

Dynamic Provenance for SPARQL Updates using Named Graphs

James Cheney and Harry Halpin

W3C/MIT (@harryhalpin harry@w3.org)

October 23, 2014

Presentation Outline

Dynamic
Provenance
for SPARQL
Updates using
Named
Graphs

James Cheney
and Harry
Halpin

What's the
point?

1 What's the point?

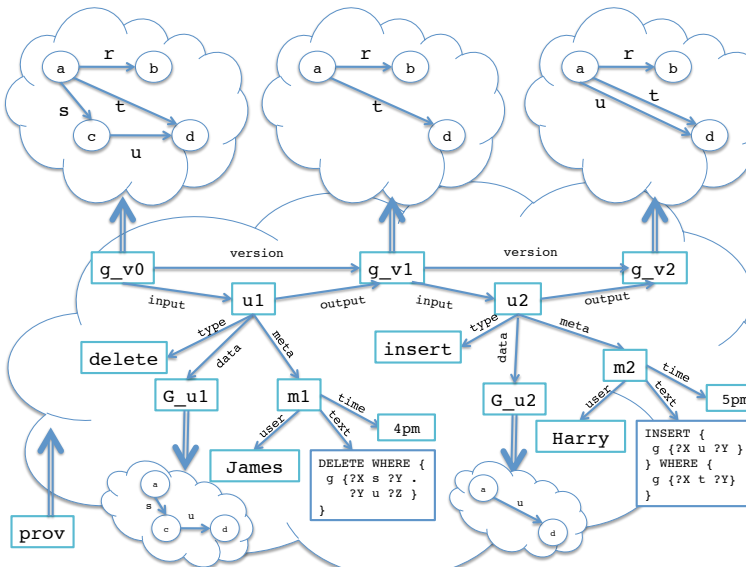
- While the (Semantic) Web currently does have a way to exhibit static provenance information in the W3C PROV standards, the Web does not have a way to describe dynamic changes to data.
- Our hypothesis is that a simple vocabulary, composed of insert, delete, and copy operations as introduced by Buneman et al., along with explicit identifiers for update steps, versioning relationships, and metadata about updates provides a flexible format for dynamic provenance on the Semantic Web.
- A primary advantage of our methodology is it keeps the changes to raw data separate from the changes in provenance metadata, so legacy applications will continue to work and the cost of storing and providing access to provenance can be isolated from that of the raw data.

An Example

Dynamic Provenance for SPARQL Updates Using Named Graphs

James Cheney and Harry Halpin

What's the point?



Terminology

Dynamic
Provenance
for SPARQL
Updates using
Named
Graphs

James Cheney
and Harry
Halpin

What's the
point?

- Any graph that records the insertion and deletion of triples from a given graph is considered a *provenance graph* (*prov*) for the given graph. which keeps track of auxiliary named graphs.
- C denotes basic graph (or dataset) patterns that may contain variables; R denotes conditions; P denotes patterns, and Q denotes queries. A graph store $\mathcal{D} = (G, \{g_i \mapsto G_1 \dots, g_n \mapsto G_n\})$ consists of a default graph G_0 together with a mapping from names g_i to graphs G_i .

SPARQL Update Model

Dynamic
Provenance
for SPARQL
Updates using
Named
Graphs

James Cheney
and Harry
Halpin

What's the
point?

For queries, we consider a simple form of provenance which calculates a set of named graphs “consulted” by the query.

$$\mathcal{S}[[C]]_G^D = \bigcup \{\text{names}(\mu(C)) \mid \mu \in [[C]]_G^D\}$$

$$\mathcal{S}[[P_1 \ . \ P_2]]_G^D = \mathcal{S}[[P_1]]_G^D \cup \mathcal{S}[[P_2]]_G^D$$

$$\mathcal{S}[[P_1 \ \text{UNION} \ P_2]]_G^D = \mathcal{S}[[P_1]]_G^D \cup \mathcal{S}[[P_2]]_G^D$$

$$\mathcal{S}[[P_1 \ \text{OPT} \ P_2]]_G^D = \mathcal{S}[[P_1]]_G^D \cup \mathcal{S}[[P_2]]_G^D$$

$$\mathcal{S}[[P \ \text{FILTER} \ R]]_G^D = \mathcal{S}[[P]]_G^D$$

$$\mathcal{S}[[\text{SELECT} \ ?\vec{X} \ \text{WHERE} \ P]]_G^D = \mathcal{S}[[P]]_G^D$$

$$\mathcal{S}[[\text{CONSTRUCT} \ C \ \text{WHERE} \ P]]_G^D = \mathcal{S}[[P]]_G^D$$

Create Graph

Dynamic
Provenance
for SPARQL
Updates using
Named
Graphs

James Cheney
and Harry
Halpin

What's the
point?

A graph creation `CREATE g` is translated to

```
CREATE  $g$ ;  
CREATE  $g_{-v_0}$ ;  
INSERT DATA {GRAPH  $prov$  {  
   $\langle g$  version  $g_{-v_0} \rangle$ ,  $\langle g$  current  $g_{-v_0} \rangle$ ,  
   $\langle u_1$  type create  $\rangle$ ,  $\langle u_1$  output  $g_{-v_0} \rangle$ ,  
   $\langle u_1$  meta  $m_i \rangle$ , (metadata)  
}}
```

where g_{-v_i} is the current version of g . Note that since this operation deletes g , after this step the URI g no longer names a graph in the store; it is possible to create a new graph named g , which will result in a new sequence of versions being created for it.

Clear Graph

A clear graph operation `CLEAR g` is handled as follows:

```
CLEAR g;  
DELETE WHERE {GRAPH prov {⟨g current g-v_i⟩}};  
INSERT DATA {GRAPH prov {  
  ⟨g version g-v_{i+1}⟩, ⟨g current g-v_{i+1}⟩,  
  ⟨u_i type clear⟩, ⟨u_i input g-v_i⟩,  
  ⟨u_i output g-v_{i+1}⟩, ⟨u_i meta m_i⟩,  
  (metadata)  
}}
```


Clear Graph

A clear graph operation `CLEAR g` is handled as follows:

```
CLEAR g;  
DELETE WHERE {GRAPH prov {⟨g current g-v_i⟩}};  
INSERT DATA {GRAPH prov {  
  ⟨g version g-v_{i+1}⟩, ⟨g current g-v_{i+1}⟩,  
  ⟨u_i type clear⟩, ⟨u_i input g-v_i⟩,  
  ⟨u_i output g-v_{i+1}⟩, ⟨u_i meta m_i⟩,  
  (metadata)  
}}
```

Load Graph

A load graph operation `LOAD h INTO g` is handled as follows:

```
LOAD  $h$  INTO  $g$ ;  
DELETE WHERE {GRAPH  $prov$  { $\langle g$  current  $g_{-v_i} \rangle$ }};  
INSERT DATA {GRAPH  $prov$  {  
   $\langle g$  version  $g_{-v_{i+1}} \rangle$ ,  $\langle g$  current  $g_{-v_{i+1}} \rangle$ ,  
   $\langle u_i$  type load  $\rangle$ ,  $\langle u_i$  input  $g_{-v_i} \rangle$ ,  
   $\langle u_i$  output  $g_{-v_{i+1}} \rangle$ ,  $\langle u_i$  source  $h_j \rangle$ ,  
   $\langle u_i$  meta  $m_i \rangle$ , (metadata)  
}}
```

where h_j is the current version of h .

Insert Graph

An insertion `INSERT {GRAPH g { C }} WHERE P` is translated to a sequence of updates that creates a new version and links it to URIs representing the update, as well as links to the source graphs identified by the query provenance semantics and a named graph containing the inserted triples:

```
CREATE  $g_{-u_i}$ ;  
INSERT {GRAPH  $g_{-u_i}$  { $C$ }} WHERE  $P$ ;  
INSERT {GRAPH  $g$  { $C$ }} WHERE  $P$ ;  
CREATE  $g_{-v_{i+1}}$ ;  
LOAD  $g$  INTO  $g_{-v_{i+1}}$ ;  
DELETE DATA {GRAPH  $prov$  { $\langle g$  current  $g_{-v_i} \rangle$ }};  
INSERT DATA {GRAPH  $prov$  {  
   $\langle g$  version  $g_{-v_{i+1}} \rangle$ ,  $\langle g$  current  $g_{-v_{i+1}} \rangle$ ,  
   $\langle u_i$  input  $g_{-v_i} \rangle$ ,  $\langle u_i$  output  $g_{-v_{i+1}} \rangle$ ,  
   $\langle u_i$  type insert  $\rangle$ ,  $\langle u_i$  data  $g_{-u_i} \rangle$   
   $\langle u_i$  source  $s_1 \rangle$ , ...,  $\langle u_i$  source  $s_m \rangle$ ,
```

Delete Graph

A deletion `DELETE {GRAPH g { C }} WHERE P` is handled similarly to an insert, except for the update type annotation.

```
CREATE  $g_{-u_i}$ ;  
INSERT {GRAPH  $g_{-u_i}$  { $C$ }} WHERE  $P$ ;  
DELETE {GRAPH  $g$  { $C$ }} WHERE  $P$ ;  
CREATE  $g_{-v_{i+1}}$ ;  
LOAD  $g$  INTO  $g_{-v_{i+1}}$ ;  
DELETE DATA {GRAPH  $prov$  { $\langle g$  current  $g_{-v_i} \rangle$ }};  
INSERT DATA {GRAPH  $prov$  {  
   $\langle g$  version  $g_{-v_{i+1}} \rangle$ ,  $\langle g$  current  $g_{-v_{i+1}} \rangle$ ,  
   $\langle u_i$  input  $g_{-v_i} \rangle$ ,  $\langle u_i$  output  $g_{-v_{i+1}} \rangle$ ,  
   $\langle u_i$  type delete  $\rangle$ ,  $\langle u_i$  data  $g_{-u_i} \rangle$   
   $\langle u_i$  source  $s_1 \rangle$ , ...,  $\langle u_i$  source  $s_m \rangle$ ,  
   $\langle u_i$  meta  $m_i \rangle$ , (metadata)}}}
```

Dynamic Provenance Vocabulary

Dynamic Provenance for SPARQL Updates using Named Graphs

James Cheney and Harry Halpin

What's the point?

Name	Description
upd:insert	Triples inserted into a graph.
upd:delete	Triples deleted from a graph.
upd:load	Triples copied from another graph.
upd:clear	All triples in graph deleted
upd:create	New graph initialized.
upd:drop	A graph is deleted.
upd:input	Linked to graph used as input in update operation.
upd:output	Link to graph of output of the operation.
upd:data	Changed data in insert/delete operation.
upd:version	Link between two versions of a graph.
upd:type	Type of update operation.
upd:current	Link to most current state of graph.
upd:source	Any other graph that was consulted by the update.
upd:meta	Link to any metadata about the graph.

Lightweight Dynamic Provenance Vocabulary

Conclusion

Dynamic
Provenance
for SPARQL
Updates using
Named
Graphs

James Cheney
and Harry
Halpin

What's the
point?

Provenance is a challenging problem for RDF. While some progress has been made on provenance and annotation for RDFS inferences and SPARQL queries, so far there has not been work on provenance for SPARQL Update. We have outlined an approach to the problem drawing on similar work in database archiving and copy-paste provenance in relational databases. In particular, the metadata carried by our technique can use the PROV data model already developed by the W3C Provenance Interchange Working Group