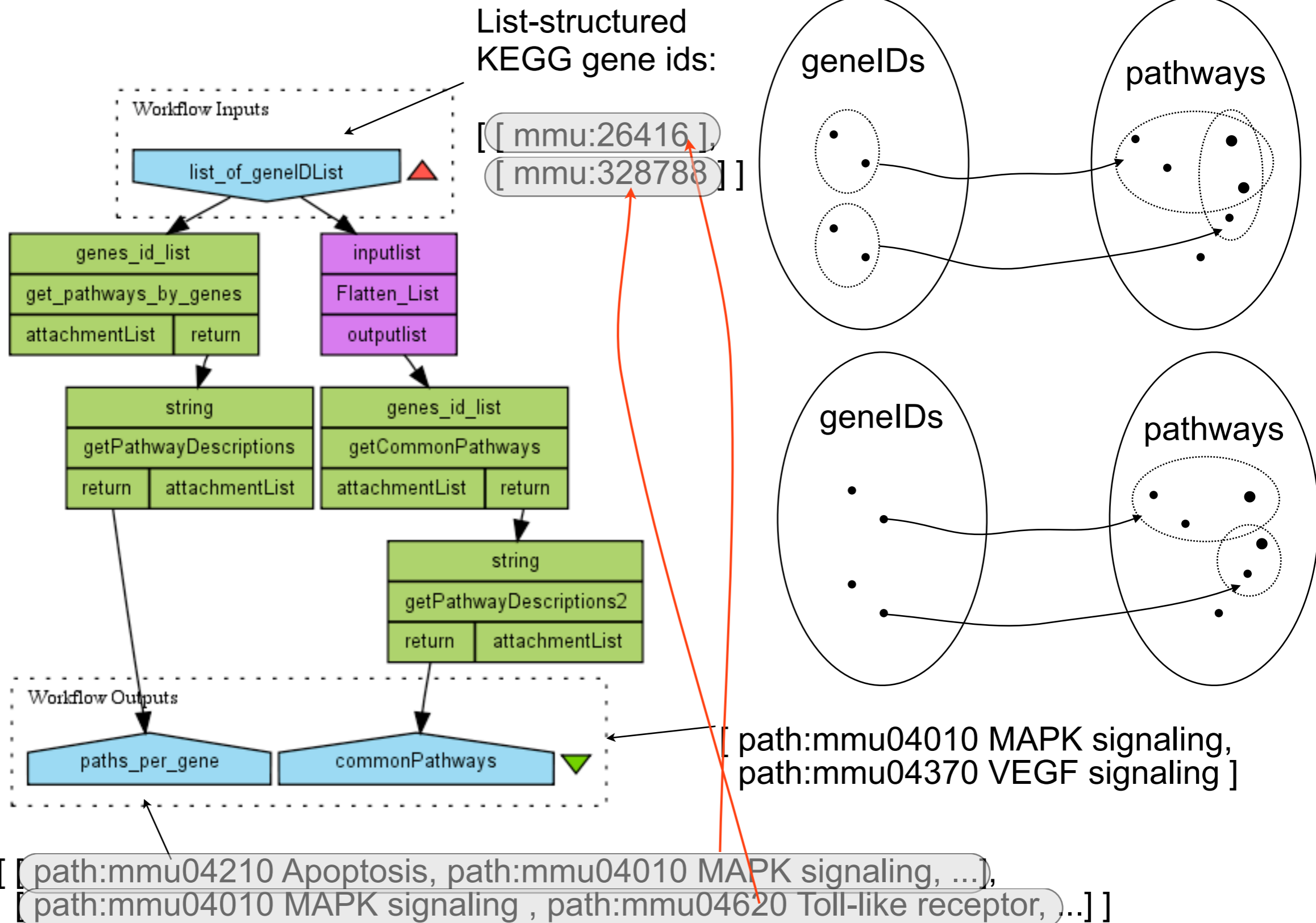[ path:mmu04010 MAPK signaling , path:mmu04620 Toll-like receptor, ...] ]

- Taverna type system: strings + nested lists
  - "cat", ["cat", "dog"], [ ["cat", "dog"], ["large", "small"] ]

- Taverna dataflow model: data-driven execution
  - services activate when input is ready

- Workflow provenance: a detailed trace of workflow execution
  - which services were executed
  - when
  - inputs used, outputs produced

- Taverna type system: strings + nested lists
  - "cat", ["cat", "dog"], [ ["cat", "dog"], ["large", "small"] ]

- Taverna dataflow model: data-driven execution
  - services activate when input is ready

- Workflow provenance: a detailed trace of workflow execution
  - which services were executed
  - when
  - inputs used, outputs produced

> Taverna dataflow model + provenance traces
> can be a powerful combination

- Causal relations:
  - ☑ which pathway sets come from which gene sets?
  - – which processes contributed to producing this image?
  - – which process(es) caused this data to be incorrect?
  - – which data caused this process to fail?

- Process and data analytics:
  - – show me the variations in output in relation to an input parameter sweep (multiple process runs)
  - – how often has my favourite service been executed?
    - on what inputs?
  - – who produced this data?
  - – how often does this pathway turn up when the input genes range over a certain set S?

Focus is on the data: the observable outcomes from a process

|  | raw provenance metadata | provenance metadata + interpretation framework |
|---|---|---|
| **design** | • process structure (workflow graph)<br>• history of process composition - reuse<br>• process versions | • service annotations:<br>• ex. get_pathways_by_genes<br>• who created /edited: attribution<br>• why: purpose, intent |
| **execution** | process events:<br>- service invocation<br>- data production / consumption<br>- causal dependency graphs<br>ex.:<br>- list_of_geneIDList = [ a, b, c]<br>- paths_per_gene = [ [d,e,f], [g,h,j]]<br>- ... in run #32 | - data annotations,<br>results interpretation in terms of conceptual data model:<br>set of pathways → gene sets |

10

| | raw provenance metadata | provenance metadata + interpretation framework |
|---|---|---|
| **design** | • exploiting semantic properties of the process structure to improve provenance exploitation<br><br>• exploring process space across versions and structural similarities<br><br>• graph matching | • semantic-based search of process space |
| **execution** | - enabling partial re-runs of resource-intensive workflows<br><br>- storing very large provenance traces that accumulate over time<br><br>- efficient query over large traces<br><br>- presentation of query answers | - semantic-based query answering over annotated traces |

| | raw provenance metadata | provenance metadata + interpretation framework |
|---|---|---|
| **design** | • exploiting semantic properties of the process structure to improve provenance exploitation<br><br>• exploring process space across versions and structural similarities<br><br>• graph matching | • semantic-based search of process space |
| **execution** | - enabling partial re-runs of resource-intensive workflows<br><br>- storing very large provenance traces that accumulate over time<br><br>- efficient query over large traces<br><br>- presentation of query answers | - semantic-based query answering over annotated traces |

The rest of this talk!

11

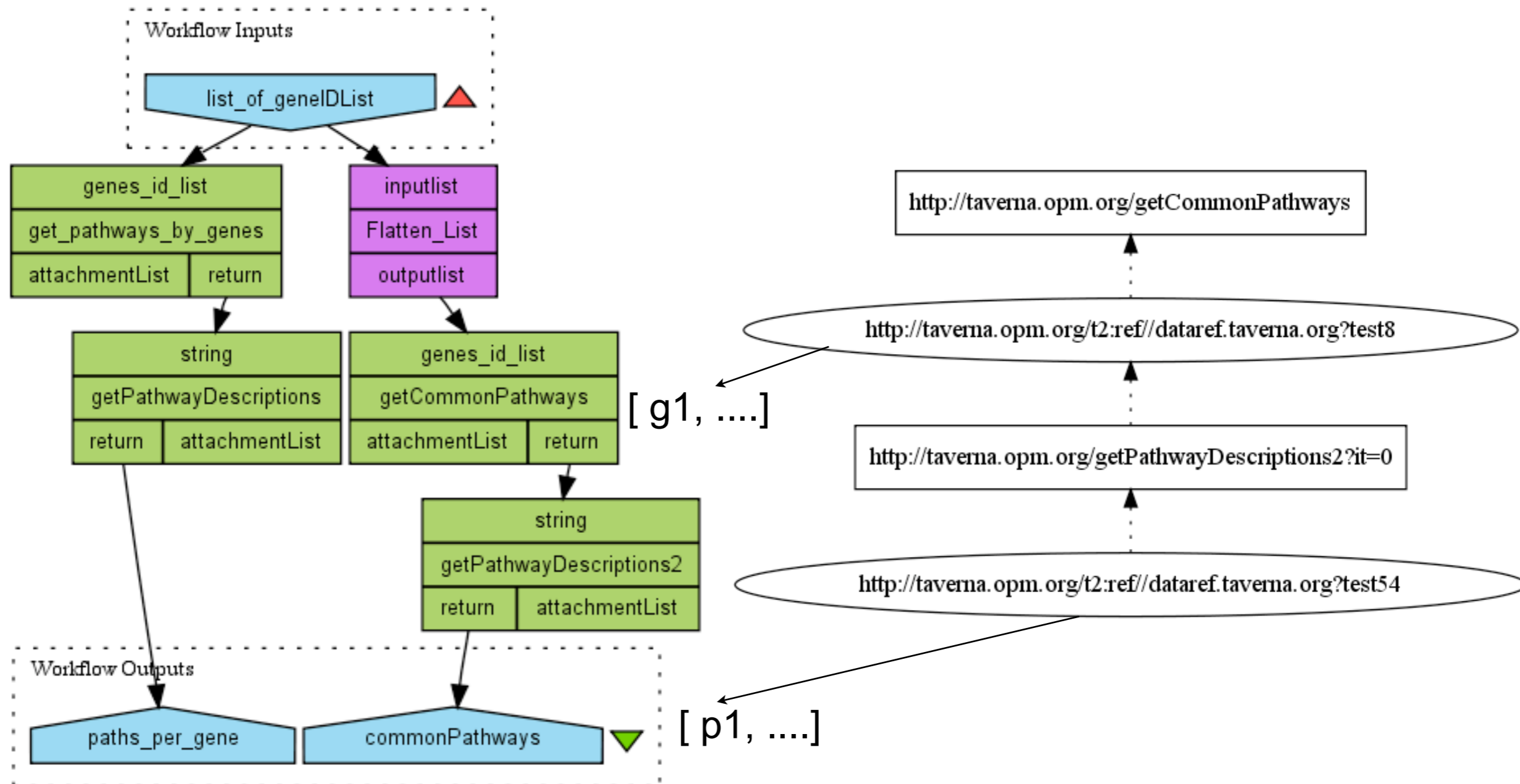| | raw provenance metadata | provenance metadata + interpretation framework |
|---|---|---|
| **design** | • exploiting semantic properties of the process structure to improve provenance exploitation<br><br>• exploring process space across versions and structural similarities<br><br>• graph matching | • semantic-based search of process space |
| **execution** | - enabling partial re-runs of resource-intensive workflows<br><br>- storing very large provenance traces that accumulate over time<br><br>- efficient query over large traces<br><br>- presentation of query answers | - semantic-based query answering over annotated traces |

The rest of this talk!

11

- **Lineage queries** involve traversing a *provenance graph* from bottom to top

- In most approaches, the originating process are not used for querying
- consequence: query requires provenance graph traversal
  - large traces → computationally complex
  - view materialization used in practice to get around the computational complexity



(a) specification

(b) one run

(c) another run

Fig. 1. Protein annotation workflow specification and runs

I - Focusing:

Not all processors are interesting:

– report lineage only at specified nodes in the graph

List-structured
KEGG gene ids:

[ [ mmu:26416 ],
[ mmu:328788 ] ]

**II - Granularity:**

Trace lineage for individual elements within collections - when possible!

[ path:mmu04010 MAPK signaling,
path:mmu04370 VEGF signaling ]

[ [ path:mmu04210 Apoptosis, path:mmu04010 MAPK signaling, ...],
[ path:mmu04010 MAPK signaling , path:mmu04620 Toll-like receptor, ...] ]

15

- III - Answer queries efficiently without special auxiliary data structures
- (and, please provide declarative query specification)



Example:

BACKTRACE
  (paths_per_gene[3,4],  paths_per_gene[1,2])
      AT get_pathway_by_genes
AND
  commonPathways[1]
      AT TOP

[ [ path:mmu04210 Apoptosis, path:mmu04010 MAPK signaling, ...],
  [ path:mmu04010 MAPK signaling , path:mmu04620 Toll-like receptor, ...]]

16

- III – Answer queries efficiently without special auxiliary data structures
- (and, please provide declarative query specification)



Example:

targets

BACKTRACE
   (paths_per_gene[3,4],  paths_per_gene[1,2])
      AT get_pathway_by_genes
AND
   commonPathways[1]
      AT TOP

[ [ path:mmu04210 Apoptosis, path:mmu04010 MAPK signaling, ...],
   [ path:mmu04010 MAPK signaling , path:mmu04620 Toll-like receptor, ...] ]

16

- III - Answer queries efficiently without special auxiliary data structures
- (and, please provide declarative query specification)

Example:

BACKTRACE
  (paths_per_gene[3,4],  paths_per_gene[1,2])
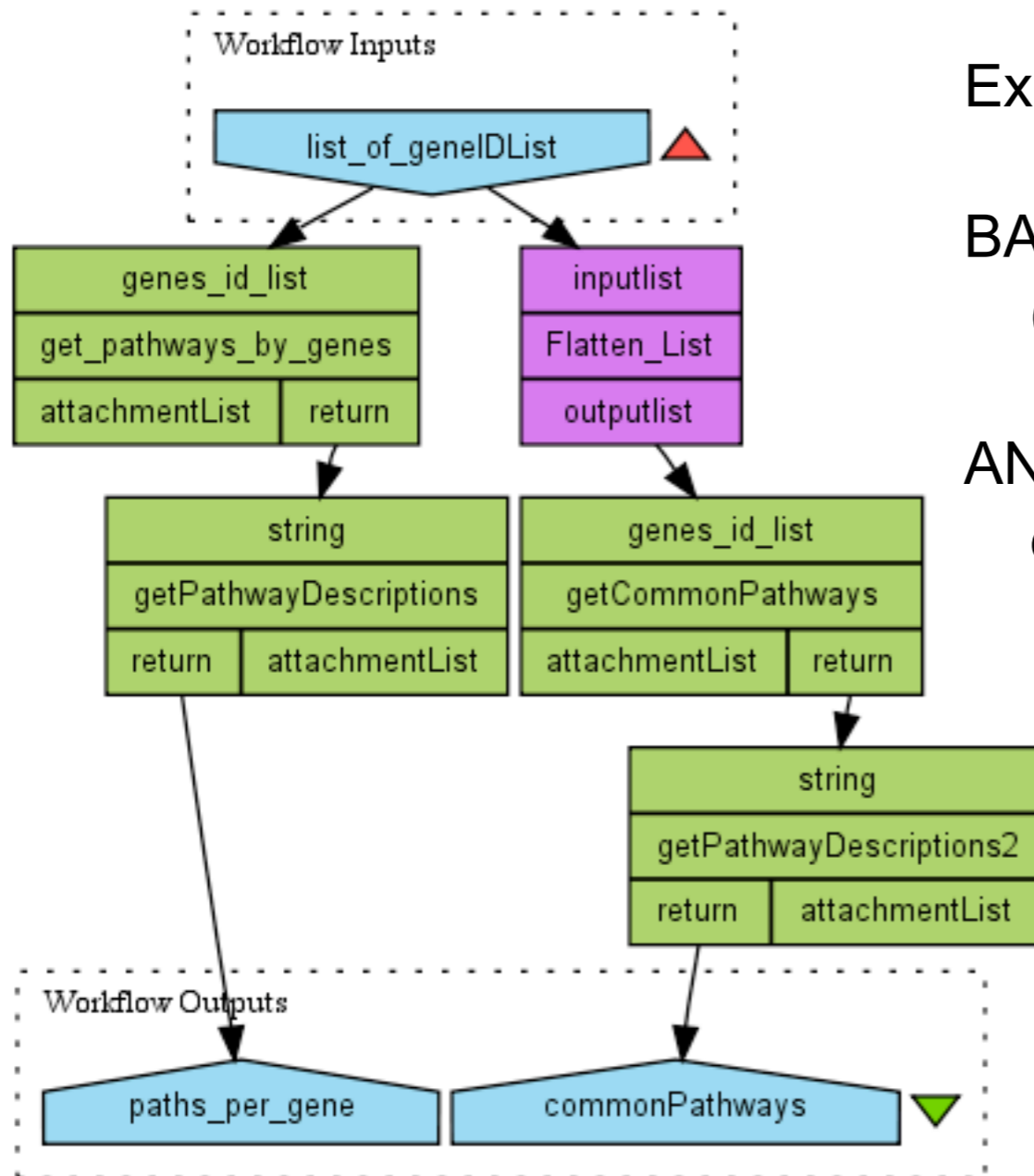    AT get_pathway_by_genes
AND
  commonPathways[1]
    AT TOP

targets

selected processors

[ [ path:mmu04210 Apoptosis, path:mmu04010 MAPK signaling, ...],
  [ path:mmu04010 MAPK signaling , path:mmu04620 Toll-like receptor, ...] ]

16

$c = [c_1 \ldots c_k]$

$a = [a_1 \ldots a_n]$

$b = [b_1 \ldots b_m]$

| $X_1$ | $X_2$ | $X_3$ |
|---|---|---|
| P | | |
| Y | | |

$c = [c_1 \ldots c_k]$

$(0,1)$    $(1,1)$    $(0,1)$

$a = [a_1 \ldots a_n]$    X_1    X_2    X_3    $b = [b_1 \ldots b_m]$

P

Y

Depth mismatch between declared / offered type:

$depth(P{:}X1) = 0$ but $depth(a) = 1$

$depth(P{:}X2) = depth(c) = 1$

$depth(P{:}X3) = 1$ but $depth(c) = 1$

$c = [c_1 \dots c_k]$

$(0,1) \quad (1,1) \quad (0,1)$

$a = [a_1 \dots a_n]$

$X_1$ | $X_2$ | $X_3$

P

Y

$b = [b_1 \dots b_m]$

Depth mismatch between declared / offered type:

depth(P:X1) = 0 but depth(a) = 1

depth(P:X2) = depth(c) = 1

depth(P:X3) = 1 but depth(c) = 1

How *y* is computed at P:

let $I = a \otimes b = [ \, [ \, <a_i, b_j> \mid b_j \in b \, ] \mid a_i \in a \, ]$   // cross product

$I' = [ \, [ \, <a_i, c, b_j> \mid b_j \in b \, ] \mid a_i \in a \, ]$ // same product but with c interleaved

$y = (\textbf{map} \ (\textbf{map} \ P) \ I') = [(\text{map } P \ [ \, <a_1,c, b_1> \ \dots \ <a_1,c, b_m>]), \dots,$
$(\text{map } P \ [ \, <a_n,c, b_1> \ \dots \ <a_n,c, b_m>]) \, ] =$
$[ \, [y_{11} \ \dots \ y_{1n}], \dots [y_{n1} \ \dots \ y_{nm}] \, ]$

$c = [c_1 \ldots c_k]$

$(0,1)$   $(1,1)$   $(0,1)$

$a = [a_1 \ldots a_n]$ ⋯⋯ $X_1$ | $X_2$ | $X_3$ ⋯⋯ $b = [b_1 \ldots b_m]$

P

Y

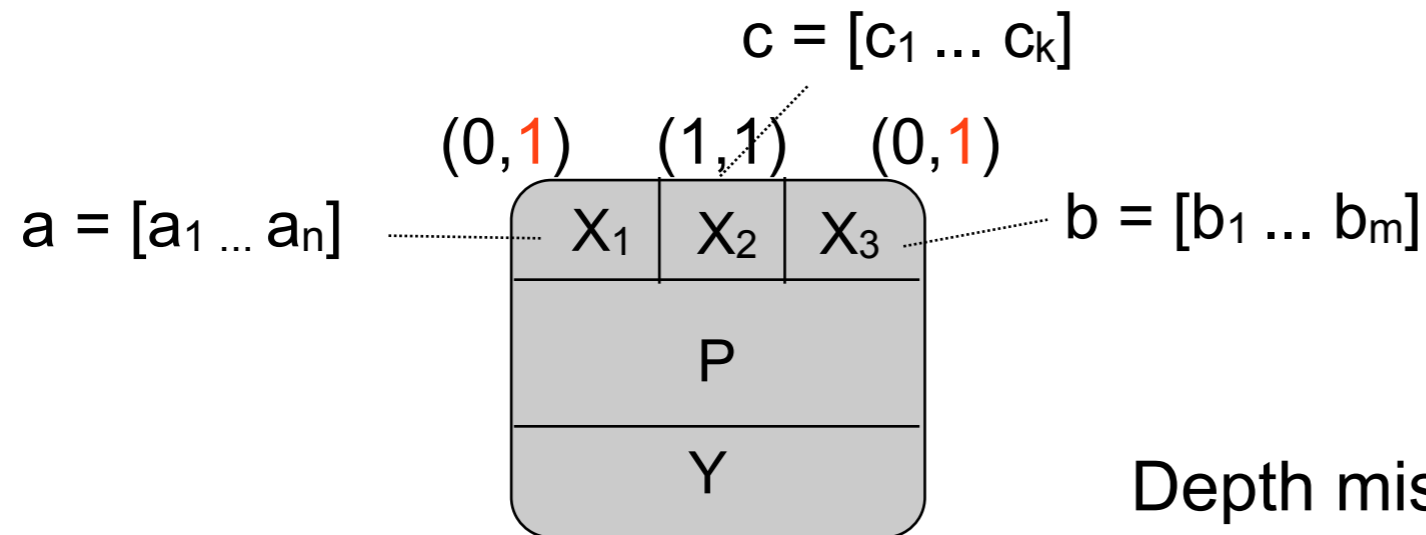$y = [\ [y_{11} \ldots y_{1n}],$
$\ldots$
$[y_{m1} \ldots y_{mn}]\ ]$

Depth mismatch between declared / offered type:

$\text{depth}(P:X1) = 0$ but $\text{depth}(a) = 1$
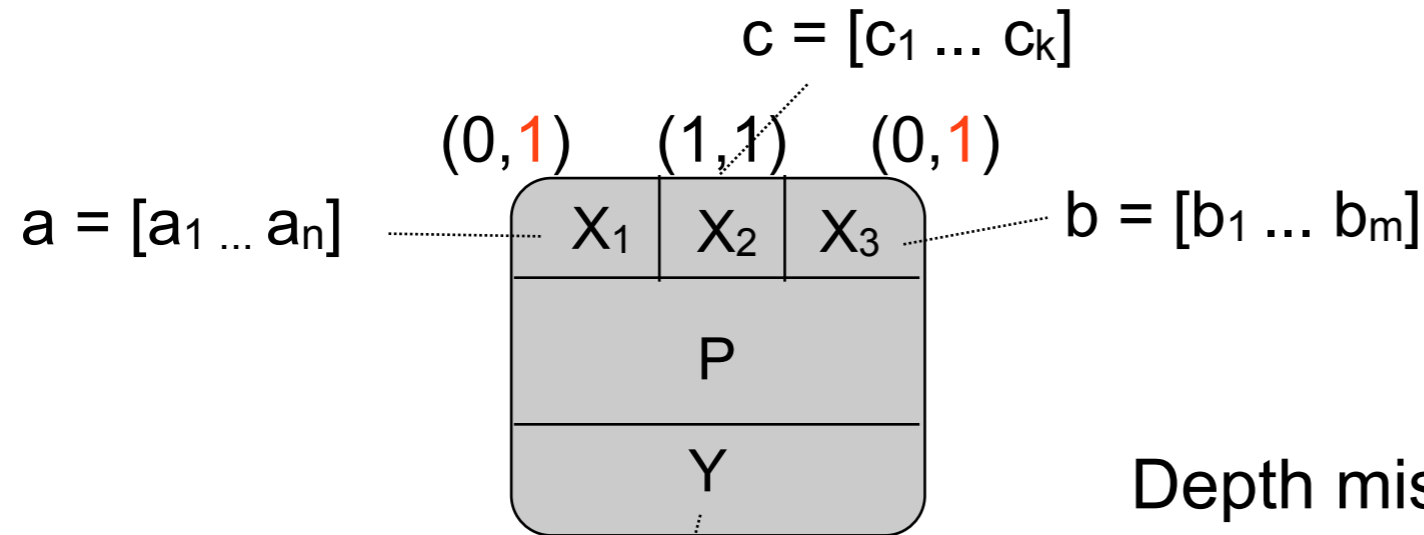
$\text{depth}(P:X2) = \text{depth}(c) = 1$

$\text{depth}(P:X3) = 1$ but $\text{depth}(c) = 1$

How *y* is computed at P:

let $I = a \otimes b = [\ [\ <a_i, b_j> \mid b_j \in b\ ] \mid a_i \in a\ ]$   // cross product

$I' = [\ [\ <a_i, c, b_j> \mid b_j \in b\ ] \mid a_i \in a\ ]$ // same product but with c interleaved

$y = (\textbf{map}\ (\textbf{map}\ P)\ I') = [(\text{map}\ P\ [\ <a_1,c, b_1> \ldots <a_1,c, b_m>]), \ldots,$
$(\text{map}\ P\ [\ <a_n,c, b_1> \ldots <a_n,c, b_m>])\ ] =$
$[\ [y_{11} \ldots y_{1n}], \ldots [y_{n1} \ldots y_{nm}]\ ]$

$c = [c_1 ... c_k]$

$(0,1)$  $(1,1)$  $(0,1)$

$a = [a_{1 ... } a_n]$

| $X_1$ | $X_2$ | $X_3$ |
|---|---|---|
| P | | |
| Y | | |

$b = [b_1 ... b_m]$

Depth mismatch between declared / offered type:

$y = [ [y_{11} ... y_{1n}],$
$...$
$[y_{m1} ... y_{mn}] ]$

bottom line:
$y_{ij}$ depends only on values $a_i$, c, $b_j$

$depth(P:X3) = 1$ but $depth(c) = 1$

How *y* is computed at P:

let $I = a \otimes b = [ [ <a_i, b_j> | b_j \in b ] | a_i \in a ]$   // cross product

$I' = [ [ <a_i, c, b_j> | b_j \in b ] | a_i \in a ]$ // same product but with c interleaved

$y = (\textbf{map} (\textbf{map} P) I') = [(\text{map } P [ <a_1,c, b_1> ... <a_1,c, b_m>]), ...,$
$(\text{map } P [ <a_n,c, b_1> ... <a_n,c, b_m>]) ] =$
$[ [y_{11} ... y_{1n}], ... [y_{n1} ... y_{nm}] ]$

Workflow structure graph

$v = [v_1 \ldots v_n]$

$w$

| X |
|---|
| Q |
| Y |

| X |
|---|
| R |
| Y |

$c = [c_1 \ldots c_k]$

$a = [a_1 \ldots a_n]$

| $X_1$ | $X_2$ | $X_3$ |
|---|---|---|
| P | | |
| Y | | |

$b = [b_1 \ldots b_m]$

$y = [ \; [y_{11} \ldots y_{1n}],$
$\ldots$
$[y_{m1} \ldots y_{mn}] \; ]$

## Workflow structure graph

$v = [v_1 ... v_n]$       $w$

| X |
|---|
| Q |
| Y |

| X |
|---|
| R |
| Y |

$c = [c_1 ... c_k]$

$a = [a_1 ... a_n]$

| $X_1$ | $X_2$ | $X_3$ |
|---|---|---|
| P | | |
| Y | | |

$b = [b_1 ... b_m]$

$y = [ \ [y_{11} ... y_{1n}],$
$...$
$[y_{m1} ... y_{mn}] \ ]$

## Provenance graph

## Workflow structure graph

$v = [v_1 \ldots v_n]$

$w$

| X |
|---|
| Q |
| Y |

| X |
|---|
| R |
| Y |

$c = [c_1 \ldots c_k]$

$a = [a_1 \ldots a_n]$

| $X_1$ | $X_2$ | $X_3$ |
|---|---|---|
| P | | |
| Y | | |

$b = [b_1 \ldots b_m]$

$y = [\ [y_{11} \ldots y_{1n}],$
$\ldots$
$[y_{m1} \ldots y_{mn}]\ ]$

## Provenance graph

$v_1$ … $v_n$ $w$

$a_1$ … $a_n$ $b_1$ … $b_m$

$y_{11}$ … $y_{mn}$

Hypothesis:
we can exploit the static workflow graph structure to avoid explicitly traversing the entire trace to answer a query

18

| X₁ | X₂ | X₃ |
|---|---|---|
| P | | |
| Y | | |

1) In general the actual depth at the output is:

$$depth(y) = depth(Y) + \Sigma \, \delta(X_i = x_i)$$

where $\delta(X_i = x_i) = depth(x_i) - depth(X_i)$

Therefore:
$\delta(X=x)$ can be computed statically on the workflow graph structure,
- given the declared $depth(X)$
- using a simple propagation algorithm

$$Y[i.j] \; \rightarrow X1[i], \; X2[], \; X3[j]$$

$$[i_1 \, . \, i_2 \, . \, ... \, . \, i_k] = \underline{\hspace{8cm}}$$

$(0,\textcolor{red}{1})$   $(1,1)$   $(0,\textcolor{red}{1})$

| $X_1$ | $X_2$ | $X_3$ |
|-------|-------|-------|
| P | | |
| Y | | |

$(0,\textcolor{red}{2})$

$Y[i.j] \rightarrow X1[i], X2[], X3[j]$

1) In general the actual depth at the output is:

$$depth(y) = depth(Y) + \Sigma\, \delta(X_i = x_i)$$

where $\delta(X_i = x_i) = depth(x_i) - depth(X_i)$

Therefore:
$\delta(X=x)$ can be computed statically on the workflow graph structure,

- given the declared $depth(X)$
- using a simple propagation algorithm

$[i_1 . i_2 . \ldots . i_k] = $ ————————————————————

$(0,1)$   $(1,1)$   $(0,1)$

| $X_1$ | $X_2$ | $X_3$ |
|-------|-------|-------|
| P | | |
| Y | | |

$(0,2)$

$Y[i.j] \rightarrow X1[i], X2[], X3[j]$

1) In general the actual depth at the output is:

$$depth(y) = depth(Y) + \Sigma\ \delta(X_i = x_i)$$

where $\delta(X_i = x_i) = depth(x_i) - depth(X_i)$

Therefore:
$\delta(X=x)$ can be computed statically on the workflow graph structure,

- given the declared $depth(X)$
- using a simple propagation algorithm

$$[i_1 . i_2 . \dots . i_k] = \frac{\delta(X_1 = x_1)}{\rule{6cm}{0.4pt}}$$

$(0,1)$   $(1,1)$   $(0,1)$

| $X_1$ | $X_2$ | $X_3$ |
|-------|-------|-------|
| P     |       |       |
| Y     |       |       |

$(0,2)$

1) In general the actual depth at the output is:

$$depth(y) = depth(Y) + \Sigma \; \delta(X_i = x_i)$$

where $\delta(X_i = x_i) = depth(x_i) - depth(X_i)$

Therefore:
$\delta(X=x)$ can be computed statically on the workflow graph structure,
- given the declared depth(X)
- using a simple propagation algorithm

$Y[i.j] \; \rightarrow X1[i], X2[], X3[j]$

$[i_1 . i_2 . ... . i_k] =$ _____

_____

$X_1$

$(0,1)$      $(1,1)$      $(0,1)$

| X₁ | X₂ | X₃ |
|----|----|----|
| P  |    |    |
| Y  |    |    |

$(0,2)$

$Y[i.j] \rightarrow X1[i], X2[], X3[j]$

1) In general the actual depth at the output is:

$depth(y) = depth(Y) + \Sigma\, \delta(X_i = x_i)$
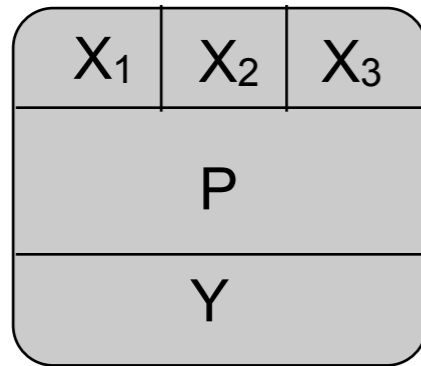
where $\delta(X_i = x_i) = depth(x_i) - depth(X_i)$

Therefore:
$\delta(X=x)$ can be computed statically on the workflow graph structure,
- given the declared $depth(X)$
- using a simple propagation algorithm

$$[i_1 . i_2 . \ldots . i_k] = \frac{\delta(X_2 = x_2)}{\phantom{xxxxxxxxxxxxxxxxxxx}}$$

$$\frac{\phantom{xxxxx}}{X_1}$$

19

$(0,1)$    $(1,1)$    $(0,1)$

| $X_1$ | $X_2$ | $X_3$ |
|---|---|---|
| P | | |
| Y | | |

$(0,2)$

$Y[i.j] \rightarrow X1[i], X2[], X3[j]$

$[i_1 . i_2 . \dots . i_k] =$

———————    ———————
$X_1$              $X_2$

1) In general the actual depth at the output is:

depth(y) = depth(Y) + $\Sigma$ $\delta(X_i = x_i)$

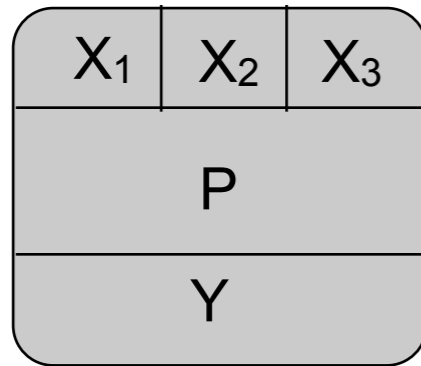where $\delta(X_i = x_i) =$ depth($x_i$) - depth($X_i$)

Therefore:
$\delta(X=x)$ can be computed statically on the workflow graph structure,
- given the declared depth(X)
- using a simple propagation algorithm

———————————————

$(0,1)$ $(1,1)$ $(0,1)$

| $X_1$ | $X_2$ | $X_3$ |
|:---:|:---:|:---:|
| P | | |
| Y | | |

$(0,2)$

$Y[i.j] \rightarrow X1[i], X2[], X3[j]$

$[i_1 . i_2 . ... . i_k] =$

$$\frac{\qquad}{X_1} \qquad \frac{\qquad}{X_2}$$

1) In general the actual depth at the output is:

$$depth(y) = depth(Y) + \Sigma \, \delta(X_i = x_i)$$
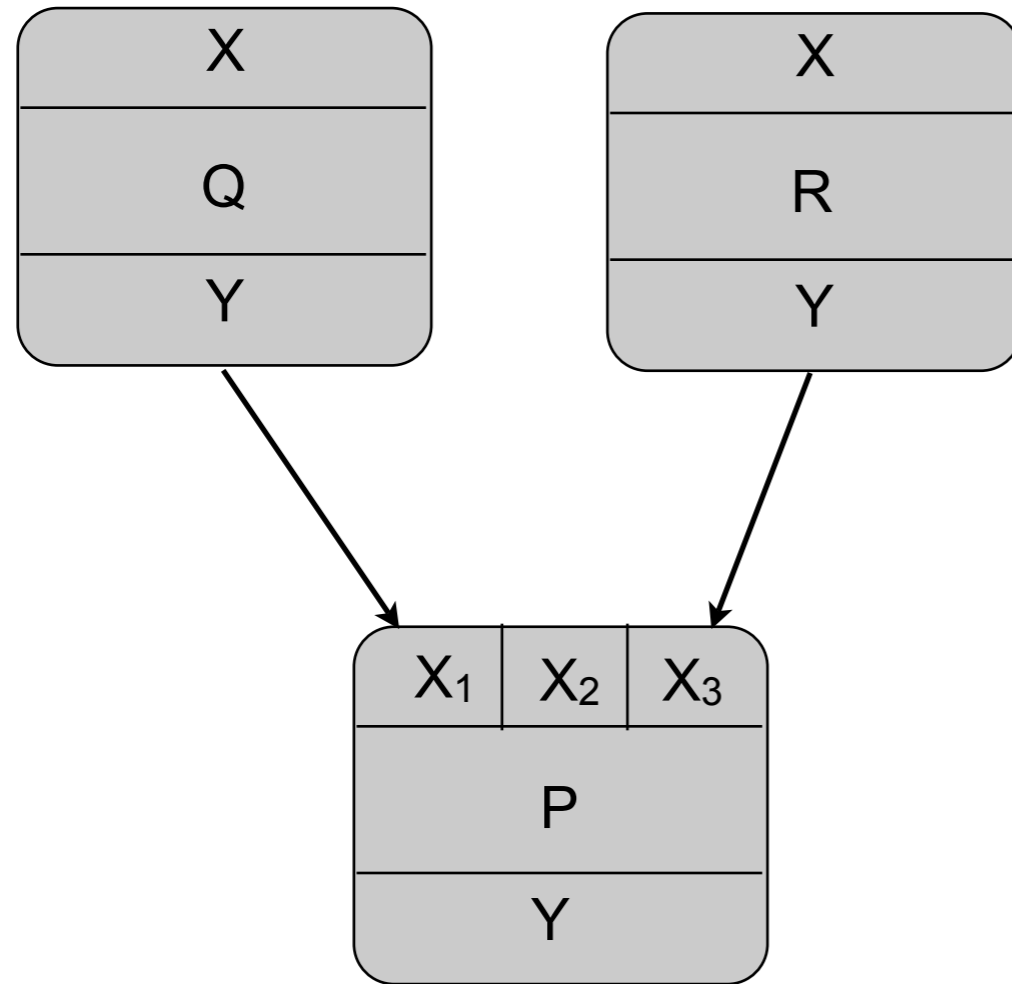
where $\delta(X_i = x_i) = depth(x_i) - depth(X_i)$
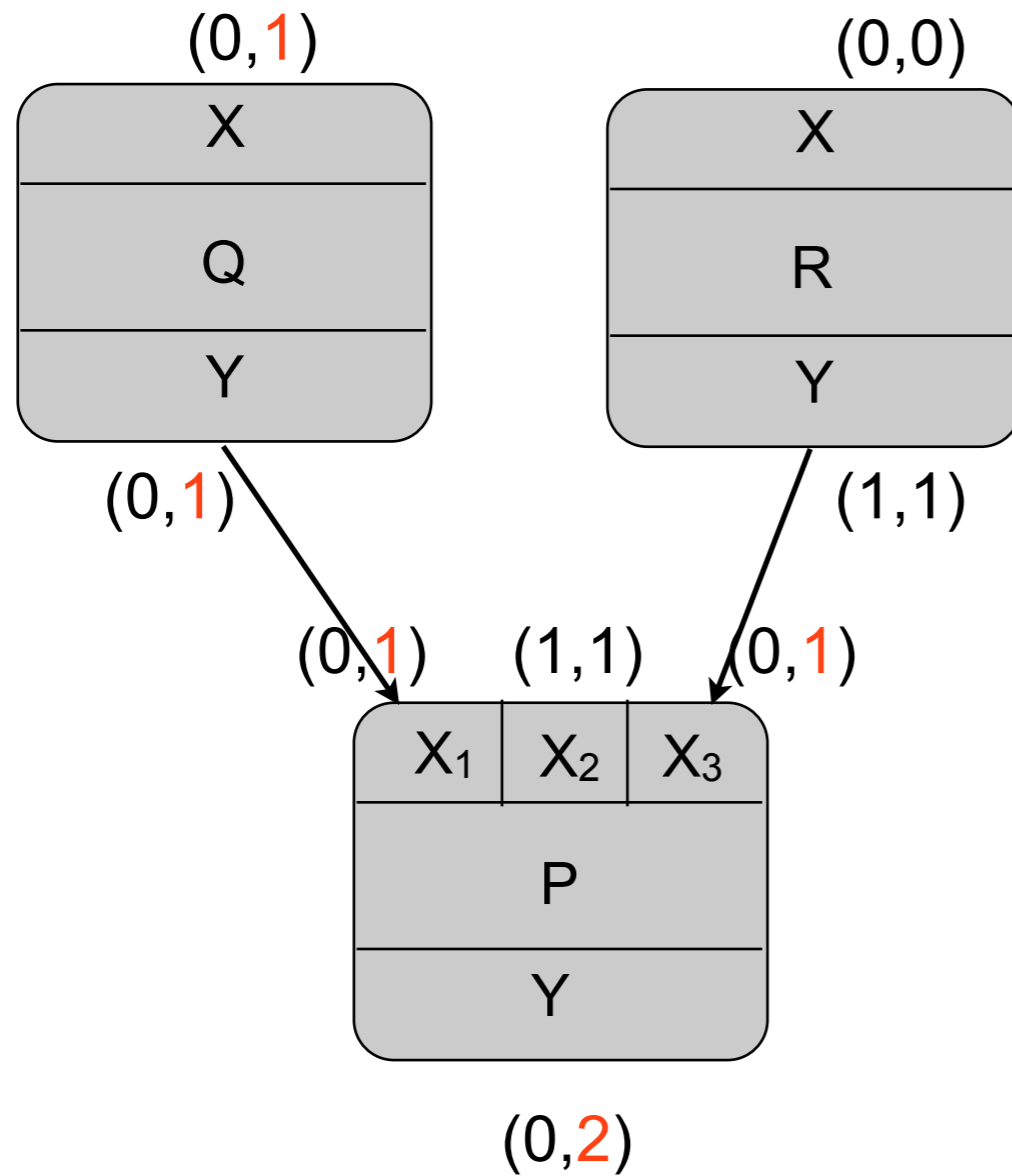
Therefore:
$\delta(X=x)$ can be computed statically on the workflow graph structure,

- given the declared depth(X)
- using a simple propagation algorithm

$$\frac{\delta(X_k = x_k)}{\qquad}$$

19

$(0,1)$    $(1,1)$    $(0,1)$

| $X_1$ | $X_2$ | $X_3$ |
|-------|-------|-------|
| P | | |
| Y | | |

$(0,2)$

$Y[i.j] \rightarrow X1[i], X2[], X3[j]$

$[i_1 . i_2 . \ldots . i_k] =$

$\underline{\phantom{XXXX}}$     $\underline{\phantom{XXXX}}$              $\underline{\phantom{XXXX}}$
   $X_1$              $X_2$                               $X_k$

1) In general the actual depth at the output is:

$\text{depth}(y) = \text{depth}(Y) + \sum \delta(X_i = x_i)$

where $\delta(X_i = x_i) = \text{depth}(x_i) - \text{depth}(X_i)$

Therefore:
$\delta(X=x)$ can be computed statically on the workflow graph structure,
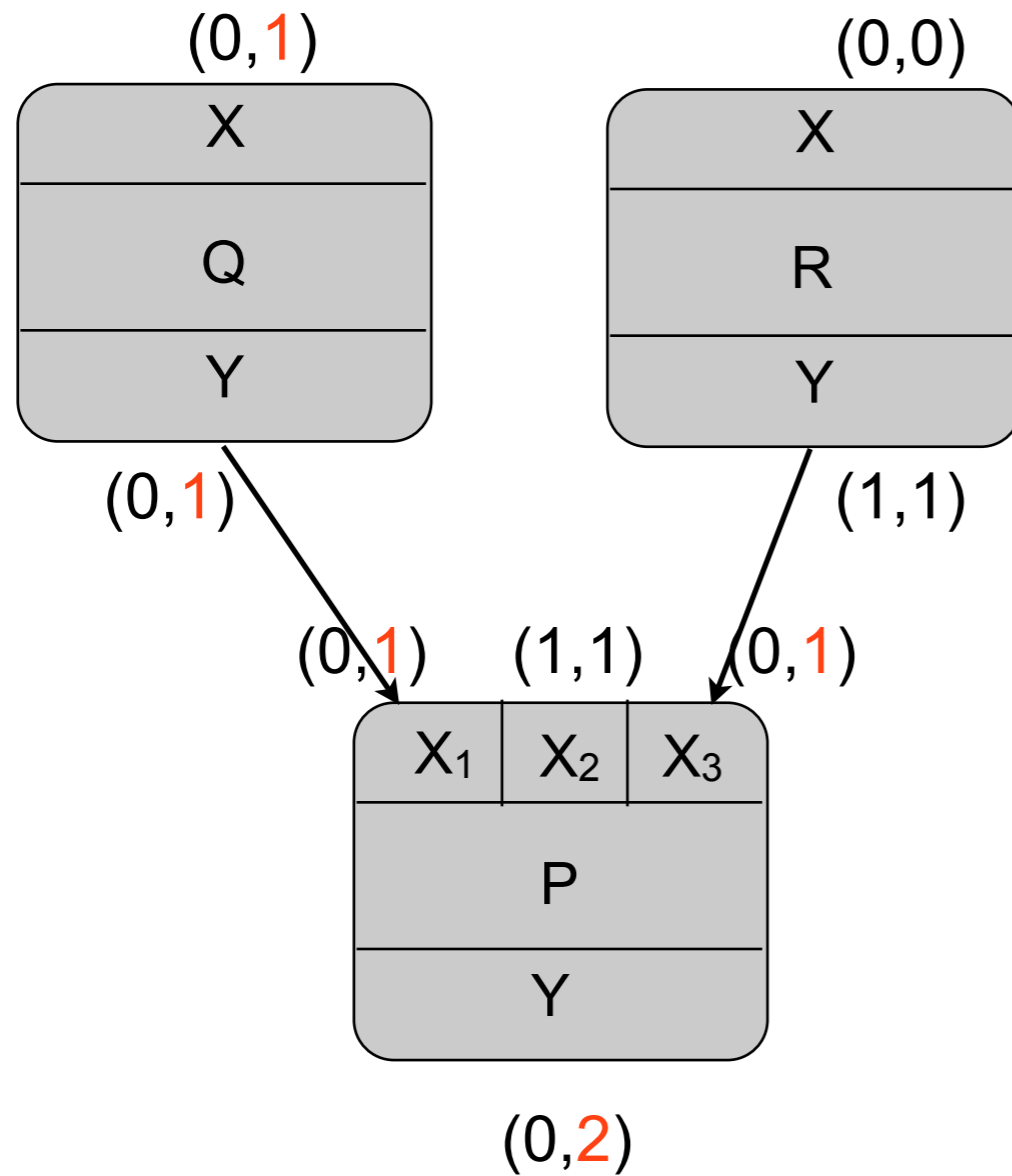- given the declared depth(X)
- using a simple propagation algorithm

$(0,1)$                                $(0,0)$

| X |
|---|
| Q |
| Y |

| X |
|---|
| R |
| Y |

$(0,1)$                     $(1,1)$

$(0,1)$     $(1,1)$     $(0,1)$

| $X_1$ | $X_2$ | $X_3$ |
|---|---|---|
| P | | |
| Y | | |

$(0,2)$

$Y = [ \ [ \ ...], \ ... \ [...] \ ]$

20

$(0,1)$

$(0,0)$

X

Q

Y

X

R

Y

$(0,1)$

$(1,1)$

$(0,1)$ $(1,1)$ $(0,1)$

$X_1$ $X_2$ $X_3$

P

Y

$(0,2)$

$Y = [ [ ...], ... [...] ]$

$lineage(P:Y[3,4]) \rightarrow lineage(P:X1[3]),$
$lineage(P:X2[]),$
$lineage(P:X2[4])$

20

$(0,1)$            $(0,0)$

| X |
|---|
| Q |
| Y |

| X |
|---|
| R |
| Y |

$(0,1)$            $(1,1)$

$(0,1)$    $(1,1)$    $(0,1)$

| $X_1$ | $X_2$ | $X_3$ |
|---|---|---|
| P | | |
| Y | | |

[3]          [4]

$(0,2)$

$Y = [ [ ...], ... [...] ]$

lineage(P:Y[3,4]) → lineage(P:X1[3]),
 lineage(P:X2[]),
 lineage(P:X2[4])

$(0,1)$        $(0,0)$

X

Q

Y

X

R

Y

$(0,1)$        $(1,1)$

$(0,1)$   $(1,1)$   $(0,1)$

$X_1$ | $X_2$ | $X_3$

[3]   P   [4]

Y

$(0,2)$

Y = [ [ ...], ... [...] ]

lineage(P:Y[3,4]) → lineage(P:X1[3]),
         lineage(P:X2[]),
         lineage(P:X2[4])

lineage(P:X1[3]) = lineage(Q:Y[3]) →
         lineage(Q:X[3])
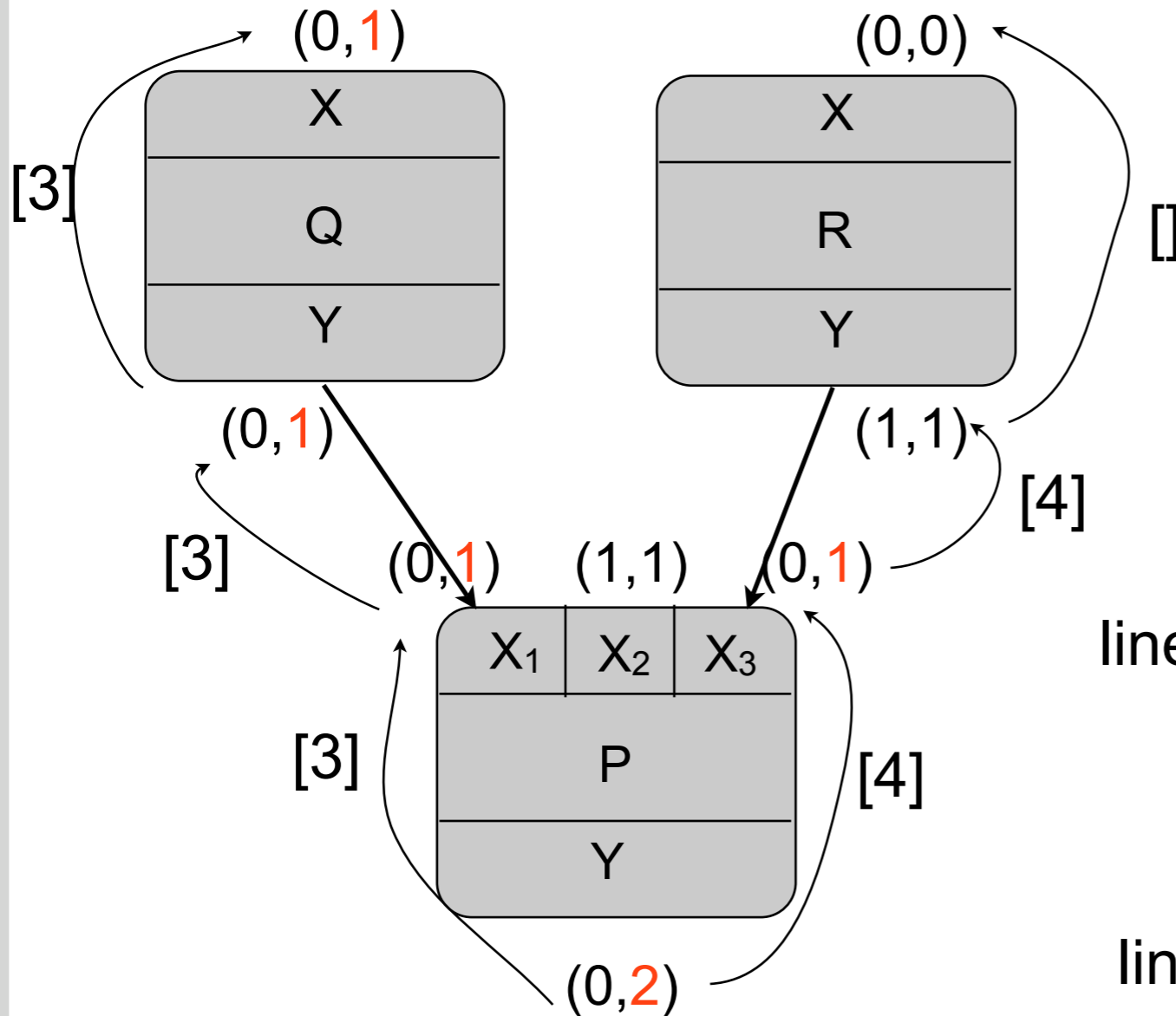
20

$(0,1)$

$(0,0)$

[3]

X

Q

Y

X

R

Y

$(0,1)$

$(1,1)$

[3]

$(0,1)$ $(1,1)$ $(0,1)$

$X_1$ $X_2$ $X_3$

[3] P [4]

Y

$(0,2)$

Y = [ [ ...], ... [...] ]

lineage(P:Y[3,4]) → lineage(P:X1[3]),
     lineage(P:X2[]),
     lineage(P:X2[4])

lineage(P:X1[3]) = lineage(Q:Y[3]) →
     lineage(Q:X[3])

$\text{lineage}(P{:}Y[3,4]) \rightarrow \text{lineage}(P{:}X1[3]),$
$\text{lineage}(P{:}X2[]),$
$\text{lineage}(P{:}X2[4])$

$\text{lineage}(P{:}X1[3]) = \text{lineage}(Q{:}Y[3]) \rightarrow$
$\text{lineage}(Q{:}X[3])$

$\text{lineage}(P{:}X3[4]) = \text{lineage}(R{:}Y[4]) \rightarrow$
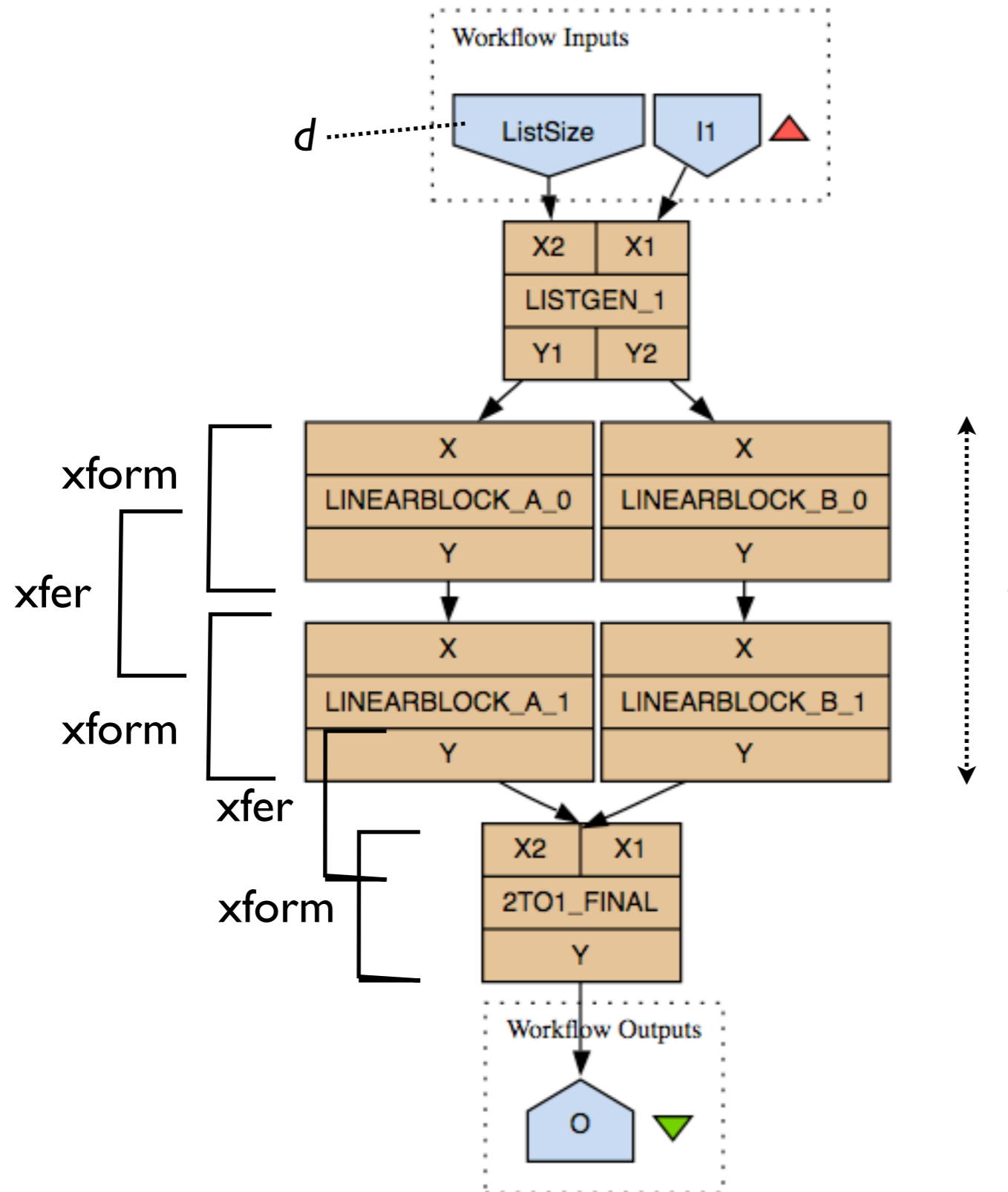$\text{lineage}(R{:}X[])$

$Y = [\ [\ ...],\ ...\ [...]\ ]$

$(0,1)$

$(0,0)$

[3]

[]

| X |
|---|
| Q |
| Y |

| X |
|---|
| R |
| Y |

$(0,1)$

$(1,1)$

[3]

[4]

$(0,1)$   $(1,1)$   $(0,1)$

| $X_1$ | $X_2$ | $X_3$ |
|---|---|---|
| P | | |
| Y | | |

[3]

[4]

$(0,2)$

Y = [ [ ...], ... [...] ]

lineage(P:Y[3,4]) → lineage(P:X1[3]),
                            lineage(P:X2[]),
                            lineage(P:X2[4])

lineage(P:X1[3]) = lineage(Q:Y[3]) →
                            lineage(Q:X[3])

lineage(P:X3[4]) = lineage(R:Y[4]) →
                            lineage(R:X[])

- Query processing:
  - alternating sequence of *xform* and *xfer* steps
- apply path projection at each *xform* step

- A complete granular and focused query can be answered by traversing the workflow graph alone
  - starting from the target vars
  - one simple query for each selected processor input port

- ## Scalability:
  - query time depends on size of workflow graph, not size of provenance graph
  - workflow graphs are small, fit in memory, can be indexed easily, etc.
  - search over a graph at least as large as the workflow graph is inevitable -- this is the baseline cost!

- ## Graceful degradation:
  - worst case is a completely unfocused query
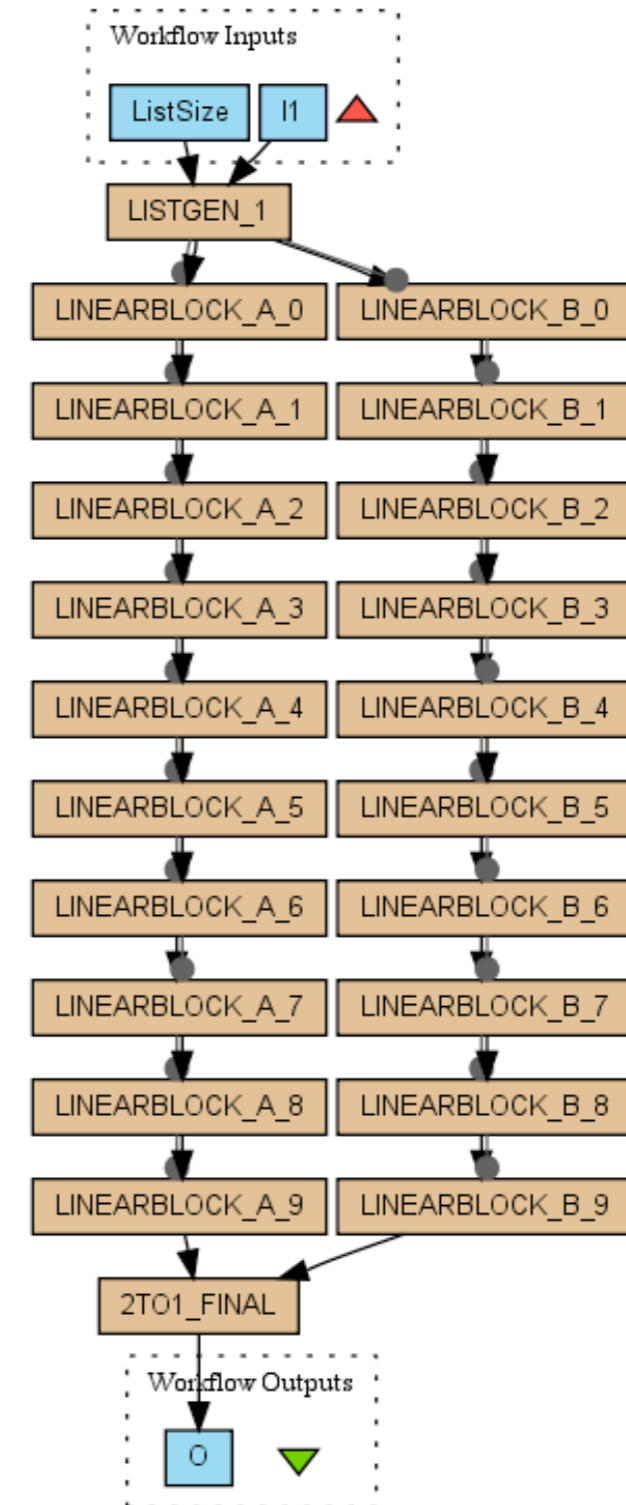  - one query to trace at each *xform* step
  - no worse than other approaches

- ## Fine-grain answers provided at no additional cost

- Performance evaluation performed on programmatically generated dataflows
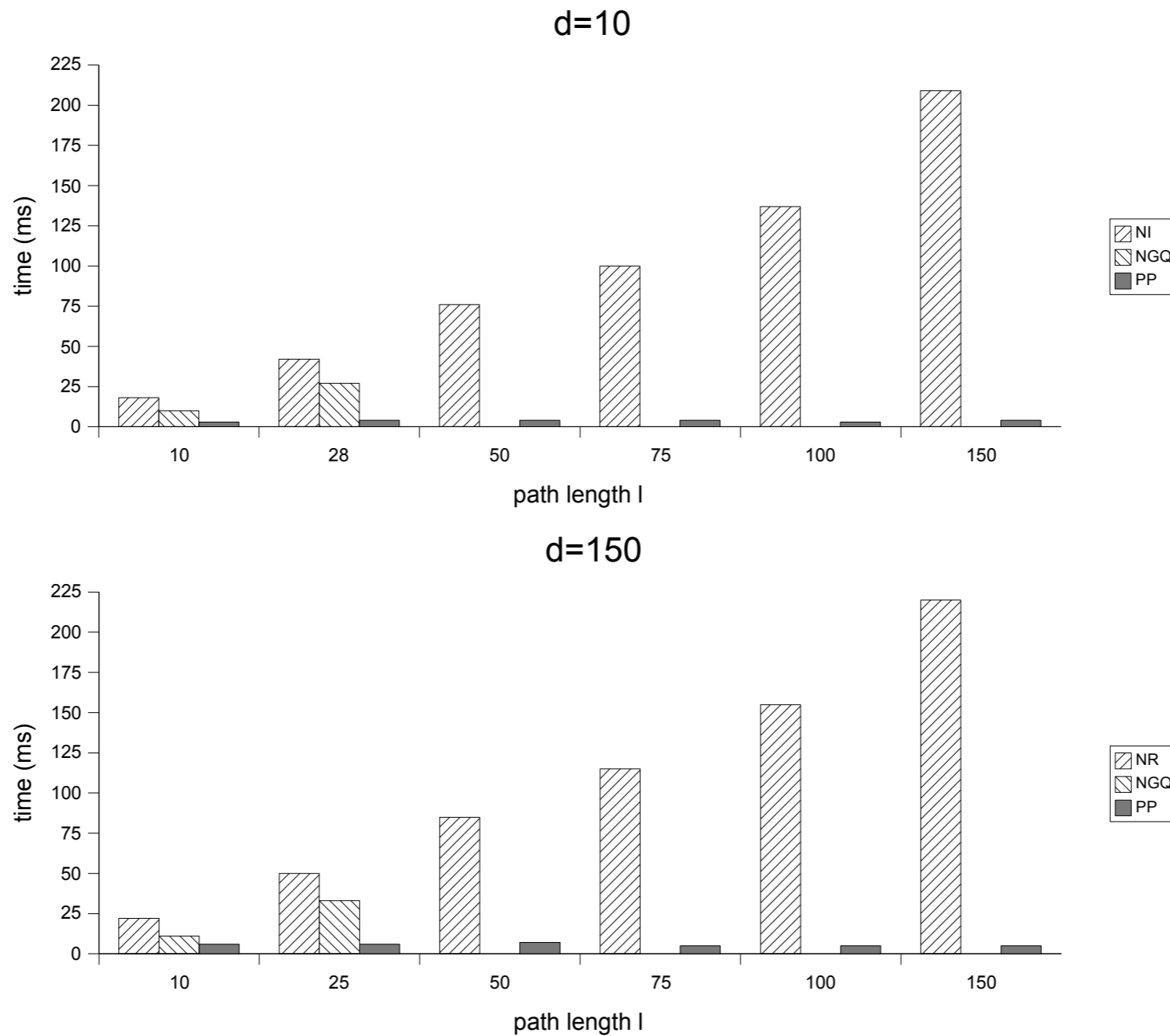
– the "T-towers"

control:
- size of the lists involved
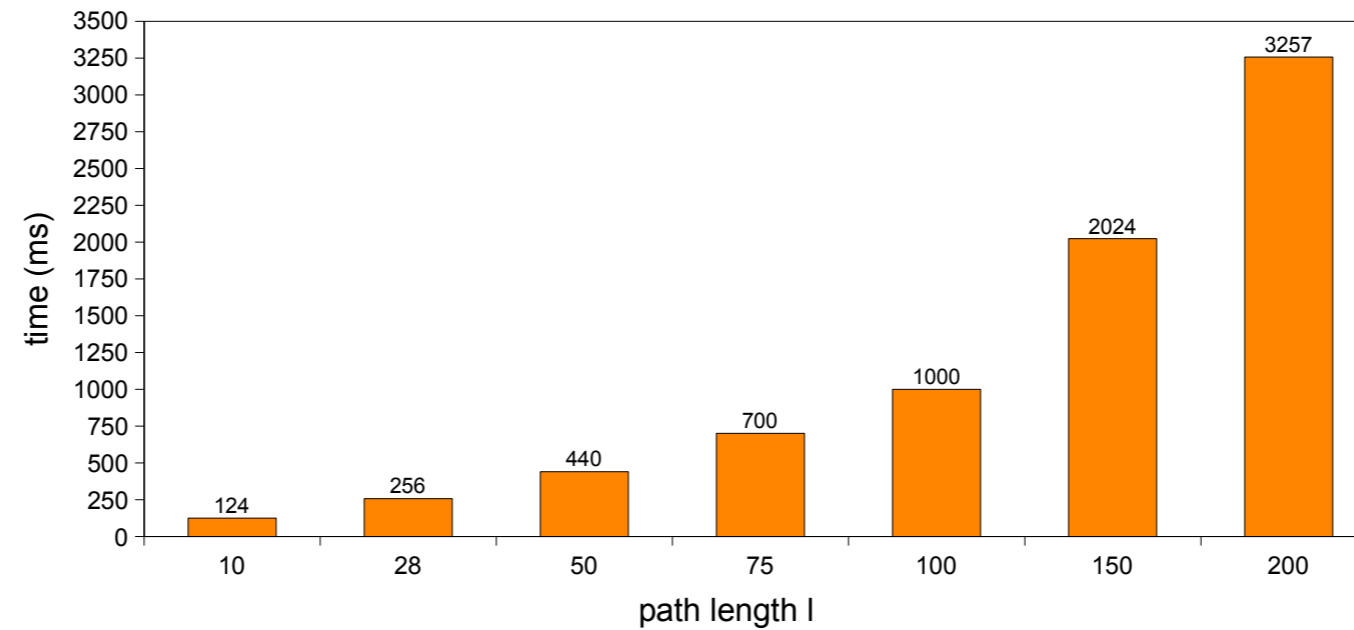- length of the paths
- includes one cross product

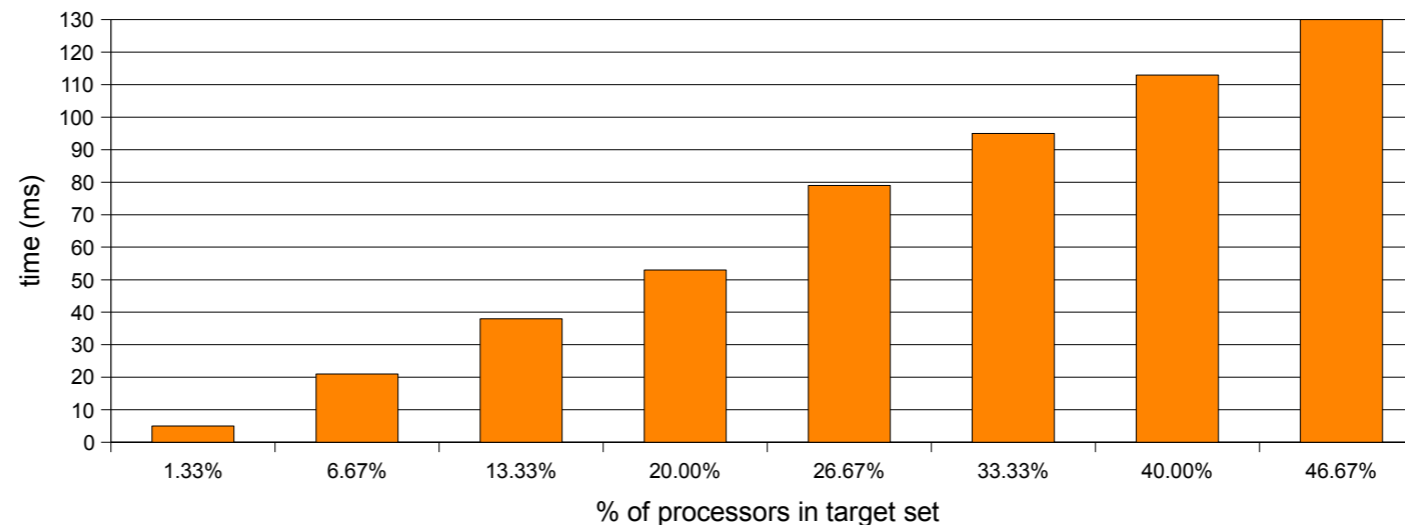- query response time: naive vs. "path projection" approaches

- workflow search time by path length ("tower height")
  - common to all strategies!

**workflow pre-processing time by graph size**



- performance degradation on fully unfocused queries

**response times for PP on unfocused queries  (l=150)**

- An original approach to lineage queries for Taverna that combines
  - efficiency and fine-granularity

- Relies on semantic properties of the Taverna dataflow model

- Further work:
  - visual specification of user query
  - visual presentation of query answer
  - space compression
  - semantic overlays, annotations

- To be bundled with some future version of Taverna...