

Data lineage model for Taverna workflows with lightweight annotation requirements

*Paolo Missier, Khalid Belhajjame, Jun Zhao, Carole
Goble*

*School of Computer Science
The University of Manchester, UK*

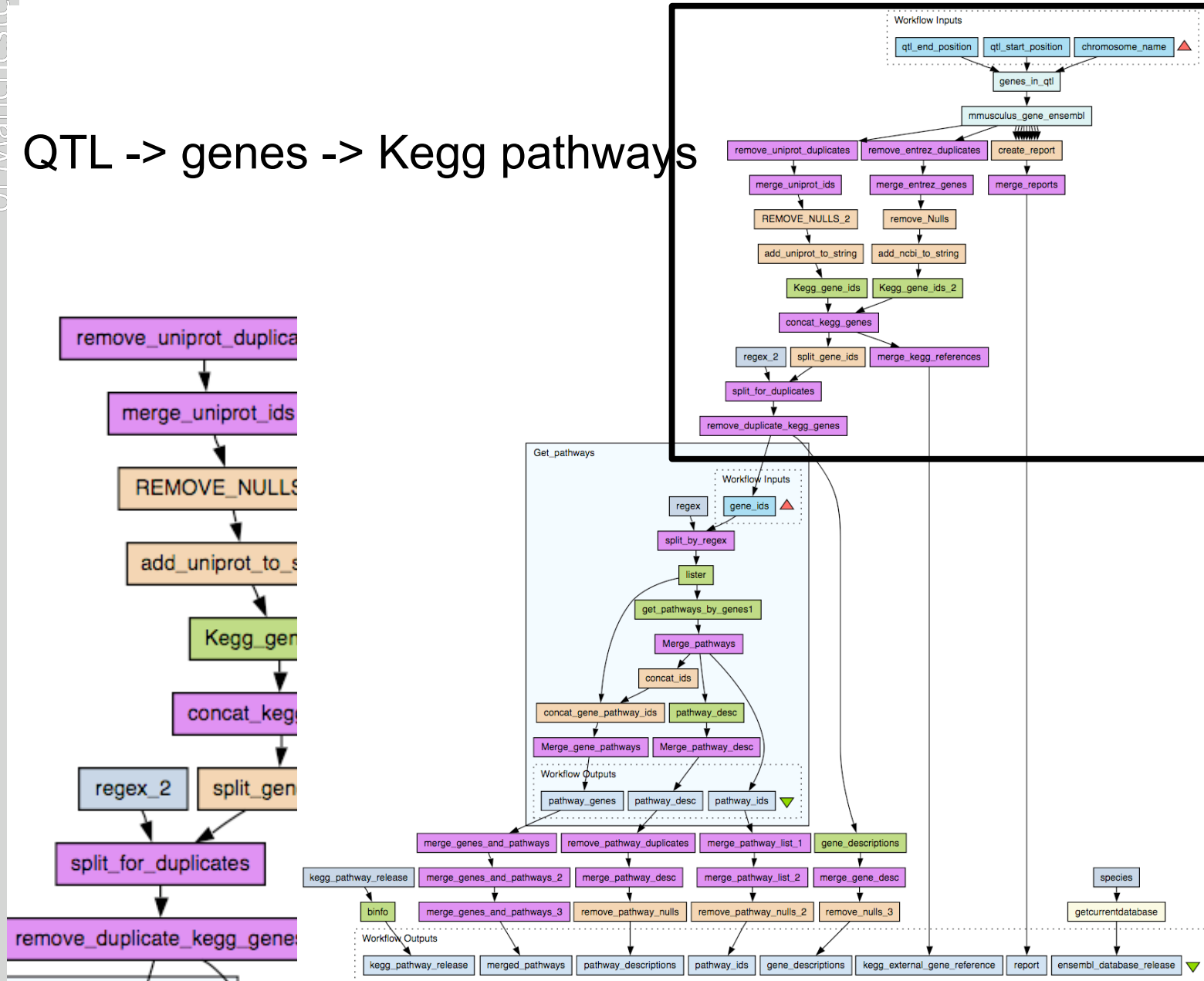
Ongoing work on a new provenance component for Taverna

- myGrid consortium

Scope:

- capture raw provenance events
 - data transformations, data transfers
- store one *lineage graph* for each dataflow execution
- query over single or multiple *lineage graphs*

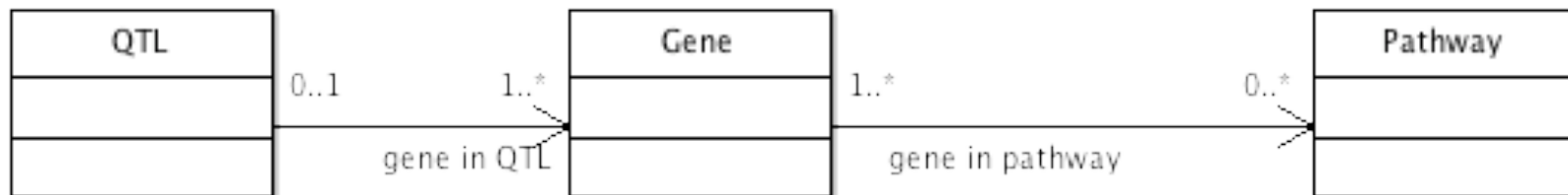
QTL -> genes -> Kegg pathways



Some user questions on lineage

- on a single workflow run:
 - find all genes that participate in some pathway p
 - find all pathways derived from Uniprot genes
 - describe the complete derivation of each pathway in which gene g is involved
- on a collection of runs:
 - find all distinct pathways produced by runs of a dataflow
 - [over a period of time,
produced by a member of my group, ...]

- Granularity
 - risk of returning trivial answers
 - “all outputs depend on all inputs”
- Semantics
 - Results not expressed in the language of the designer
- Abstraction level, noise – the “latent data model”
 - many processors are irrelevant – shims, mundane tasks



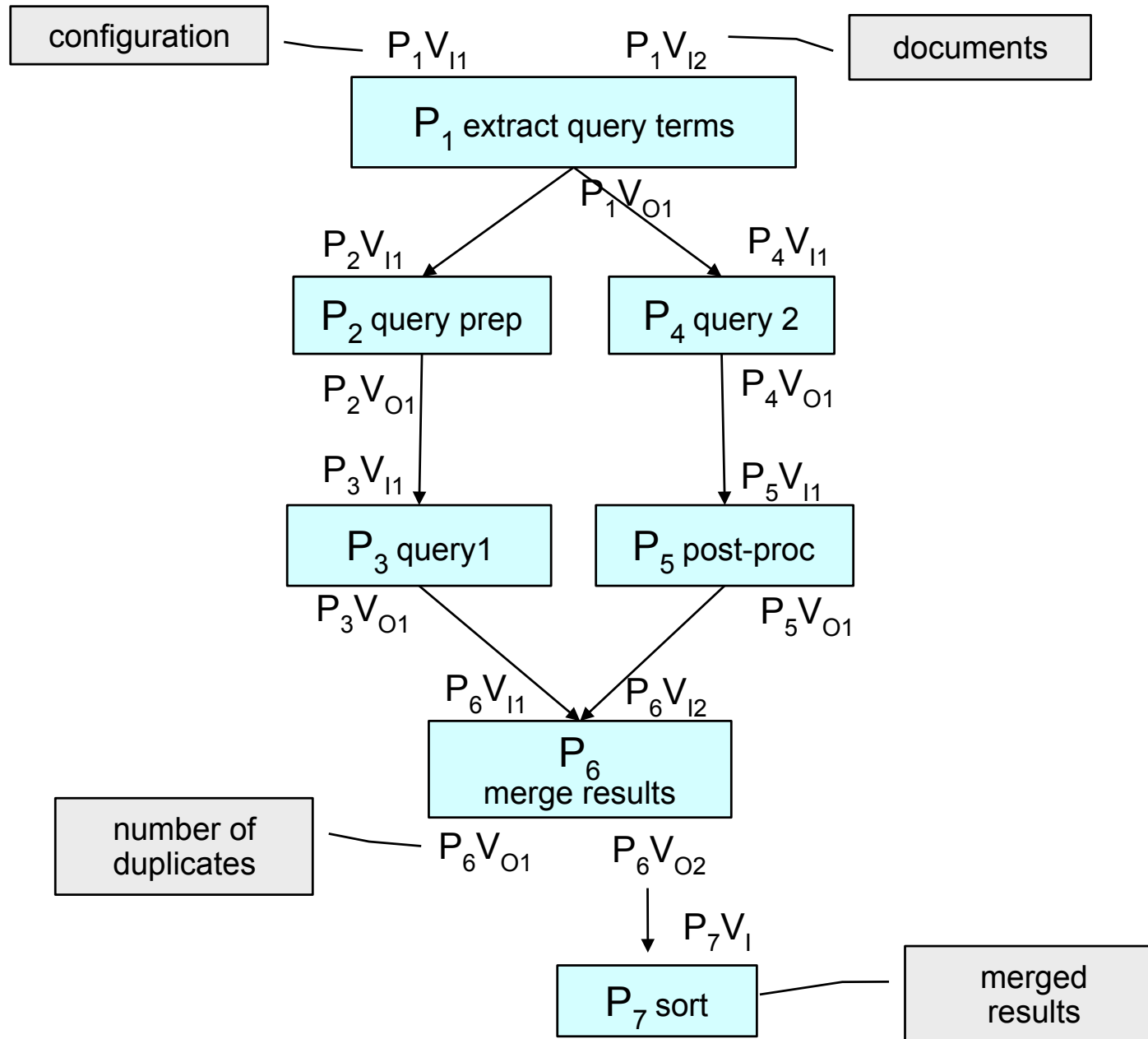
The need for selective annotations

- As long as processors are black boxes, these remain difficult problems
- Adding annotations to processors is tempting

Scope of this work:

to explore the “gray box” region

- simple annotations with minimal semantics
- driving principle: justified by technical benefits
 - precision of query results
 - efficiency of query processing



Two main annotation types

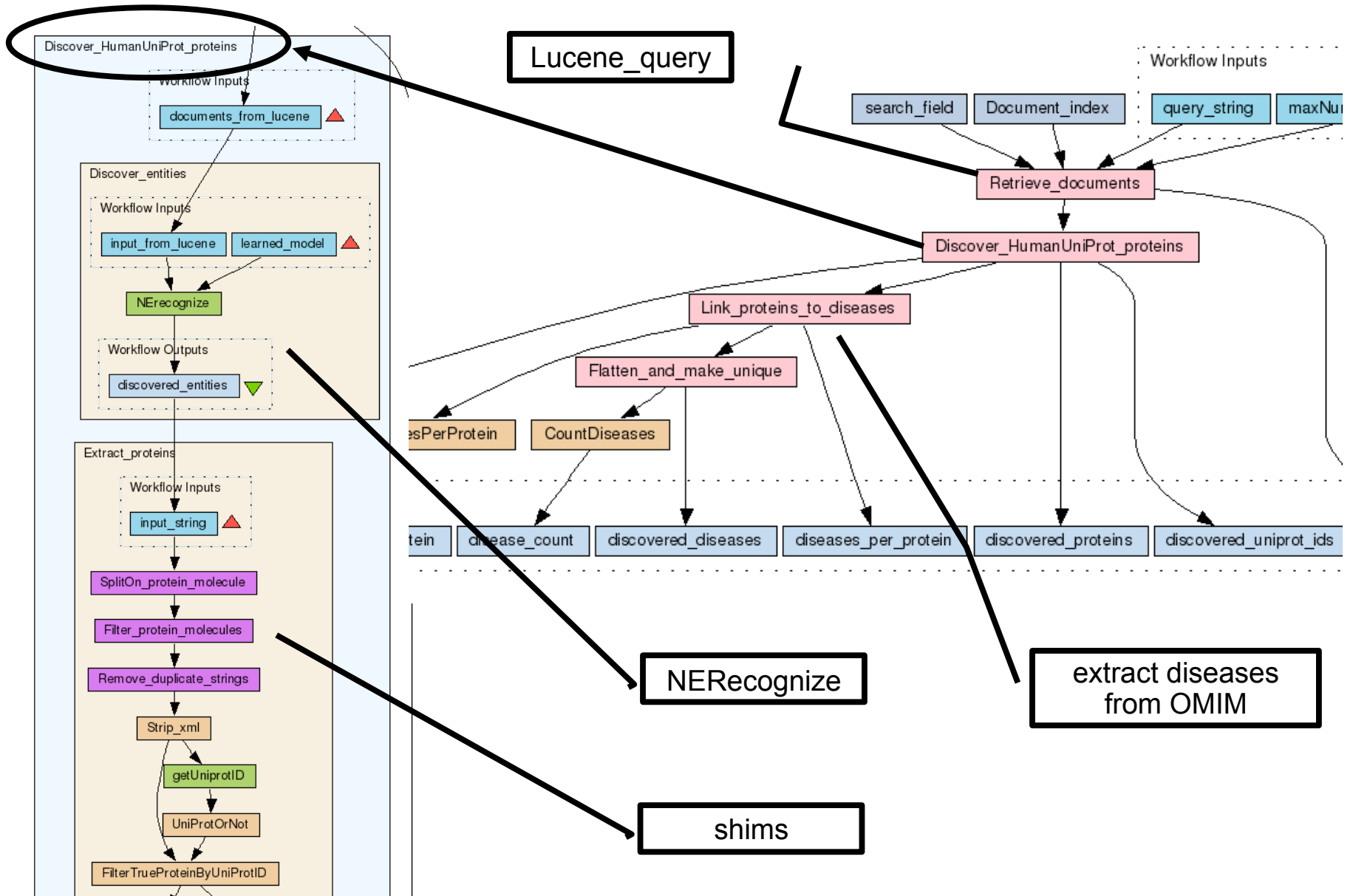
Focusing: processor selection

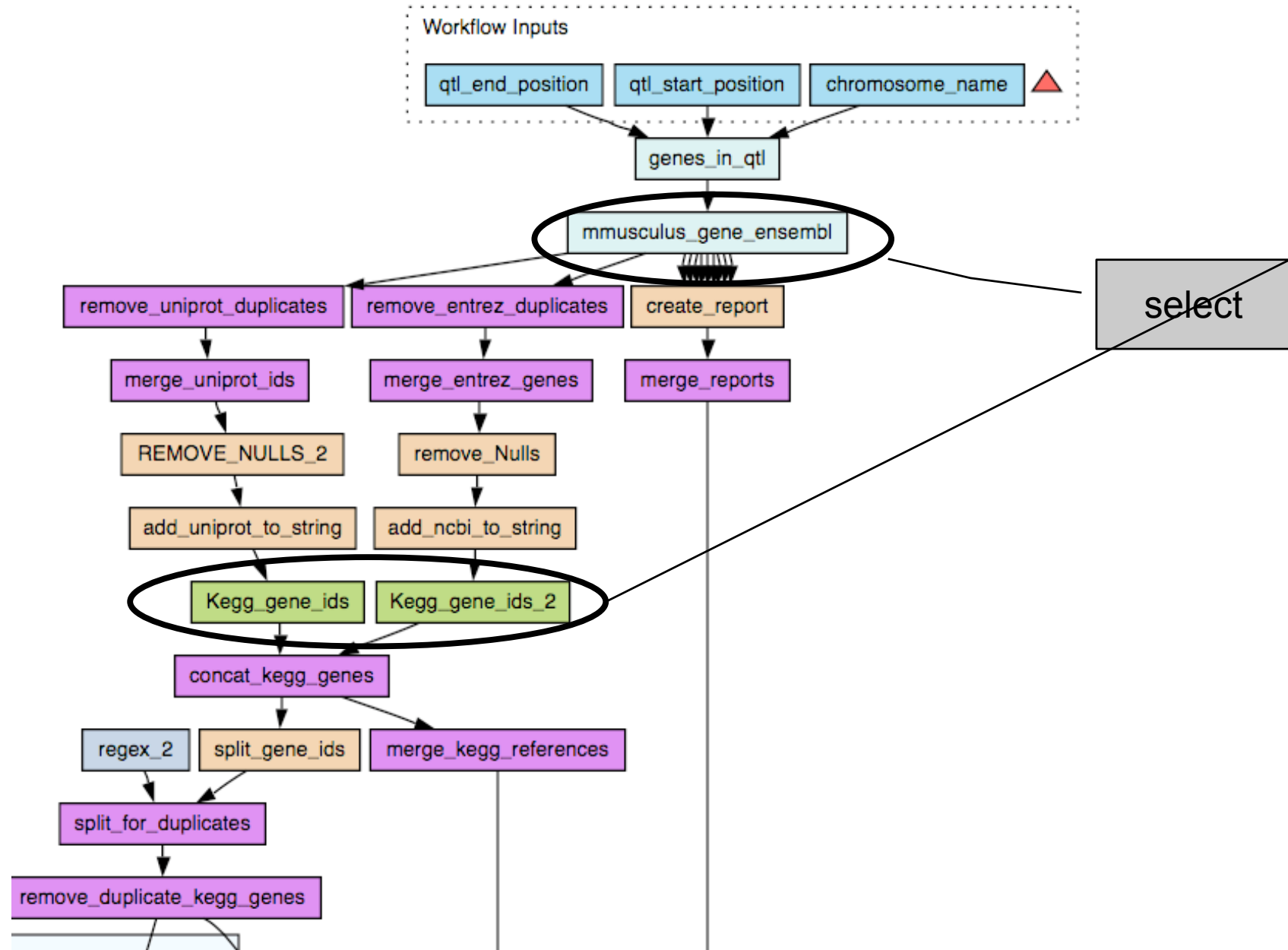
- some processors are more interesting than others
 - “boring” annotations
 - query-time user selection of interesting processors

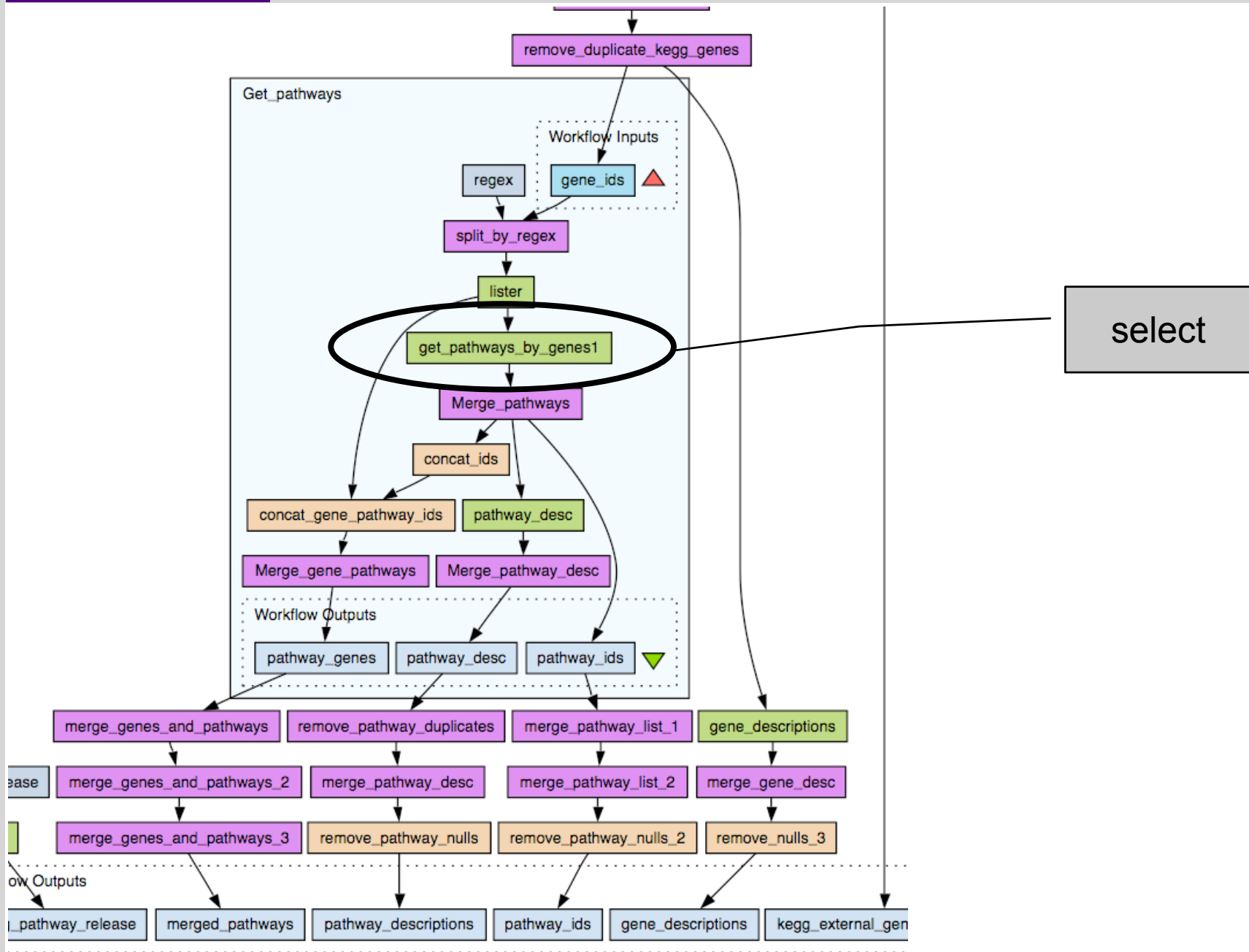
Precision: fine-grained lineage tracing

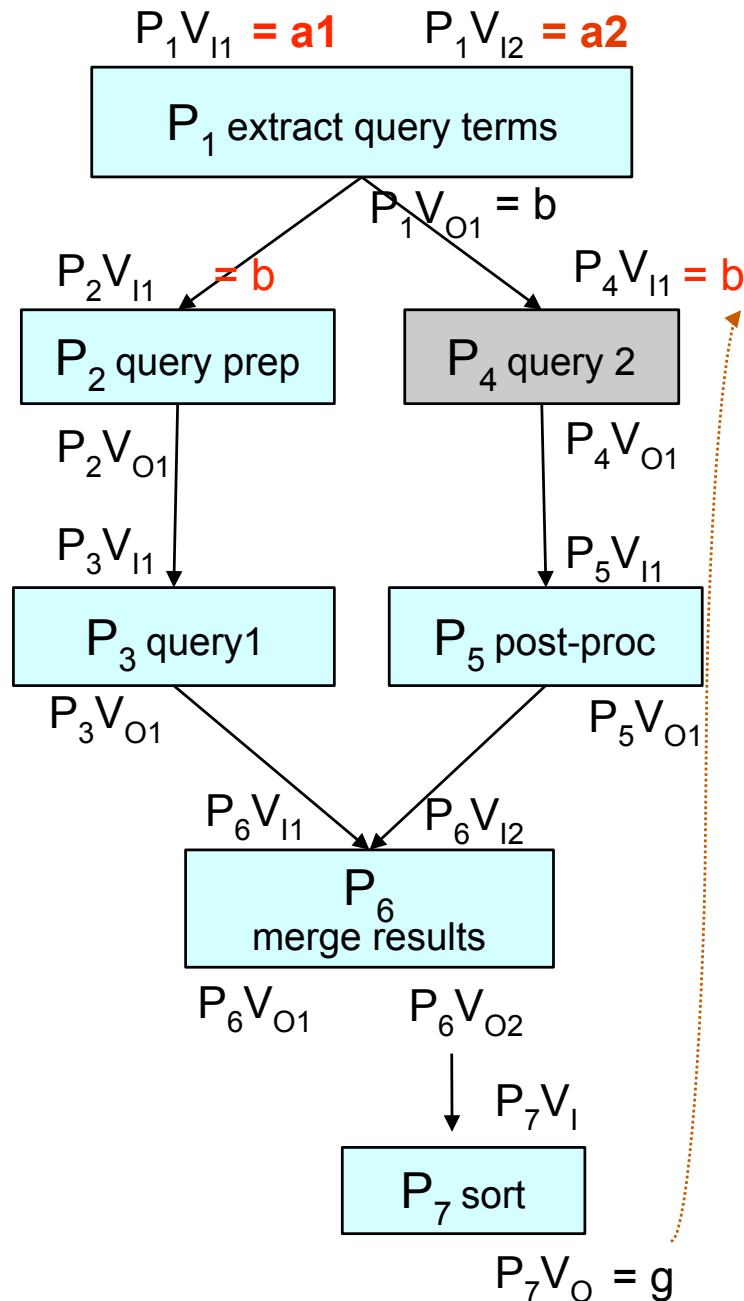
- goal: trace lineage of individual items within a collection

Abstraction by modularization









var	value
P1/V11	a1
P1/V12	a2
P1/V01	b
P2/V11	b
P4/V11	b
P2/V01	c1
P4/V01	c2
P7/V0	g

only interesting processor

Assume all values atomic

Query: $\text{lineage}(P_7V_0, \{P_4\})$

Goal:

- avoid recursive queries on instance tables

Idea:

- use recursion on static model to generate a targeted query
- execute query only once

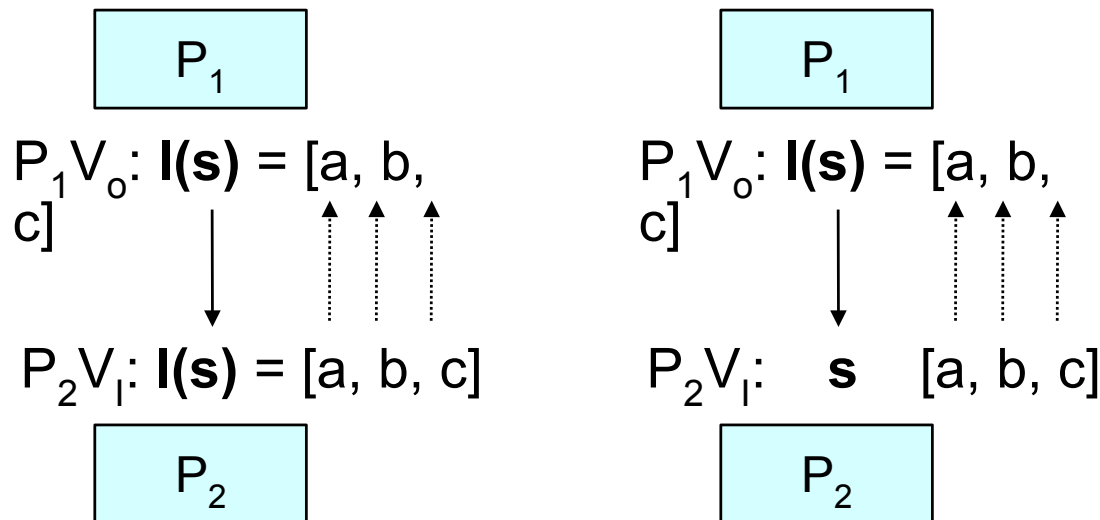
IN	OUT	P
P4/V11	P4/V01	P4
P5/V11	P5/V01	P5
P6/V12	P6/V02	P6
P7/V1	P7/V0	P7

Problem: `xform()` also applies to list values

- It may be impossible to trace individual elements
 - “which pathways (out) depend on which genes (in)”?

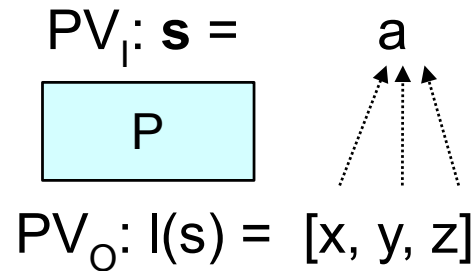
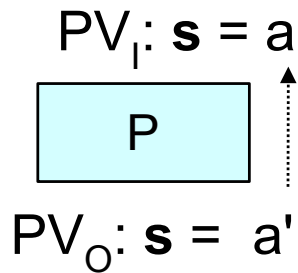
Goal: extend the query generation idea just sketched to trace element-level lineage within collections

Approach: exploit static typing of Taverna processors

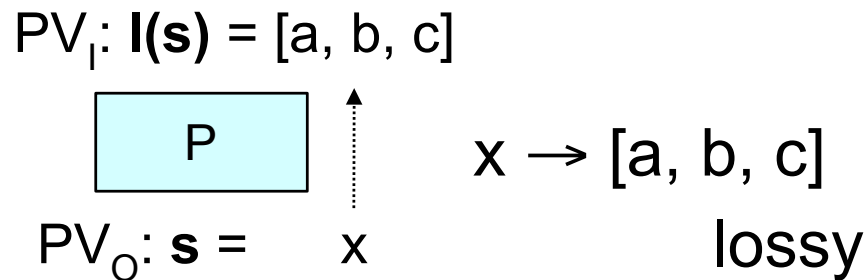


Taverna resolves mismatches on nesting levels:
`(map P2 [a,b,c])`

Loss of precision in transformations



“lossless”
transformations



possible behaviours:

- selection of an element
- **aggregation function**

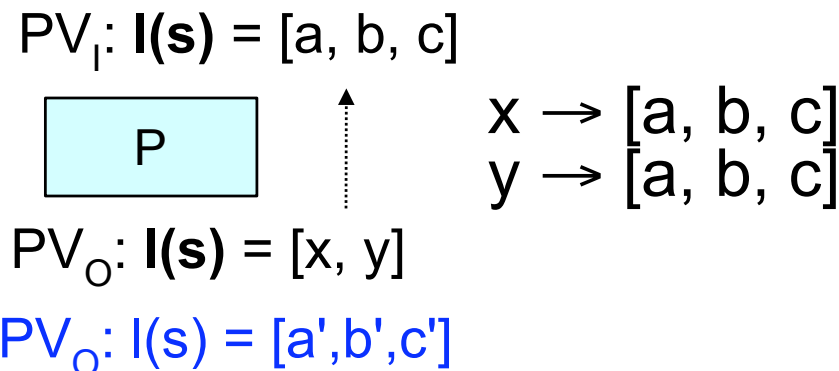
c
tion f() useful annotation:

$lineage(PV_0) = f(PV_1)$
only useful annotation:

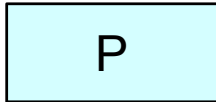
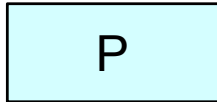
P is **index-preserving**:

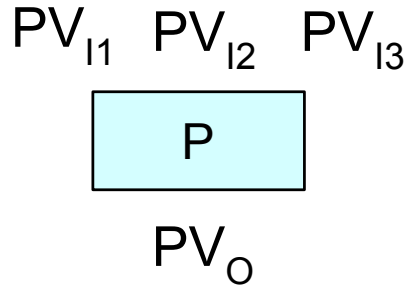
$$PV_0[i] = PV_1[i]$$

$$lineage(PV_0[i]) = PV_1[i]$$



- *Passive* processors do not contribute explicit provenance info
- **Cooperative** processors actively feed metadata to the lineage service

	$PV_1: I(\mathbf{s}) = [a, b, c]$  $PV_0: \mathbf{s} = x$	$PV_1: I(\mathbf{s}) = [a, b, c]$  $PV_0: I(\mathbf{s}) = [x, y]$
Static annotations:	aggregation $f()$	$PV_0[i] = PV_1[i]$
Dynamic annotations:	selection: $x = PV_1[i]$	sorting: $PV_0 = \Pi(PV_1)$



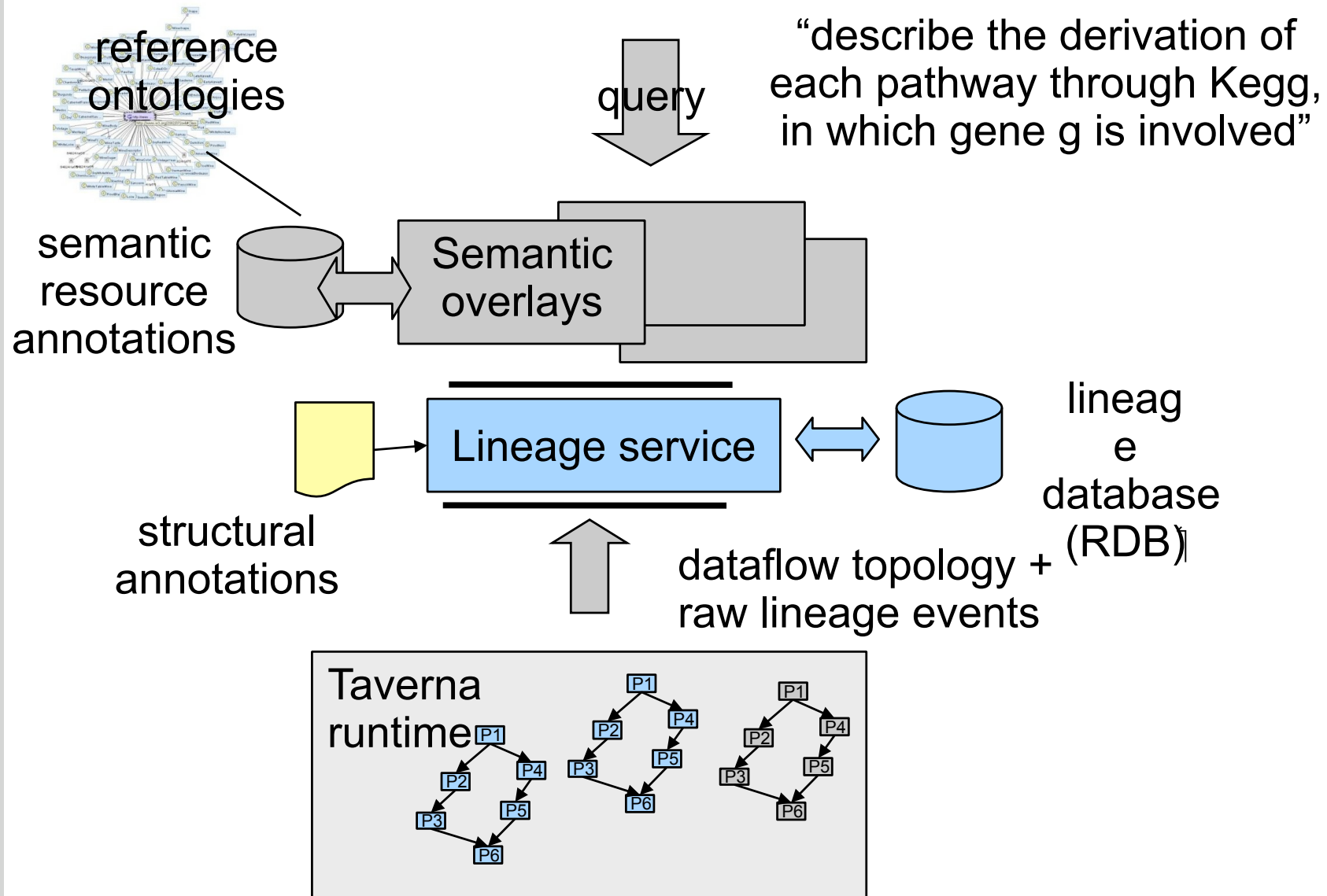
- Distinction between configuration and input data
 - PV_{I3} is a configuration parameter
 - compare effect of different config. across multiple runs

- specific functional dependencies
 - $[PV_{I1}, PV_{I2}] \rightarrow PV_O$

- stateless processor
 - execute process \leftrightarrow retrieve provenance

More evaluation needed on these

Towards a 2 tier provenance model



A data lineage model for Taverna workflows

- Raw lineage data has shortcomings
- A few, selected lightweight annotations added in a principled way
 - win-win:
 - helpful to users
 - *and* enable query optimization
- Form the base layer in a broader approach to efficient querying of semantic provenance for e-science
- Ongoing implementation