

Provenance of workflow data products

Paolo Missier
School of Computing,
Newcastle University, UK

TAPP'11 workshop
Heraklion, Crete
June 20-21, 2011

Taverna type system:

- strings + nested lists
- “cat”, [“cat”, “dog”], [[“cat”, “dog”], [“large”, “small”]]

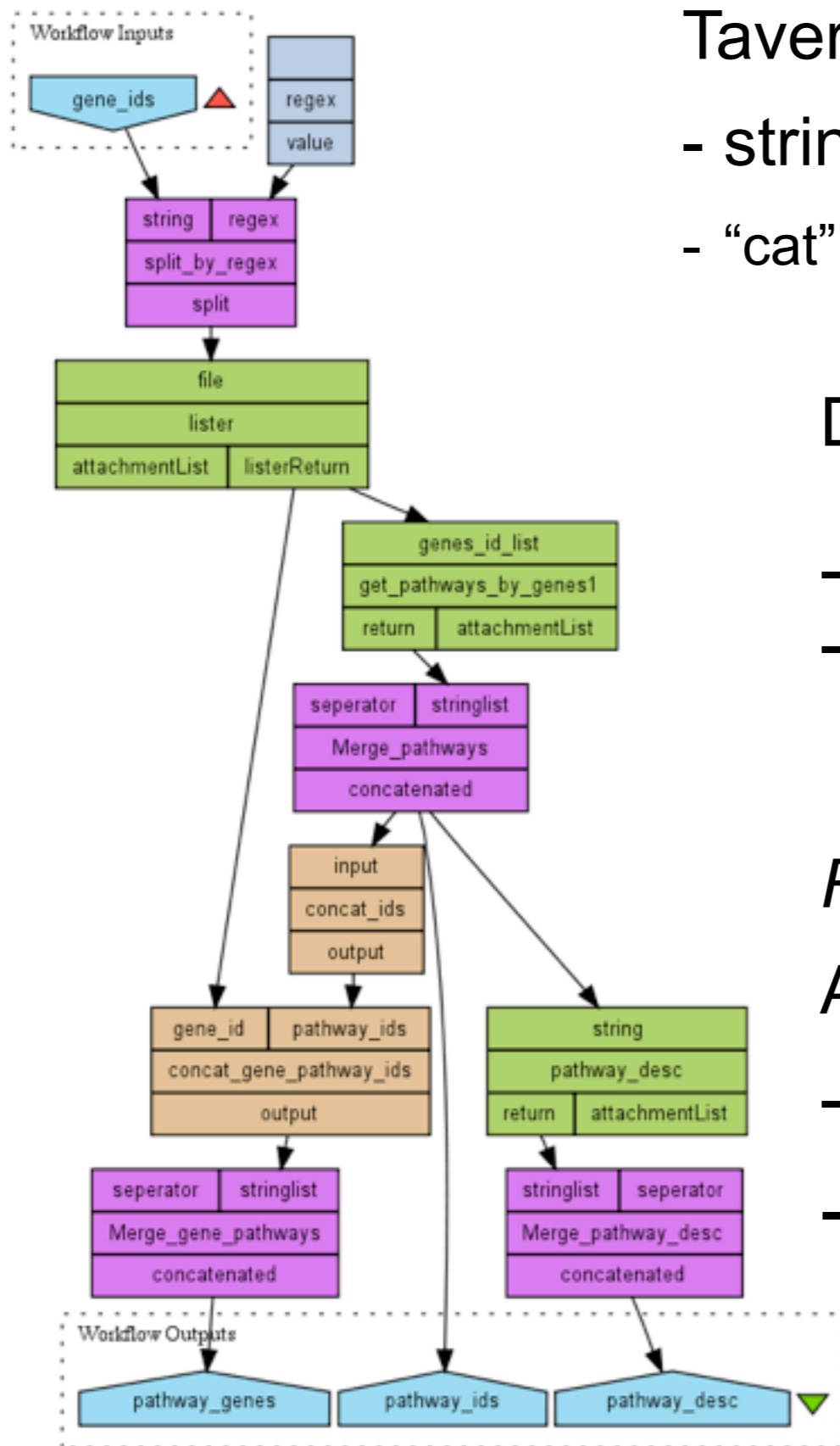
Dataflow model:

- data-driven execution
- services activate when input is ready

Raw provenance:

A detailed trace of workflow execution

- tasks performed, data transformations
- inputs used, outputs produced



Taverna type system:

- strings + nested lists
- "cat", ["cat", "dog"], [["cat", "dog"], ["large", "small"]]

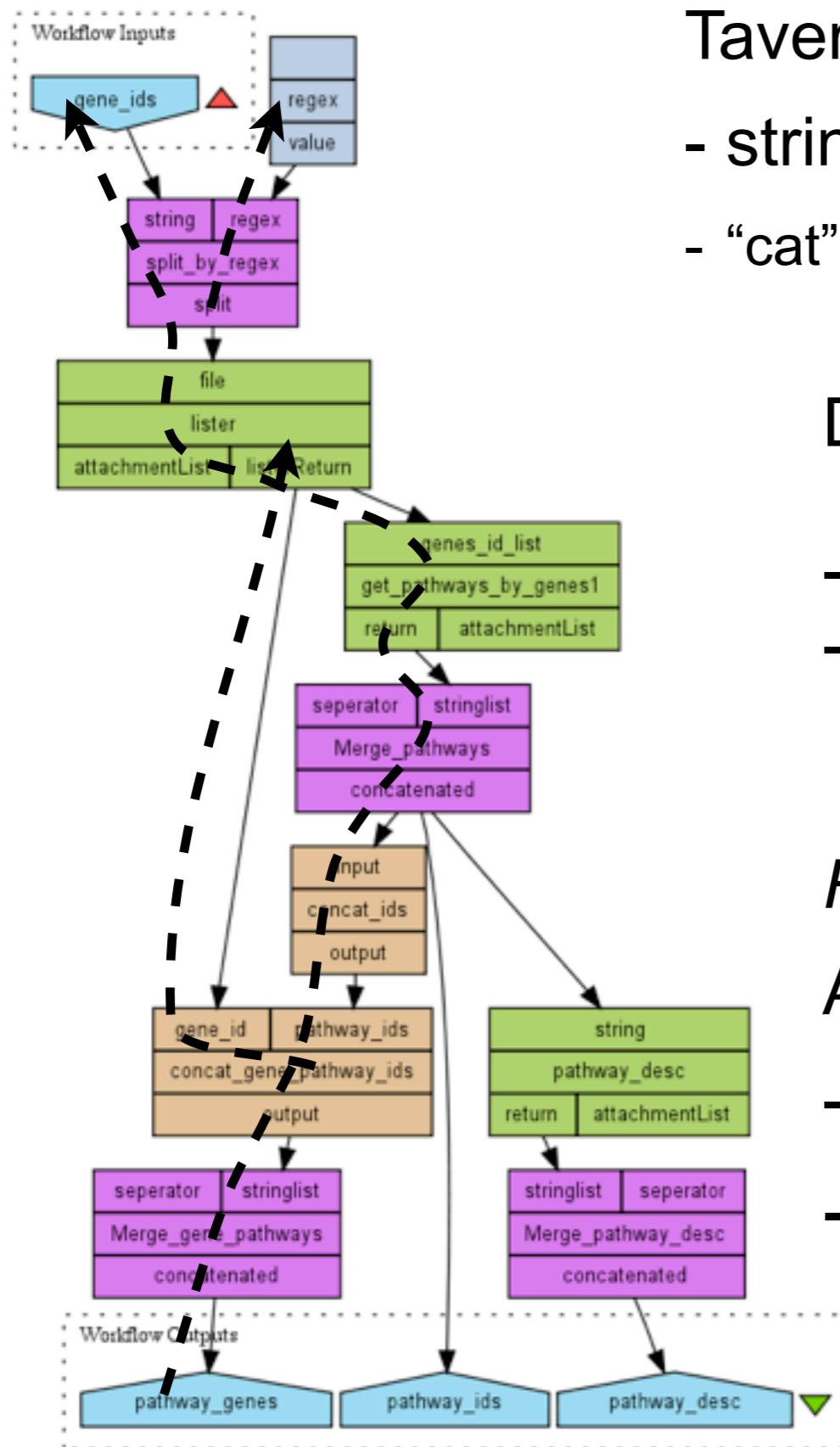
Dataflow model:

- data-driven execution
- services activate when input is ready

Raw provenance:

A detailed trace of workflow execution

- tasks performed, data transformations
- inputs used, outputs produced



Taverna type system:

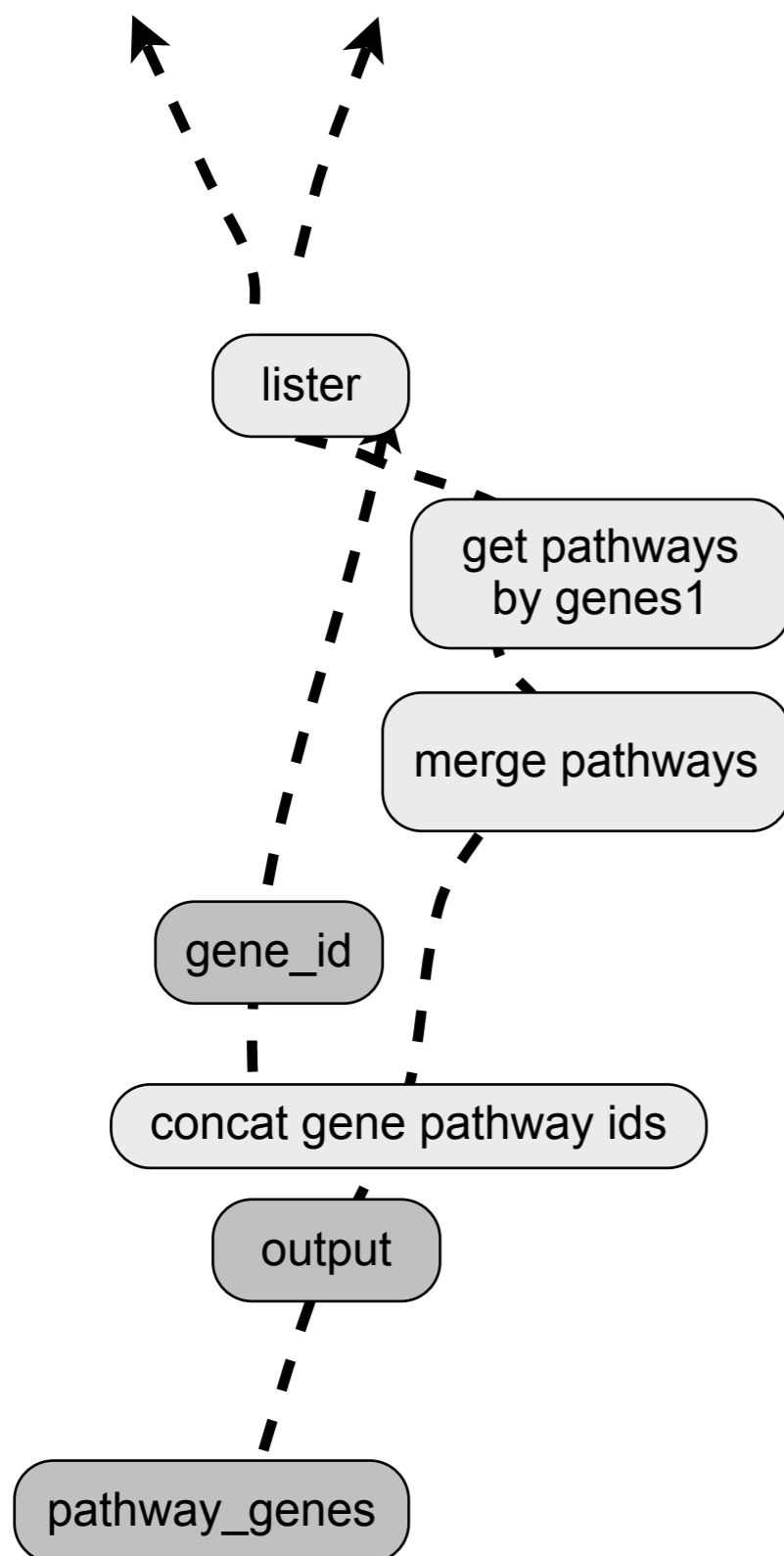
- strings + nested lists
- “cat”, [“cat”, “dog”], [[“cat”, “dog”], [“large”, “small”]]

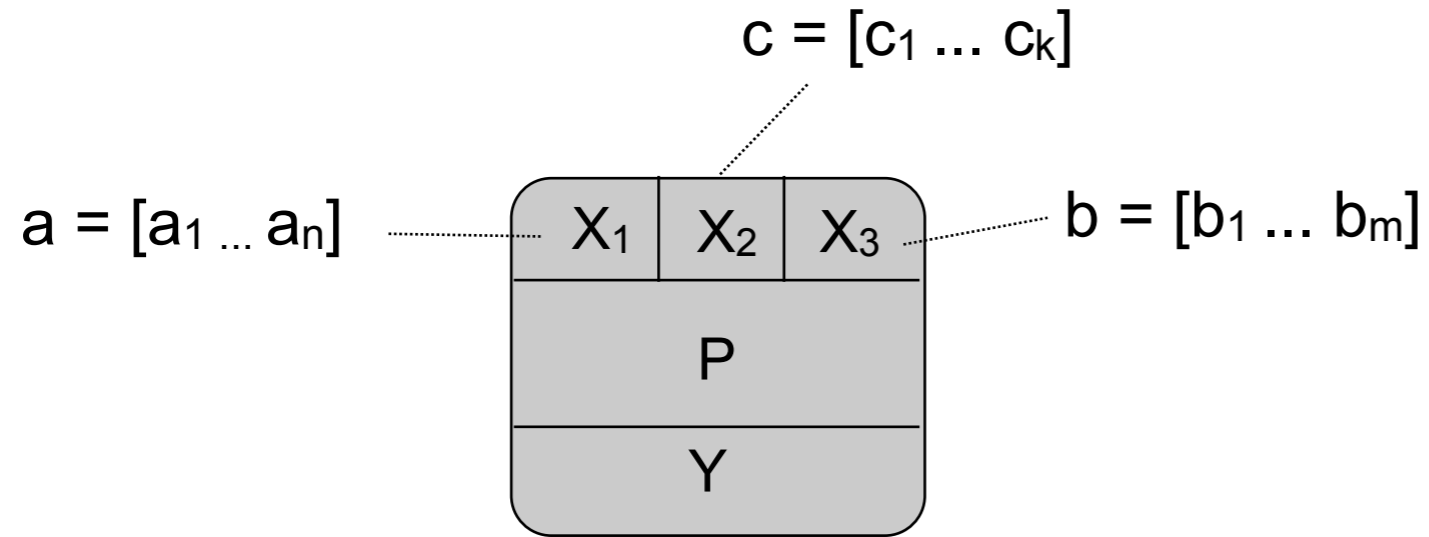
Dataflow model:

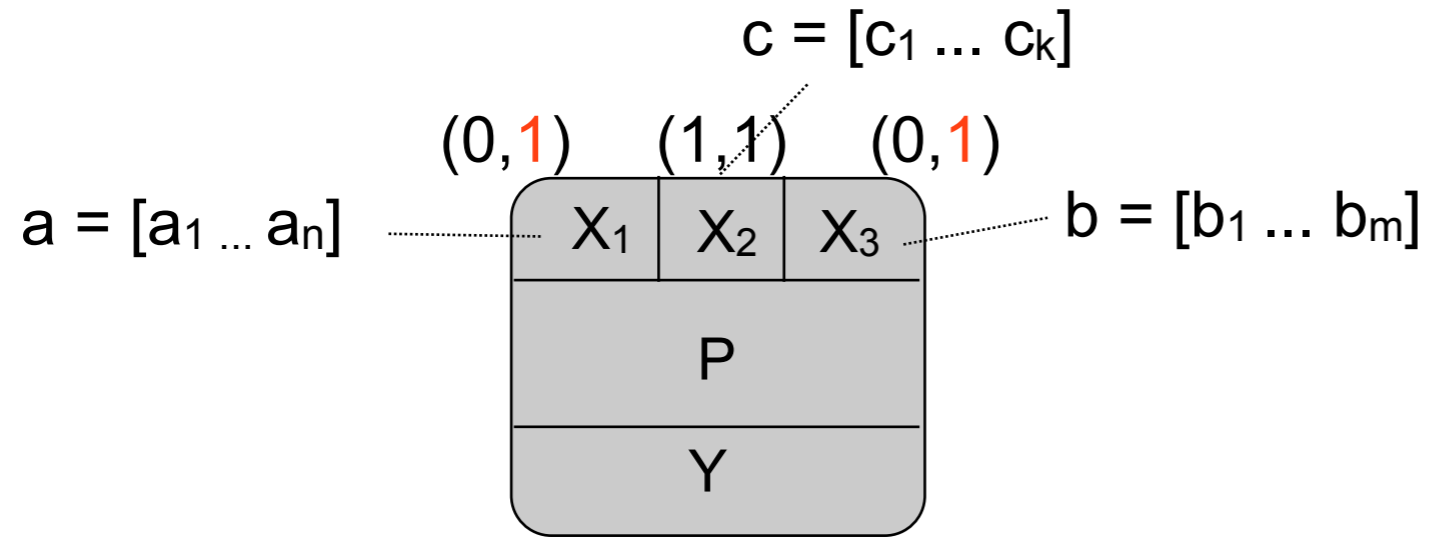
- data-driven execution
- services activate when input is ready

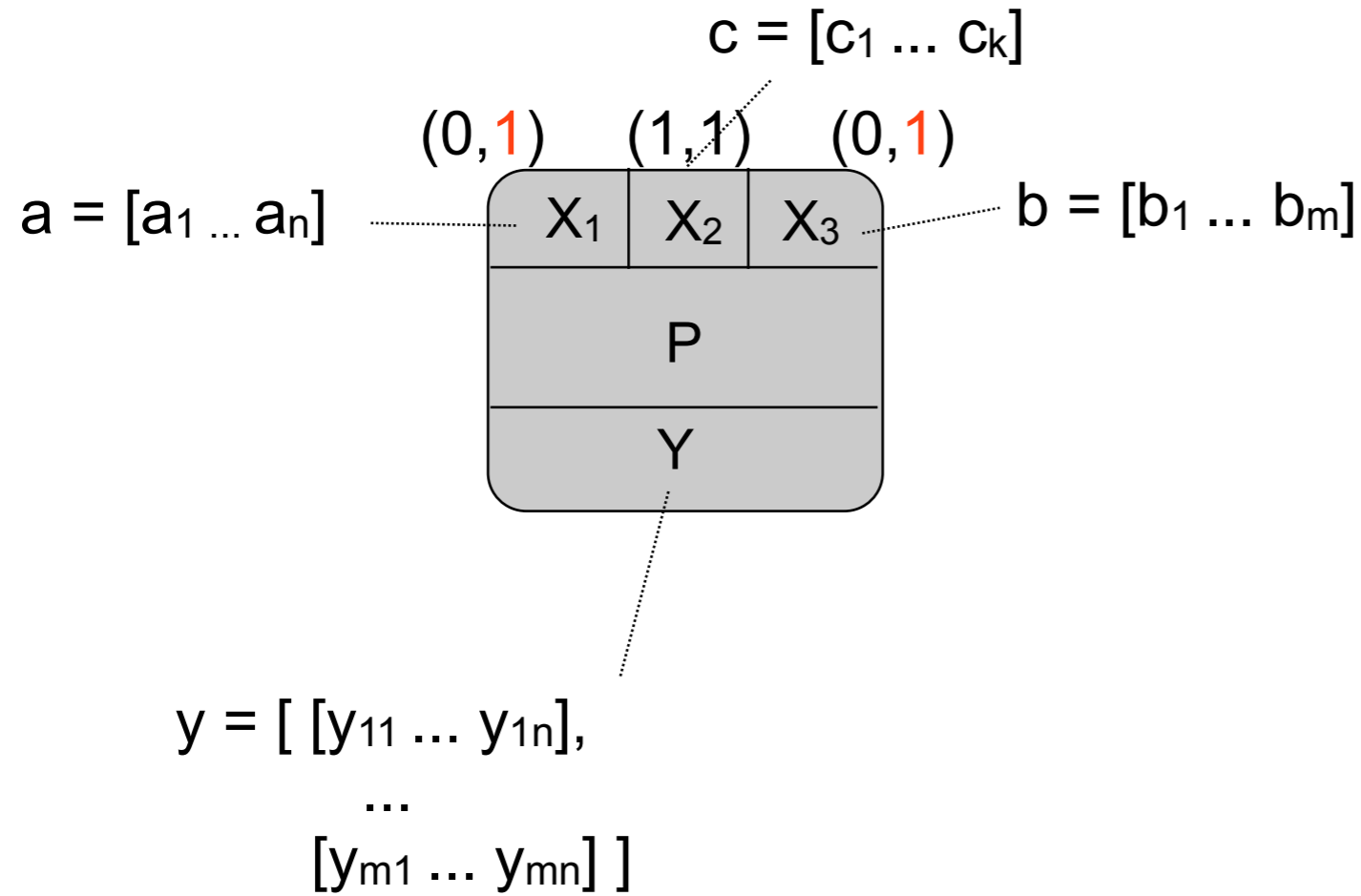
Raw provenance:

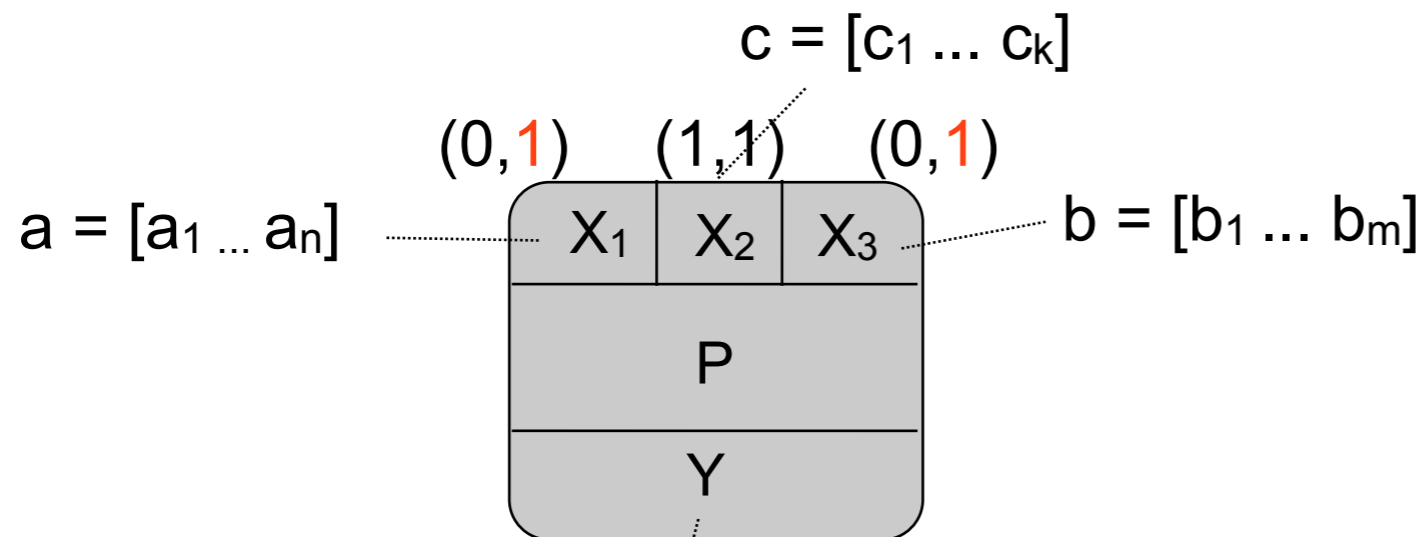
- A detailed trace of workflow execution
- tasks performed, data transformations
 - inputs used, outputs produced











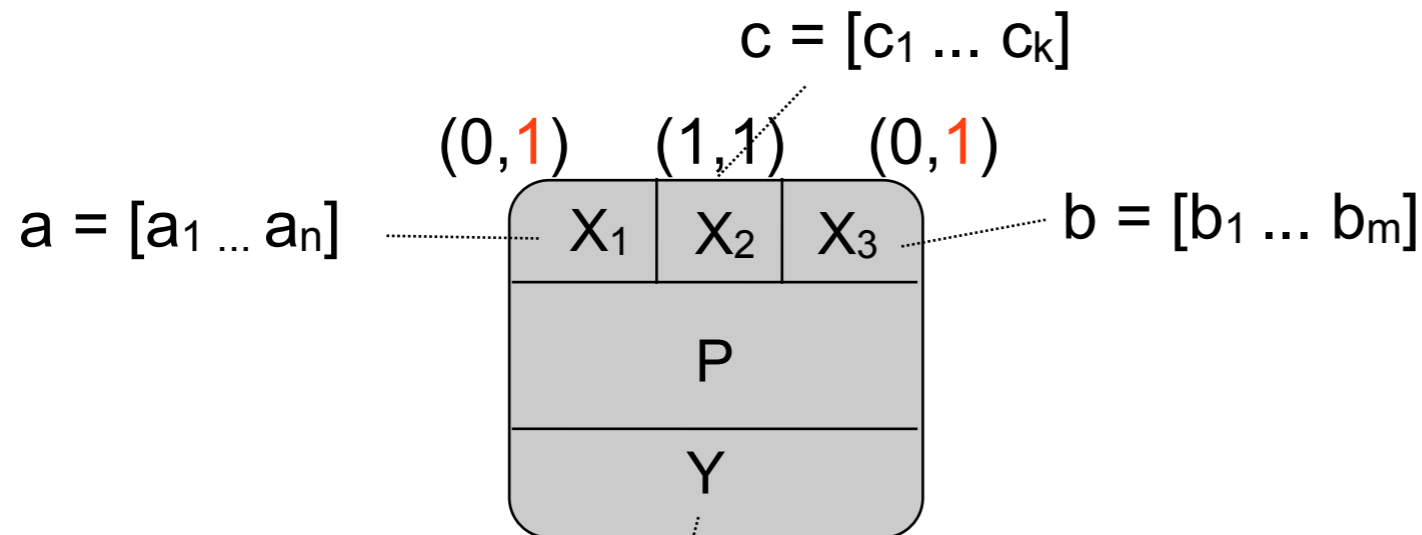
$$y = [[y_{11} \dots y_{1n}], \dots [y_{m1} \dots y_{mn}]]$$

How y is computed at P :

let $l = a \otimes b = [[\langle a_i, b_j \rangle \mid b_j \in b] \mid a_i \in a]$ // cross product

$l' = [[\langle a_i, c, b_j \rangle \mid b_j \in b] \mid a_i \in a]$ // same product but with c interleaved

$$y = (\text{map} (\text{map} P) l') = [(\text{map} P [\langle a_1, c, b_1 \rangle \dots \langle a_1, c, b_m \rangle]), \dots, (\text{map} P [\langle a_n, c, b_1 \rangle \dots \langle a_n, c, b_m \rangle])] = [[y_{11} \dots y_{1n}], \dots [y_{n1} \dots y_{nm}]]$$



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

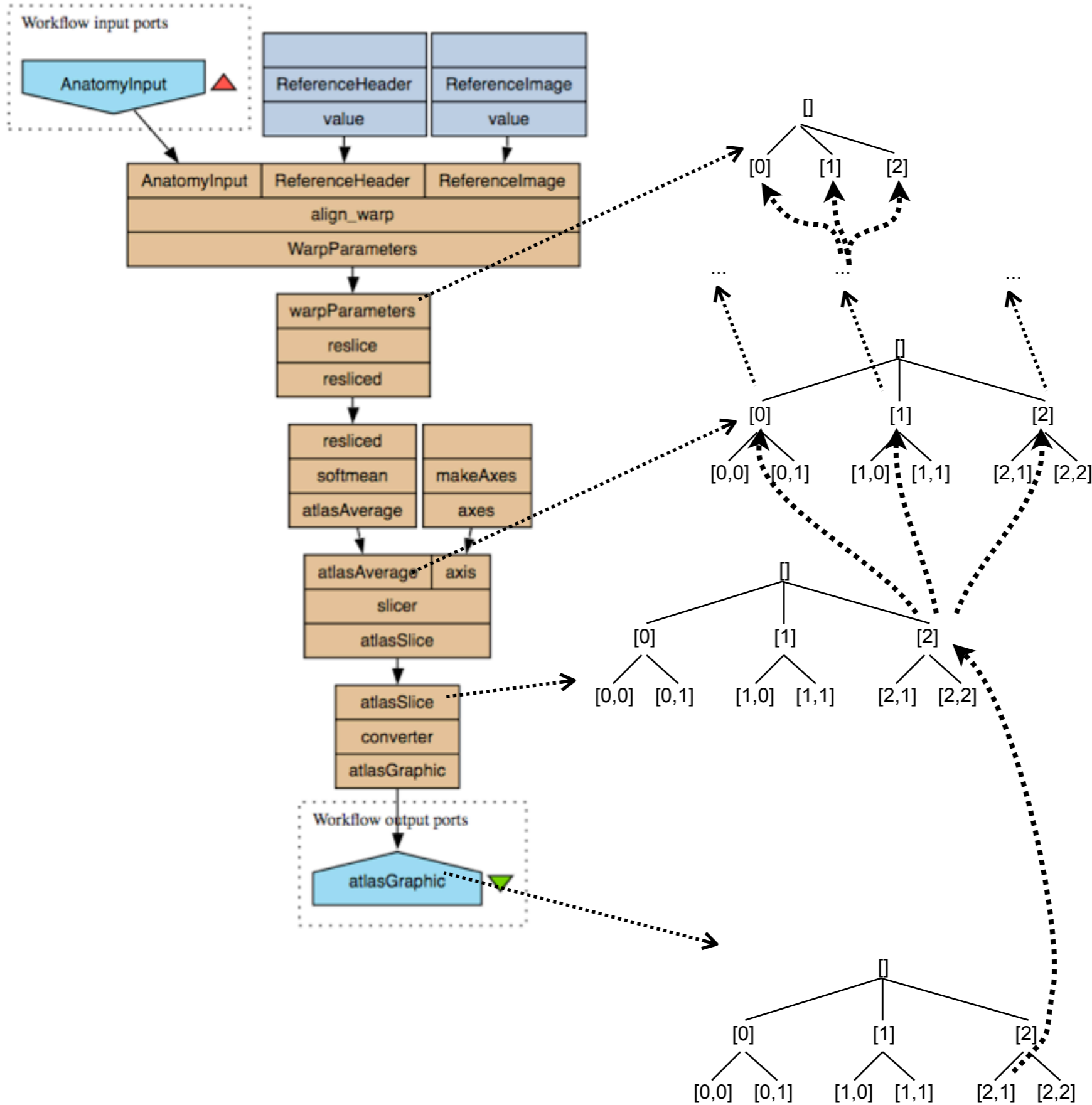
bottom line:
y_{ij} depends only on values a_i, c, b_j

How y is computed at P:

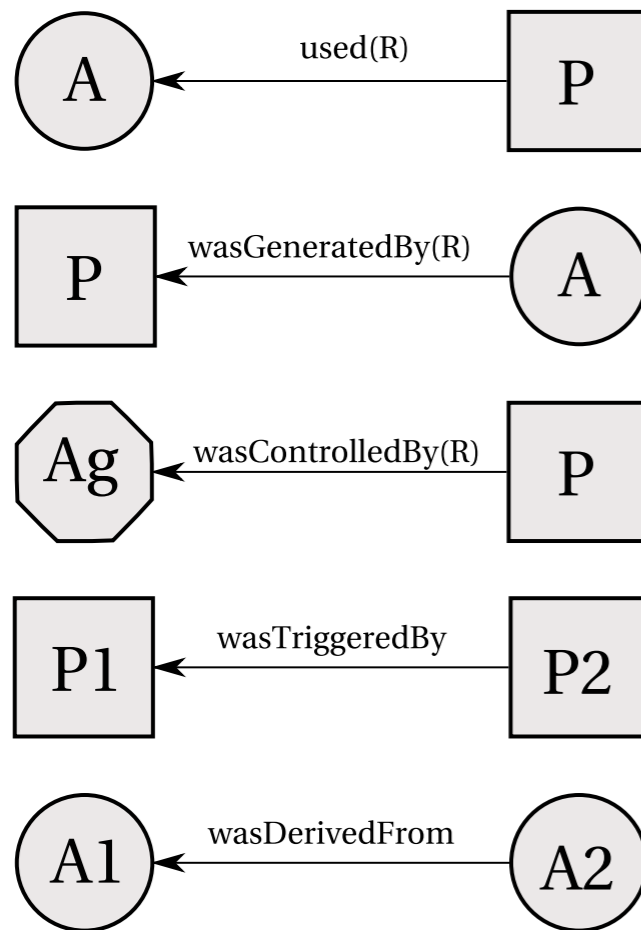
let $l = a \otimes b = [[\langle a_i, b_j \rangle \mid b_j \in b] \mid a_i \in a]$ // cross product

$l' = [[\langle a_i, c, b_j \rangle \mid b_j \in b] \mid a_i \in a]$ // same product but with c interleaved

$$y = (\text{map} (\text{map} P) l') = [(\text{map} P [\langle a_1, c, b_1 \rangle \dots \langle a_1, c, b_m \rangle]), \dots, (\text{map} P [\langle a_n, c, b_1 \rangle \dots \langle a_n, c, b_m \rangle])] = [[y_{11} \dots y_{1n}], \dots [y_{n1} \dots y_{nm}]]$$



- Let's look at the Open Provenance Model as a starting point

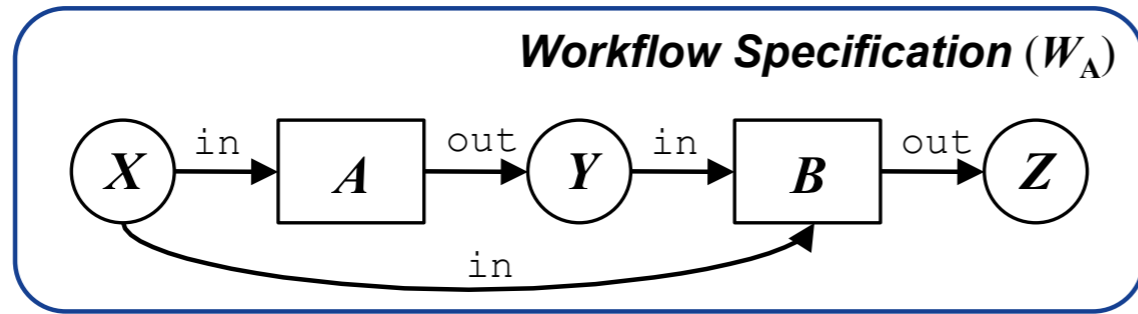


Core OPM

- agnostic wrt **Artifact, Processor types**
- roles: annotations on binary relations
- extensions by subclassing
 - node types,
 - relation types

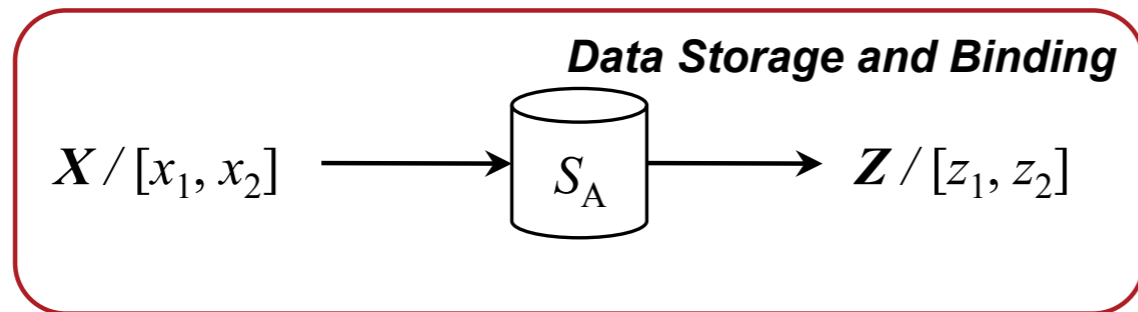
Formal (temporal) semantics hopefully available soon

Single User's View (Alice)



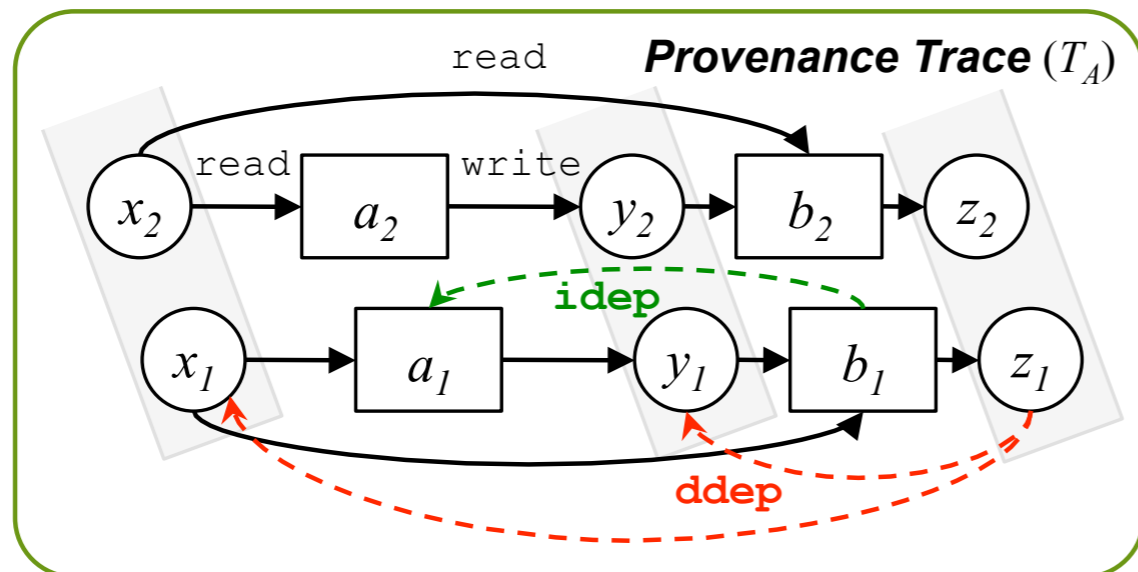
Process space

↓ Data Binding & Enactment



Data space

↓ Execution & Trace Capture



Provenance space

$\text{ans}(X) \text{ :- ddep}^*(X, z_1).$ **Provenance Queries**

- **read**, **write** are natural observables for a workflow run
- possible additional relations (recorded or inferred):

• **invocation dependencies**: $a_1 \overset{\text{idep}}{\leftarrow} a_2$

Explicit or via: $\text{idep}(a_1, a_2) :- \text{write}(a_1, d), \text{read}(d, a_2)$

“ a_2 depends on a_1 ” because a_1 has written data d , a_2 has read d

• **data dependencies**: $d_1 \overset{\text{ddep}}{\leftarrow} d_2$

Explicit or via: $\text{ddep}(d_1, d_2) :- \text{read}(d_1, a), \text{write}(a, d_2)$

“ d_2 depends on d_1 ”

... because some **actor invocation** a read d_1 prior to writing d_2

- Closure queries:
- operate on the transitive closure $ddep^*$ over $ddep$:

$$ddep^*(d_1, d_2) :- ddep(d_1, d_2)$$

$$ddep^*(d_1, d_2) :- ddep(d_1, d), ddep^*(d, d_2)$$

But also:

- queries on the workflow structure
- queries on the data structures (e.g. collections)

and importantly:

use workflow graphs to justify/explain the provenance graph for one workflow run:

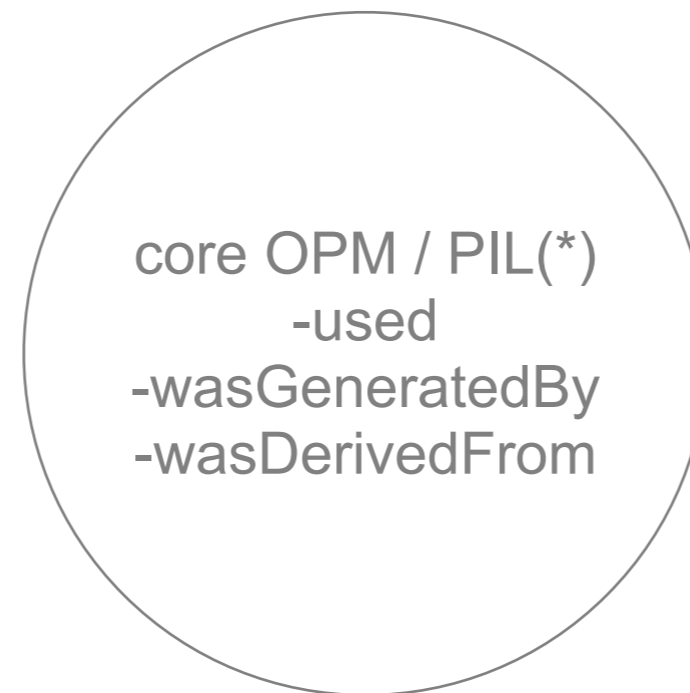
T_A *trace instance of* W_A :

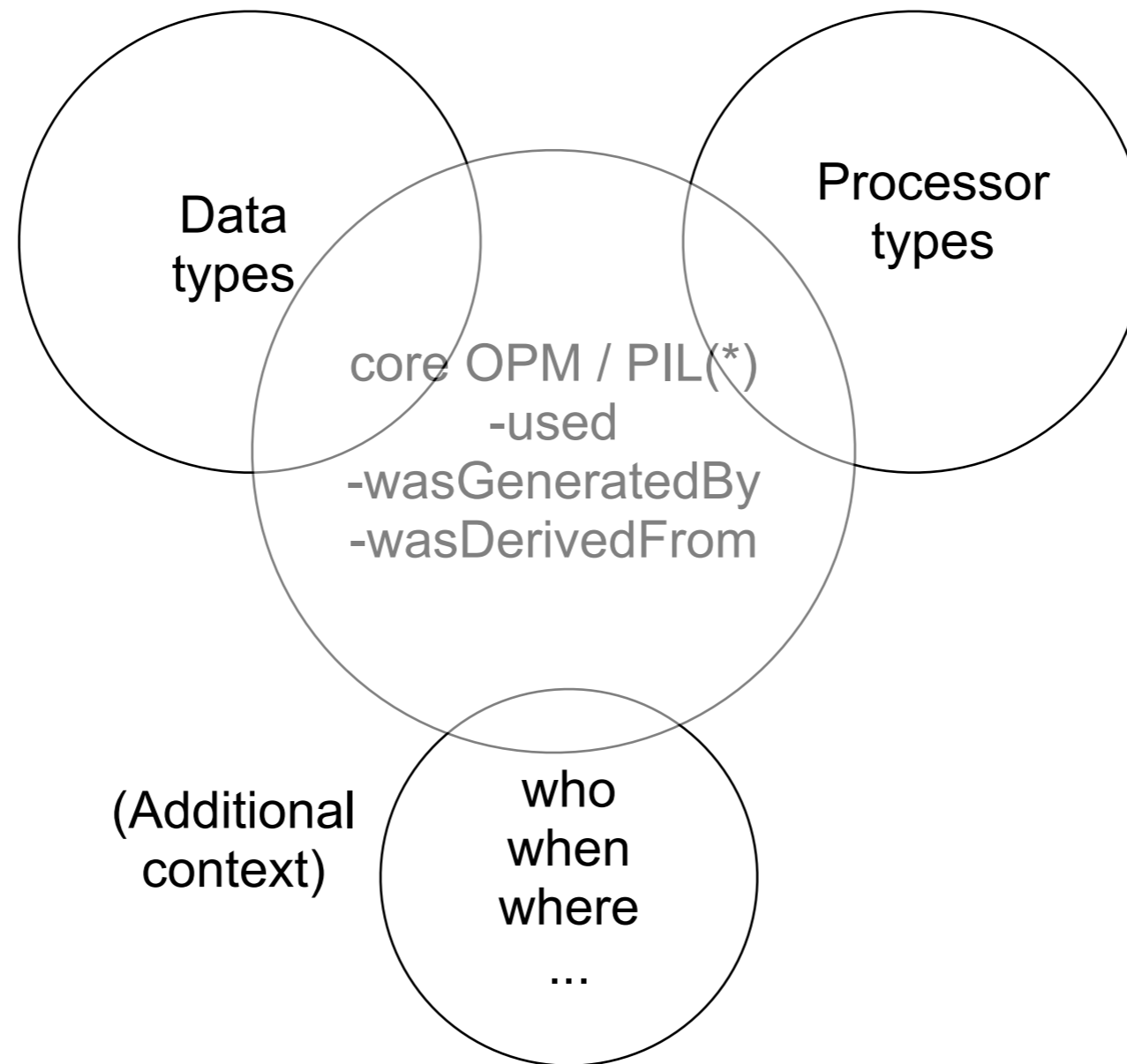
$h: T_A \rightarrow W_A$ *homomorphism*

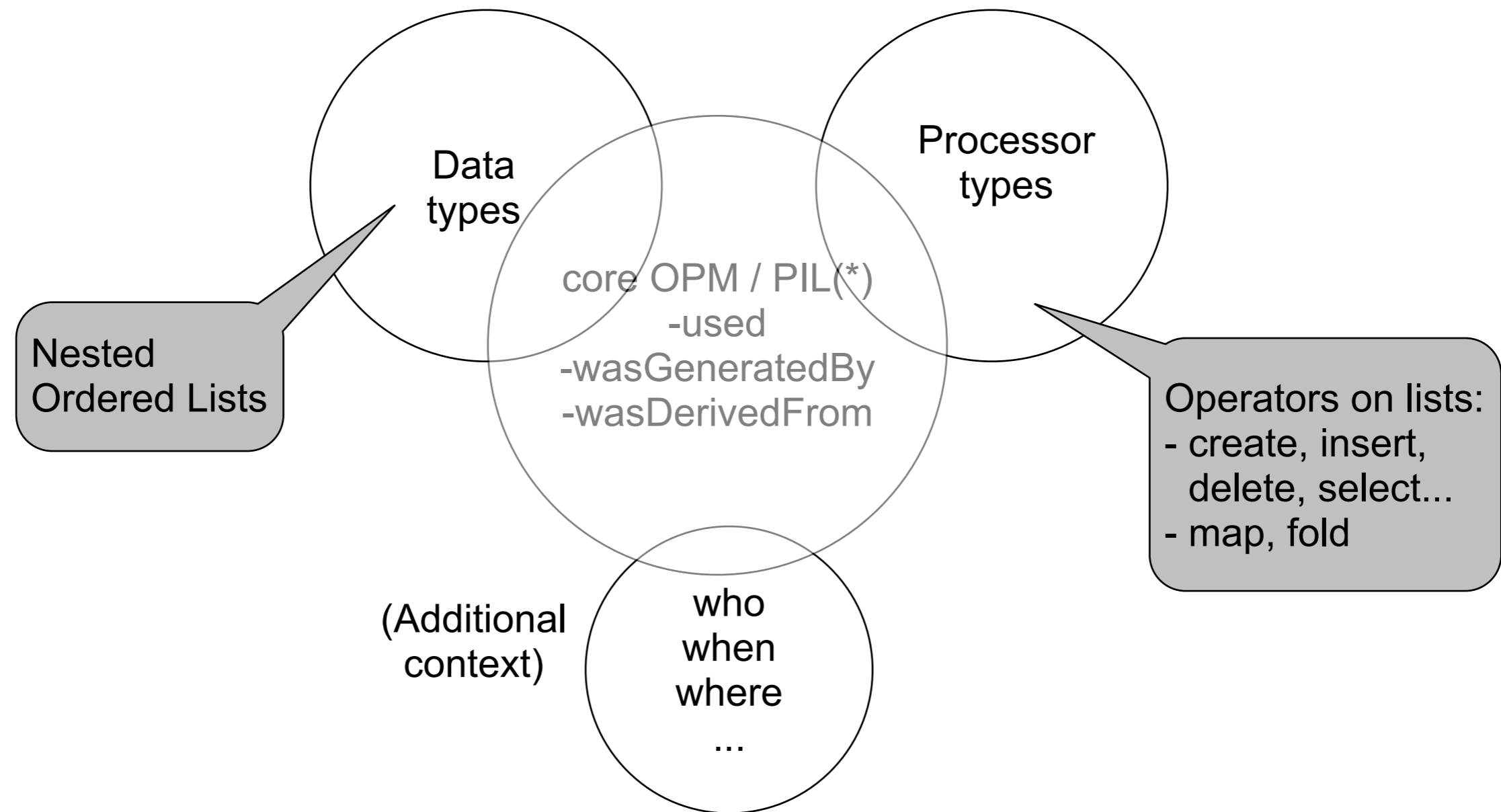
$h(x_1 \rightarrow a_1) = h(x_2 \rightarrow a_2) = X \rightarrow A,$

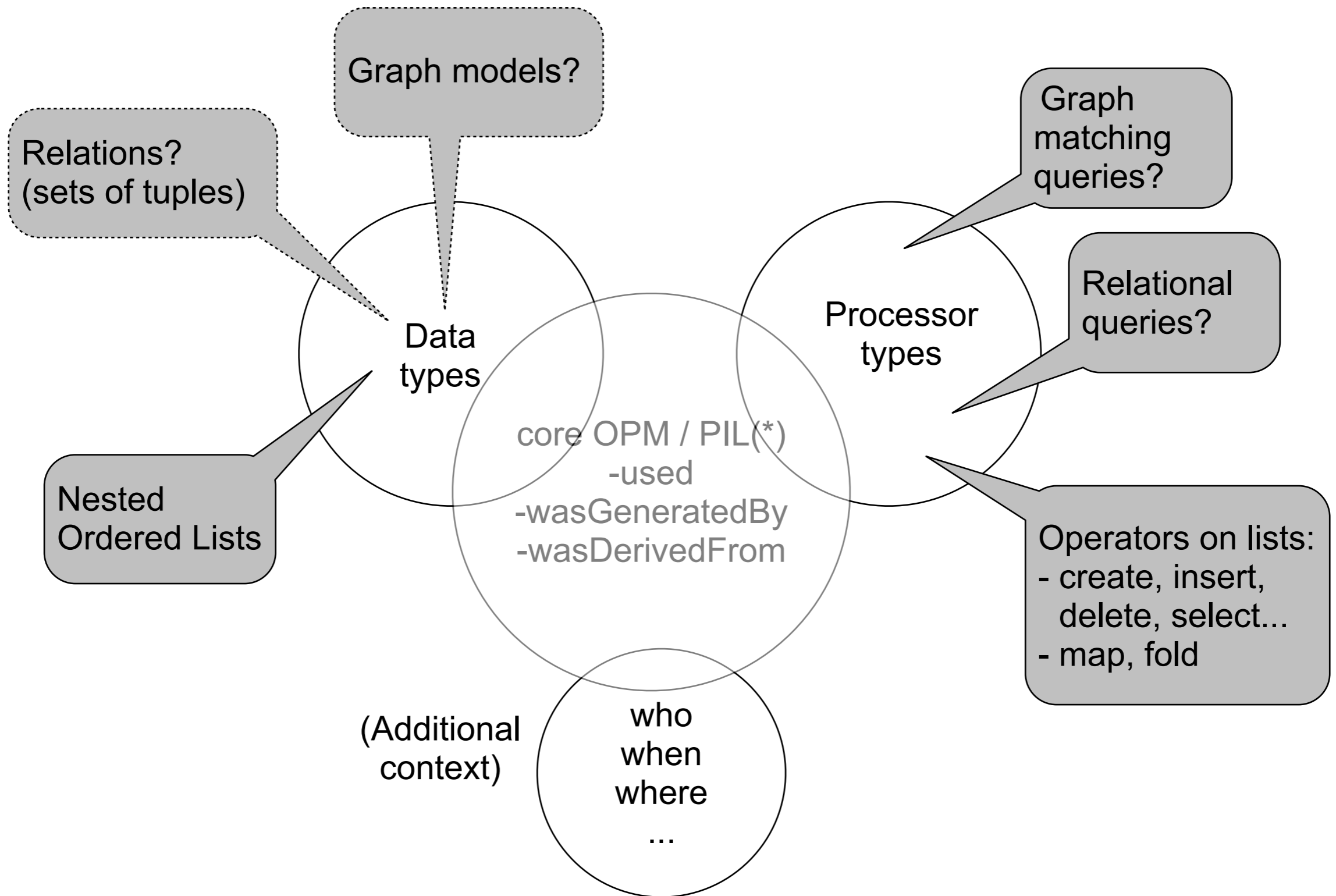
$h(a_1 \rightarrow y_1) = h(a_2 \rightarrow y_2) = A \rightarrow Y$

...





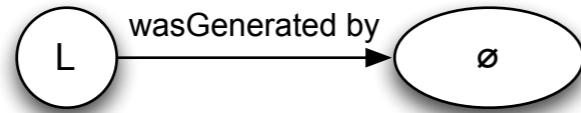




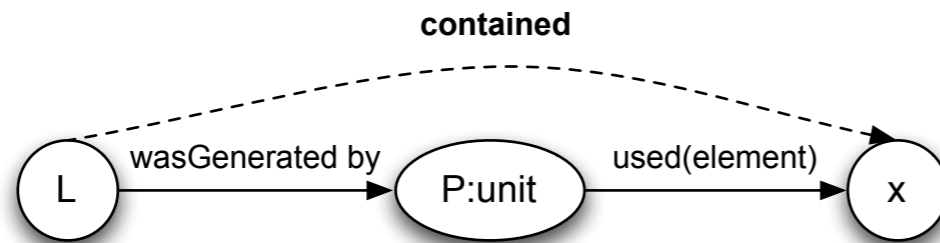
| Role | used in relation: | Context of use |
|-----------------------------|-------------------|---|
| element | <i>Contained</i> | $L \text{ Contained}(\text{element}) x$ |
| list | <i>Used</i> | $P \text{ Used}(\text{list}) L$ |
| position | <i>Used</i> | $P \text{ Used}(\text{position}) p$ |
| term generator filter | <i>Used</i> | List comprehensions, see Sec. 4.2 |
| function operand | <i>Used</i> | map, see Sec. 4.3 |

| Causal relation | Example |
|--|---|
| $\text{Contained}(R) \subseteq [\tau] \times A \times [\text{Int}]$ | $L' \text{ Contained}([i_1 \dots i_n]) x$ x was inserted into L at position $[i_1 \dots i_n]$ |
| $\text{wasSelectedFrom}(R) \subseteq [\tau] \times [\tau] \times [\text{Int}]$ | $L' \text{ wasSelectedFrom}([i_1 \dots i_n]) L$ L' at position $[i_1 \dots i_n]$ was selected from L |
| $\text{wasRemovedFrom}(R) \subseteq [\tau] \times [\tau] \times [\text{Int}]$ | $L' \text{ wasRemovedFrom}([i_1 \dots i_n]) L$ L' at position $[i_1 \dots i_n]$ was deleted from L |
| $\text{wasSameAs} \subseteq A \times A$ | $L \text{ wasSameAs } x$ inferred (various contexts) |

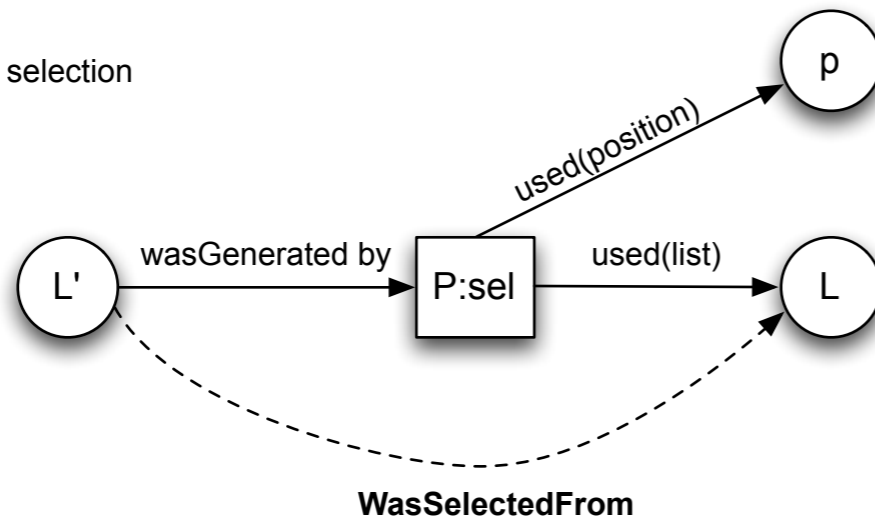
empty list



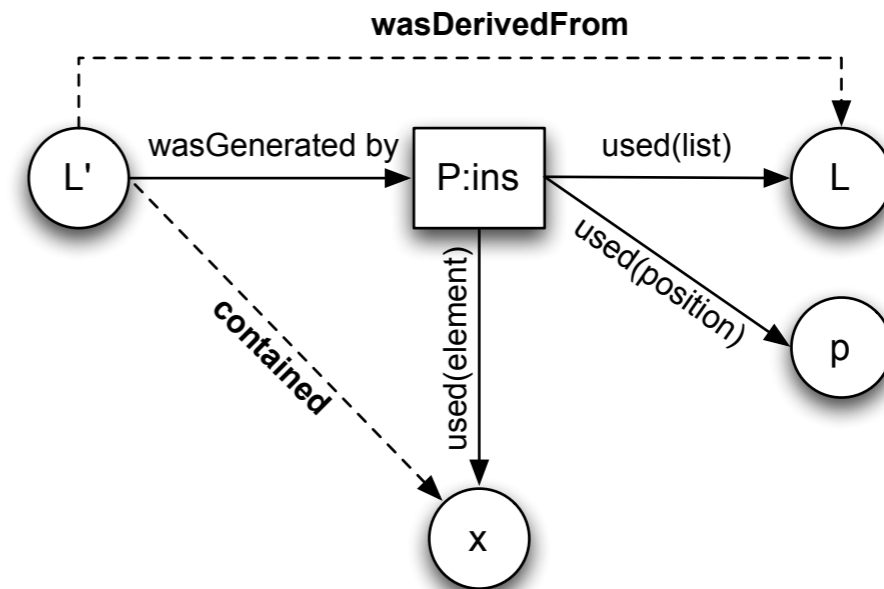
unit



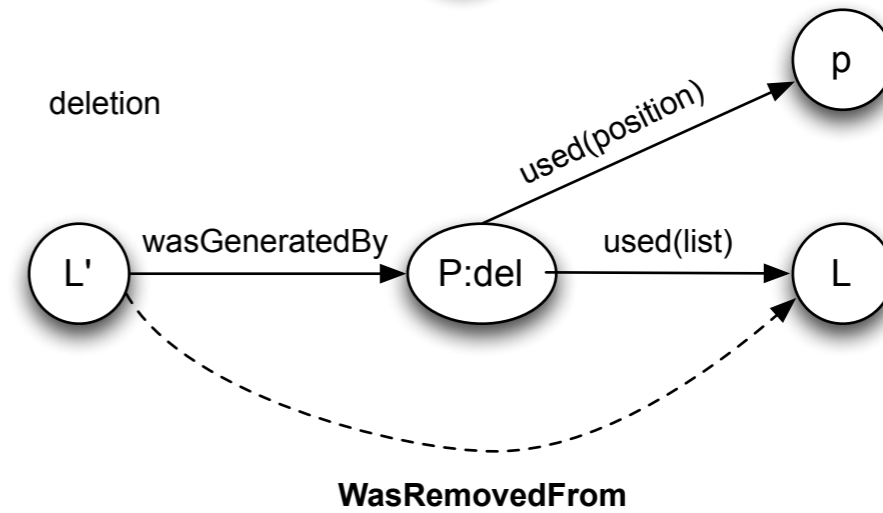
selection



insertion



deletion



- Equalities on operators may translate into inferences on the graphs

$$\text{sel}(\text{ins } x \ L \ p) \ p = x$$

sameAs(L1,L) :-

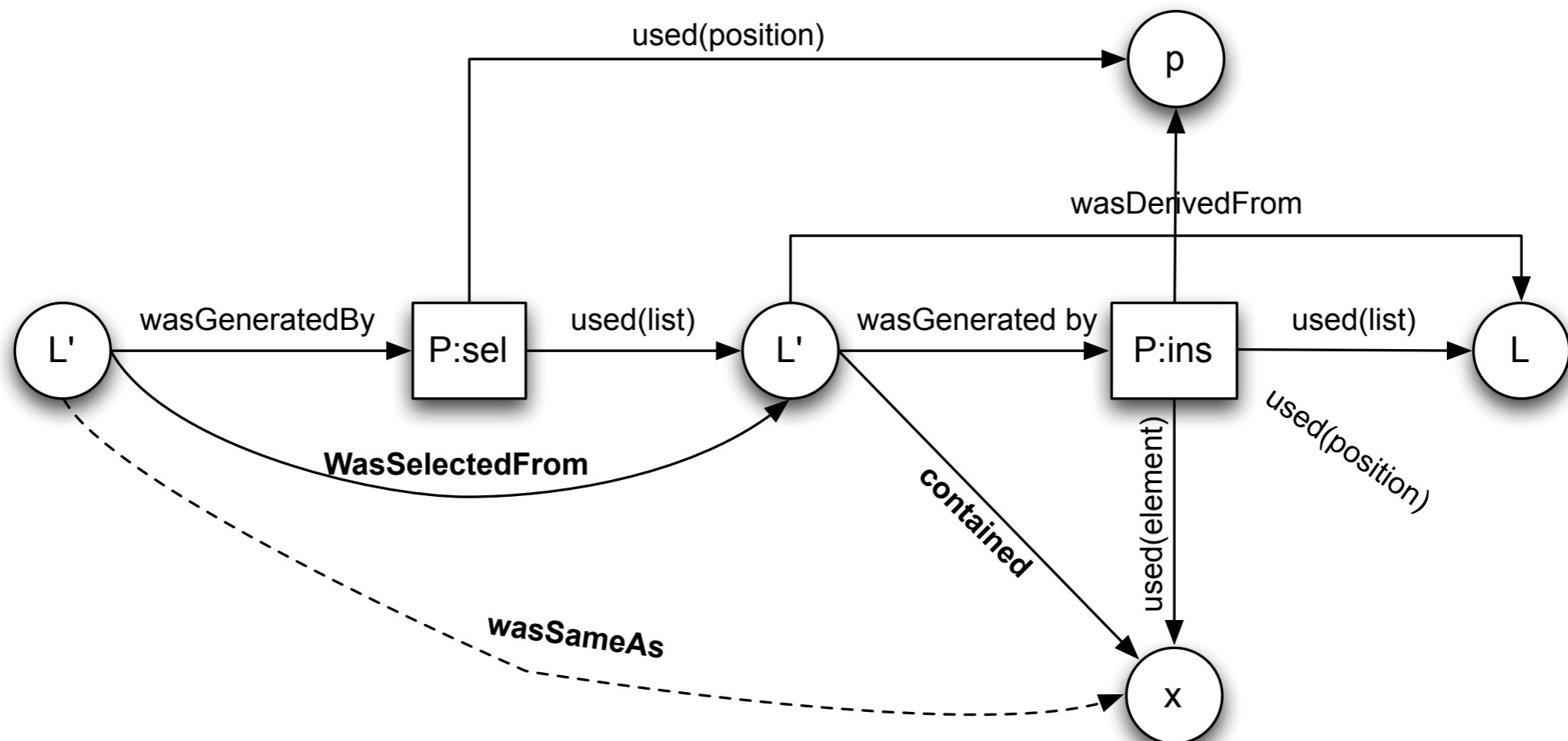
*pType(P1, ins), used(P1, X, element), used(P1, Pos, position), wgby(L,P1),
pType(P2, sel), used(P2, L, list), used(P2, Pos, position), wgby(L1,P2).*

- Equalities on operators may translate into inferences on the graphs

$$sel(ins\ x\ L\ p)\ p = x$$

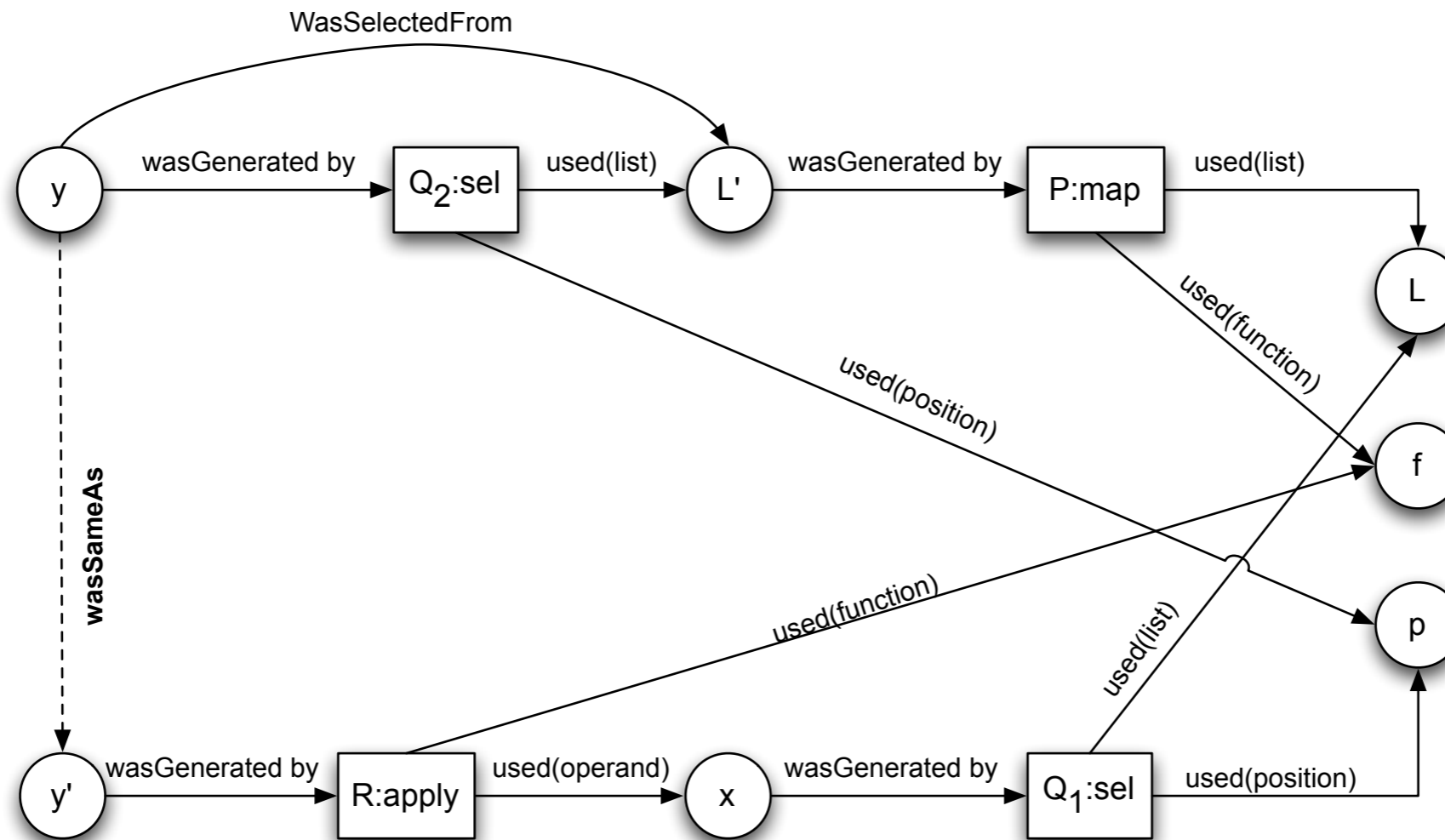
sameAs(L1,L) :-

*pType(P1, ins), used(P1, X, element), used(P1, Pos, position), wgby(L,P1),
 pType(P2, sel), used(P2, L, list), used(P2, Pos, position), wgby(L1,P2).*



Approach applies to map, fold (reduce)...

$$sel (map f x) p = f (sel x p)$$



- OPM (PIL) a candidate starting point for workflow-based provenance
- extension mechanisms are provided, but they must be used sensibly
 - data types
 - processor types
- Provenance of nested ordered lists used as a prototypical example
 - semantics of provenance graphs and graph composition grounded in the semantics of lists
- Can this approach be useful for other interesting data types?
 - sets of tuples / relational algebra