

















Β.	Ludaescher.	ECS289F-W05.	Topics in Scientific Data Management	

(1)	0	_	
(1)	OBSERVATION		V PROPERTY. OBSPROPERTY = 1 PROPERTY. OBSPROPERTY
			THEM.OBSERVABLEFTEM = 1 HEM.OBSERVABLEFTEM Y SDITLU CONTEXT □ ∃ SDITLU CONTEXT □
			V SPARALCONTEXT. SPARALCONTEXT IT I SPARALCONTEXT. SPARALCONTEXT IT
(2)	OCCURRENCE	_	V TEMPORALCONTEXT. TEMPORALCONTEXT + D TEMPORALCONTEXT. TEMPORALCONTEXT
(2)	ABUNDANCE	-	
(3)	BIOMASS		ADINDANCE
(5)	BIOMASSGSM	=	BIOMASS TO GRAMSSOLIADE METER
6	BIOMASSPPH	=	BIOMASS TO POUNDSPERHECTARE
(7)	RICHNESS		ABUNDANCE
(8)	PRODUCTIVITY		BIOMASS
(9)	PRODUCTIVITYGSM		BIOMASSGSM IT GRAMSSQUAREMETER
(10)	PRODUCTIVITYPPH		BIOMASSPPH POUNDSPERHECTARE
(11)	Species		ObservableItem
(12)	Plot		SpatialContext
(13)	QUADRAT	\subseteq	SpatialContext
(14)	Year	\Box	TemporalContext
(15)	Season		TemporalContext
(16)	WINTER	\subseteq	Season
(17)	Spring		Season
(18)	Summer		Season
(19)	UNITTYPE		VForUnit.Unit ⊓ ≤1 ForUnit
(20)	GRAMSSQUAREMETER	≡	UNITTYPE T FORUNIT. (GSMUNIT)
(21)	POUNDSPERHECTARE	≡	UNIT TYPE T FORUNIT. {PPHUNIT}
(22)	SIUNIT		Unit 🗆





We consider query annotations q and semantic annotations α of the form:

• $q = \forall \bar{u} \exists \bar{v} P_{S'}(\bar{u}) := \varphi_{S}(\bar{u}, \bar{v})$

• $\alpha = \forall \bar{x} \exists \bar{y} \alpha_{\mathsf{S}}(\bar{x}) \rightarrow \alpha_{O}(\bar{x}, \bar{y})$

Semantic Annotations

Here, $P_{S'}$ is a logic atom over the output schema S' and φ_S is a query over the input schema(s) S. Similarly, in a semantic annotation α we relate instances of a schema S with those of an ontology O via subformulas α_S and α_O . The basic idea of computing $\alpha' = q(\alpha)$ is as follows: We would like to relate instances of the output schema S'

with instances of *O*. Assume we can find instances of φ_S (in *q*) that imply certain instances of α_S (in *α*) to be true. For these we then have established the desired relation between S' and *O*; we can denote this abstractly as S' $\stackrel{q}{\sim}$ S $\stackrel{a}{\sim}$ *O*. More precisely, consider *q* as a constraint of the form

$$q(\bar{u}) = P_{S'} \to \underbrace{Q_1 \wedge \dots \wedge Q_n \wedge \psi}_{\varphi_S}$$

and α of the form

$$\alpha(\bar{x}) = \underbrace{A_1 \wedge \cdots \wedge A_k}_{\alpha_k} \to \alpha_{\mathcal{C}}$$

Here, $P_{S'}$ is a logic atom over the output schema S'. Similarly, the Q_i and A_j are logic atoms over the input schema(s) S; ψ and α_0 are first-order formulas. Note that we assume that all \exists -quantified variables (\bar{v} and \bar{y} above) have been eliminated through Skolemization⁶, so that q and α can be seen as (implicitely) \forall -quantified formulas with variables \bar{u} and \bar{x} , respectively. In particular, we can assume that the variables \bar{u} in $q(\bar{u})$ are disjoint from the variables \bar{x} in $\alpha(\bar{x})$.

Observe that q can be written as a conjunction $q_1 \wedge \cdots \wedge q_n \wedge q_{\psi}$ with $q_i = P_{S^*} \rightarrow Q_i$, and $q_{\psi} = P_{S^*} \rightarrow \psi$. Now assume that there is a substitution σ that unifies some atom Q_{i_0} and some A_{j_0} , i.e., $Q_{i_0}^{\sigma} = A_{j_0}^{\sigma, 7}$. Then we can infer from $q_{i_0} = P_{S^*} \rightarrow Q_{i_0}$ and α a new annotation constraint α'_{i_0} of the form:

$$\alpha'_{i_0} = P^{\sigma}_{S'} \wedge (\alpha_S \setminus A_{i_0})^{\sigma} \to \alpha^{\sigma}_{O}$$

where $(\alpha_{S} \setminus A_{j_0})$ is the conjunction $A_1 \wedge \cdots \wedge A_k$ with A_{j_0} removed. It is easy to show that the annotation α'_{i_0} is implied by q_i and α . In this way, by successively "resolving away" atoms A_j from α_{S} with matching atoms Q_i from φ_{S} , we can obtain new semantic annotations α' that relate elements of the output schema S' to those in the ontology O.

Example (Biodiversity Workflow)

Example 6 Consider the following query annotation expressed as a skolemized first-forder formula for the *Seasonal Community* actor in Figure 2

(1) $a:biom[yr=r, seas=s, plt=t, qd=q, spp=p, bm=b] \rightarrow a:biom[yr=r, seas=s, plt=t, qd=q, spp=p, bm=b] \land f(a):sscd[spp=g(a)] \land g(a)=p$

and semantic annotation (9) in Figure 4, expressed as the first-order formula

(2) $x:biom[bm=y] \rightarrow x[PROPERTY=y:BIOMASS]$

We can resolve (1) with (2) using the substitution $\sigma = \{x \mapsto a, y \mapsto b\}^8$, which results in the formula

(3) $a:\text{biom1[yr=r, seas=s, plt=t, qd=q, spp=p, bm=b]} \rightarrow a[\text{item=b:Biomass}] \land f(a):\text{sscd[spp=g(a)]} \land g(a)=p$

Observe that we now have bm values for the output schema biom1 semantically annotated as BIOMASS instances, linked through the PROPERTY role to corresponding tuples. We also get some additional information, which can be viewed as constraints over the input schema. These additional constraints can be ignored, bringing (3) into our standard form for semantic annotations, which results in

(4) $a:\text{biom1[yr=}r, \text{seas}=s, \text{plt}=t, \text{qd}=q, \text{spp}=p, \text{bm}=b] \rightarrow a[\text{property}=b:\text{Biomass}]$

We can further simplify (4) by dropping the attribute variables not used elsewhere in the annotation, giving

(5) $a:biom1[bm=b] \rightarrow a[property=p:Biomass]$

B. Ludaescher, ECS289F-W05, Topics in Scientific Data Management

Family Example	
Conjunctive Queries. We start with a semantic annotation with the signature $\alpha: S \rightarrow$	0
• $p(c, p) \rightarrow CHILD(c) \land PARENT(p)$	(α)
and a query annotation with the signature $q:\ S'\coloneqq S$	
• $gc(y, x) := p(x, z), p(z, y)$	(q)
To propatate the semantic type α through $q,$ we need the "backward" (i.e., left-to-right) di	irection of q^4 :
• $gc(y, x) \rightarrow \exists z.p(x, z) \land p(z, y)$	(q_b)
Before we chase q_b with α , we eliminate the \exists -quantifier via Skolemization:	
• $gc(y, x) \rightarrow p(x, f_z(x, y)) \wedge p(f_z(x, y), y)$	(q_b')
We can resolve (α) and (q'_b) via the substitution $\sigma_1 = \{c \rightarrow x, \ p \rightarrow f_z(x, y)\}$ resulting in:	
• $gc(y,x) \to CHILD(x) \land PARENT(f_z(x,y)) \land p(f_z(x,y),y)$	$(res_{\sigma_1}(\alpha,q_b'))$
We can apply α again, now with $\sigma_2 = \{c \rightarrow f_z(x, y), p \rightarrow y\}$ resulting in:	
• $gc(y, x) \to Child(x) \land Parent(f_z(x, y)) \land Child(f_z(x, y)) \land Parent(y)$ $(res_{\sigma_2}) \land res_{\sigma_2} \land$	$(\alpha, res_{\sigma_1}(\alpha, q_b')))$
Observe how this result is in the desired form $\alpha': S' \to O$. We now know that x and y are in and PARENT, respectively. Indeed x is the grandchild, y is the grandparent. We also get information (that we might choose to ignore), i.e., that there is some instance $f_z(x, y)$ who and a PARENT (as it turns out, this is the intermediate parent, establishing the link between and grandparent). B. Ludaescher, EGS289F-W05, Topics in Scientific Data Management	stances of CHILD some additional is both a CHILD en the grandchild