

Default-all is dangerous!

Wolfgang Gatterbauer
Alexandra Meliou
Dan Suciu

3rd USENIX Workshop on the Theory and Praxis of Provenance (Tapp'11)

Overview Provenance Definitions

	Why?	Where?
<i>Naive</i>	Witness	"SQL interpretation"
<i>Provenance definition</i>	<p>Why-provenance = <u>w</u>itness basis (α_w)</p> <p>Buneman et al. [ICDT'01]</p>	<p>Where-provenance = <u>p</u>ropagation (α_p)</p> <p>Buneman et al. [PODS'02]</p>
<i>QRI definition (Query-Rewrite-Insensitive)</i>	<p><u>M</u>inimal <u>w</u>itness basis (α_w^m)</p> <p>Buneman et al. [ICDT'01]</p>	<p><u>D</u>efault-all <u>p</u>ropagation (α_p^d)</p> <p>Bhagwat et al. [VLDB'04]</p>

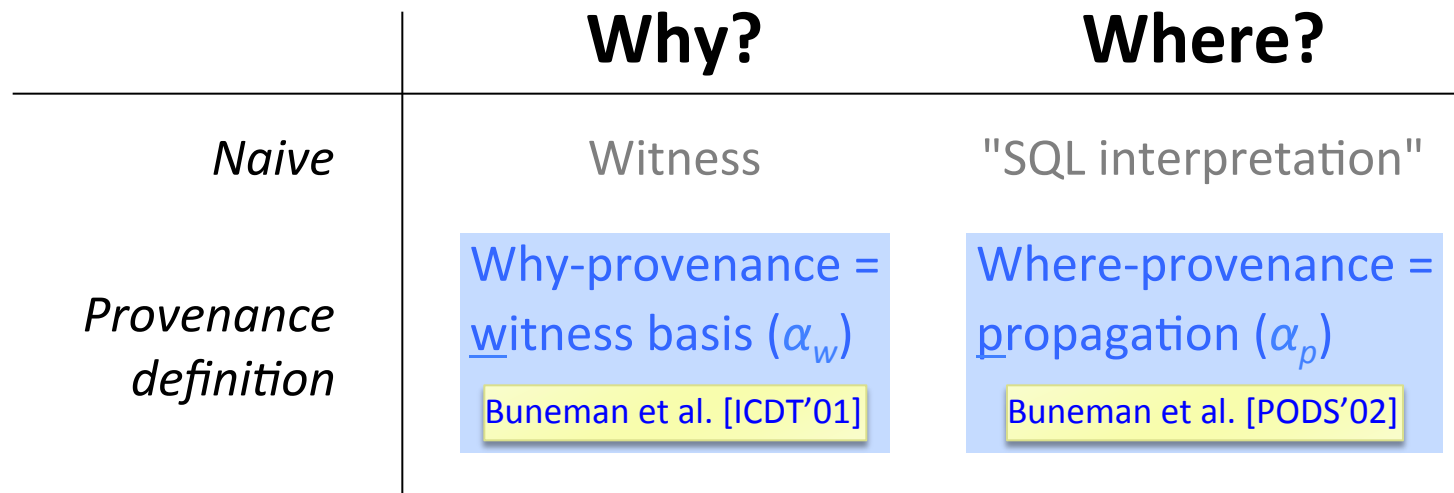
Has problems if one interprets annotations on attribute values

Minimal propagation (α_p^m)
Proposed in this paper!

Independent work presented at this WS

We do not discuss here whether QRI is desirable (see also Glavic, Miller [Tapp'11]), but merely point out that, if aiming for QRI, care has to be taken about the ramifications of the proposed semantics.

Overview Provenance Definitions



Glavic, Miller [Tapp'11]

Semantics		Sound	Complete	Responsible	Insensitive (set)	Insensitive (bag)	Stable
Why	Wit	-	X	-	X	X	X
	Why	-	X	-	-	X	X
	IWhy	-	X	X	X	X	X
Where	Where	-	-	-	-	?	X
	IWhere	-	-	-	X	X	-
How		-	X	-	-	X	X
Lineage-based	Lineage	X	X	-	-	-	X
	PI-CS	X	X	-	-	-	X
	C-CS	X	-	-	-	-	X
Causality		-	X	X	X	X	X

(α_w^m)
[ICDT'01]

Default-all propagation (α_p^d)
Bhagwat et al. [VLDB'04]

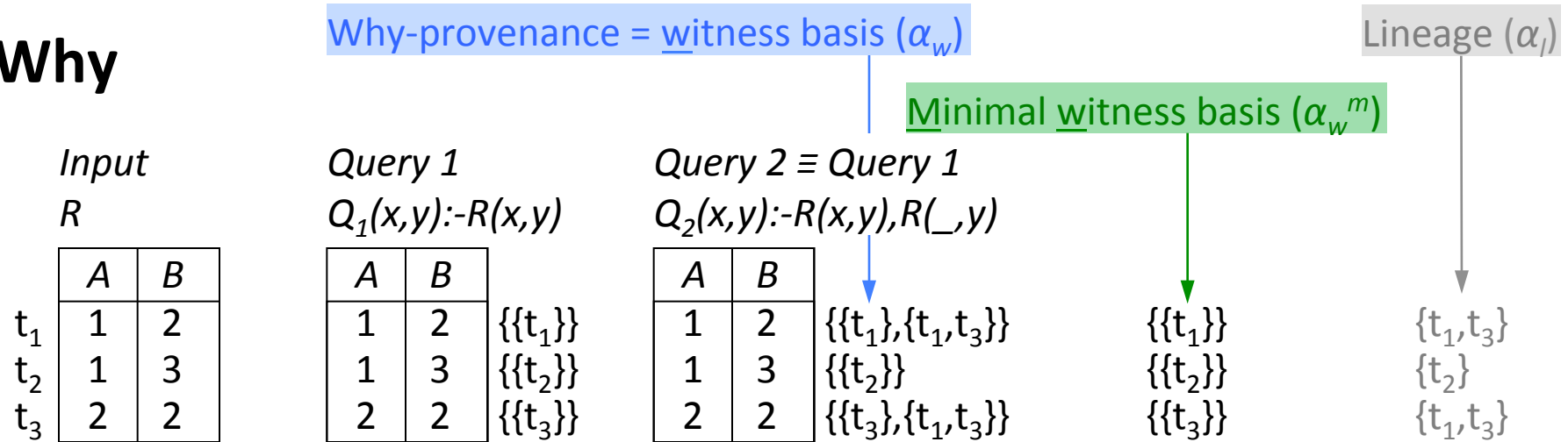
Minimal propagation (α_p^m)
Proposed in this paper!

Has problems if one interprets annotations on attribute values

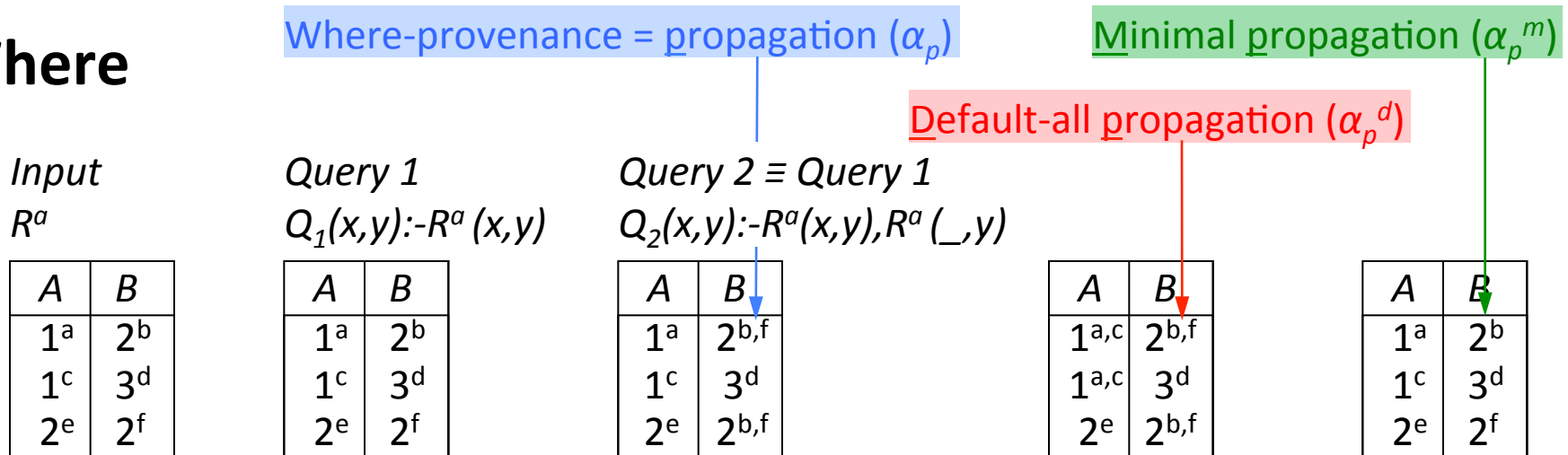
Note that Minimal propagation is "stable", in contrast to Default-all

Example 1: Query-Rewrite-Insensitivity (QRI)

Why



Where



Real example: Why Default-all is dangerous

Hanako queries a community DB for contents of LF-milk*:

Community Database

R^a

Food	Content
LF Milk	Cesium-137 ^b
LF Milk	Calcium ^d
SC Water	Cesium-137 ^f

^b Bob, March 18, 2011
Don't drink, lots of Cesium!

^f Fuyumi, March 19, 2011
No Cesium, save to drink!

Hanako's query

$Q(y) :- R^a('LF Milk', y)$

Content
Cesium-137 ^b ???
Calcium ^d

Default-all propagation makes her drink the milk:

Default-all propagation (α_p^d)

Content
Cesium-137 ^b ^f
Calcium ^d

^b Bob, March 18, 2011
Don't drink, lots of Cesium!

^f Fuyumi, March 19, 2011
No Cesium, save to drink!

"semantically irrelevant information": annotations leak over from SC Water tuple to LF Milk

Minimal propagation (α_p^m)

Content
Cesium-137 ^b
Calcium ^d

^b Bob, March 18, 2011
Don't drink, lots of Cesium!

"all relevant and only relevant"

* Note the one-to-one correspondence of this example with example 1

Definition Minimal propagation (α_p^m)

$$\alpha_p^m(t, A, Q) := \bigcup_{\substack{t' \in \cup \alpha_w^m(t, Q) \\ A' \in \text{attributes of } t' \text{ propagating to cell}(t, A)}} \alpha_p(t', A')$$

\cup transforms 'sets of sets' into 'sets', hence something like QRI lineage

Intuition:

Return the intersection between:

- query-specific where-provenance (α_p)
- and QRI minimal witness basis (α_w^m)

"all relevant ... and only relevant"

Example 1

Input

R^a

	A	B
t_1	1^a	2^b
t_2	1^c	3^d
t_3	2^e	2^f

Where provenance (α_p)

Query 2

$Q_2(x, y) :- R^a(x, y), R^a(_, y)$

A	B	
1^a	$2^{b,f}$	$\{\{t_1\}\}$
1^c	3^d	$\{\{t_2\}\}$
2^e	$2^{b,f}$	$\{\{t_3\}\}$

$\{t_1\}$
 $\{t_2\}$
 $\{t_3\}$

$\cup \alpha_w^m$

Minimal witness basis (α_w^m)

Minimal propagation (α_p^m)

	A	B
t_4	1^a	2^b
t_5	1^c	3^d
t_6	2^e	2^f

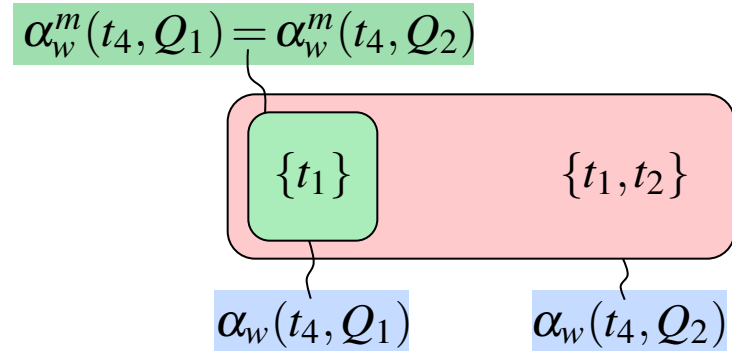
$$\alpha_p^m(t_4, B, Q_2) = \bigcup_{t' \in \{t_1\}, A'} \alpha_p(t', A') = \alpha_p(t_1, B) = \{b\}$$

Example 1: Illustration of "minimal" versus "all"

Why-provenance

Why-provenance (α_w)

Minimal witness basis (α_w^m)

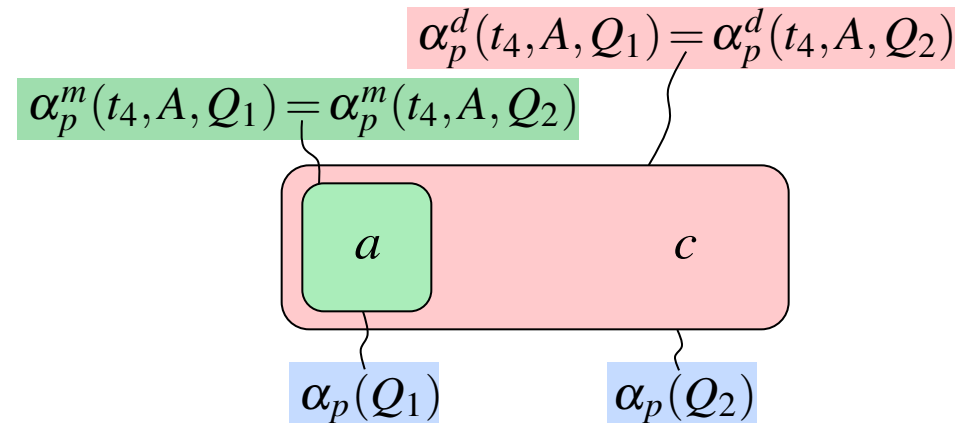


Where-provenance

Where-provenance (α_p)

Default-all propagation (α_p^d)

Minimal propagation (α_p^m)



Interpretation of Annotations 1: Attribute Value*



athens heraklion chania			
Item Name	Description	Population	Add columns
athens	PIRAEUS (Athens) - HERAKLION (Crete) - PIRAEUS (Athens) . PIRAEUS (Athens) - CHANIA (Crete) - PIRAEUS (Athens)	4 possible values	
heraklion	Heraklion or Iraklion is the largest city and capital of Crete. It is also the 4th largest city in Greece. Heraklion is the capital of	1 possible value	
kania	Chania confusingly is sometimes written Hania though it can also be written Khania, Cania, Canea and Kanis and in Greek is Χανιά	1 possible value	
Crete	A superb way of enjoying the journey to Crete is to fly to Athens and take the ferry from Piraeus (Direce) the port serving Athens	623,666	
Mykonos	Heraklion and Chania are international airports, Sitia airport is currently receiving domestic flights only (charter flights are expected to	9,320	
Istanbul	14 Days - Depart USA, stops include, Istanbul, Mount Athos, Skithos, Samos, Kusadasi, Delos,	8,260,000	

* Interpretation of annotations on entity attribute values favored by us and underlying our model

Interpretation of Annotations 1: Attribute Value*

Google squared labs

athens heraklion chania

Square it Add to this Square

Item Name	Description	Population
athens	PIRAEUS (Athens) - HERAKLION (Crete) - PIRAEUS (Athens) . PIRAEUS (Athens) - CHANIA (Crete) - PIRAEUS (Athens)	
heraklion	Heraklion or Iraklion is the largest city and capital of Crete. It is also the 4th largest city in Greece. Heraklion is the capital of	Possible values <input type="radio"/> 750000 Low confidence Greece. LOCATION. Official Website: http://www.cityofathens.gr/ . Population: 750000. Population of Athens metropolitan area, 3.7 million www.ndb.com - all 2 sources »
kania	Chania confusingly is sometimes written Hania though it can also be written Khania, Cania, Canea and Kanis and in Greek is Χανιά	<input checked="" type="radio"/> 22936, 24234 Low confidence Population for Athens www.freebase.com
Crete	A superb way of enjoying the journey to Crete is to fly to Athens and take the ferry from Piraeus (Greece) - the port serving Athens	<input type="radio"/> 1,102 Low confidence pop. for Athens www.citytowninfo.com
Mykonos	Heraklion and Chania are international airports, Sitia airport currently receiving domestic flights only (charter flights are expected)	<input type="radio"/> 18,967 Low confidence pop. for Athens www.citytowninfo.com - all 2 sources »
Istanbul	14 Days - Depart USA, stops include, Istanbul, Mount Athos, Skithos, Samos, Kusadasi, Delos,	Search for more values »

Annotations on values of an attribute (here "population") for a particular entity (here "Athens")

Argument: Interpreting cell annotations as relevant to the tuple (entity) adds something that is not trivially modeled with normalized tables.

* Interpretation of annotations on entity attribute values favored by us and underlying our model

Interpretation of Annotations 2: Domain Value*

Domain value annotations*

Input R^a :

A	B
1 ^a	2 ^b
1 ^c	3 ^d
2 ^e	2 ^f

b Bob, March 18, 2011
This number is a prime number.

f Fuyumi, March 19, 2011
Two is not a prime number because it is even.

Input S^a :

...	Date
...	Dec 25
...	...
...	Dec 25

b This is a holiday.

f This is a holiday too !!!

Argument for default-all: If annotations are on domain values, then retrieving all annotations are relevant.

Alternative representation

Annotation table S^a :

B	annotation
2	b: Bob, March 18, 2011 This number is a prime number.
2	f: Fuyumi, March 19, 2011 Two is not a prime number because it is even

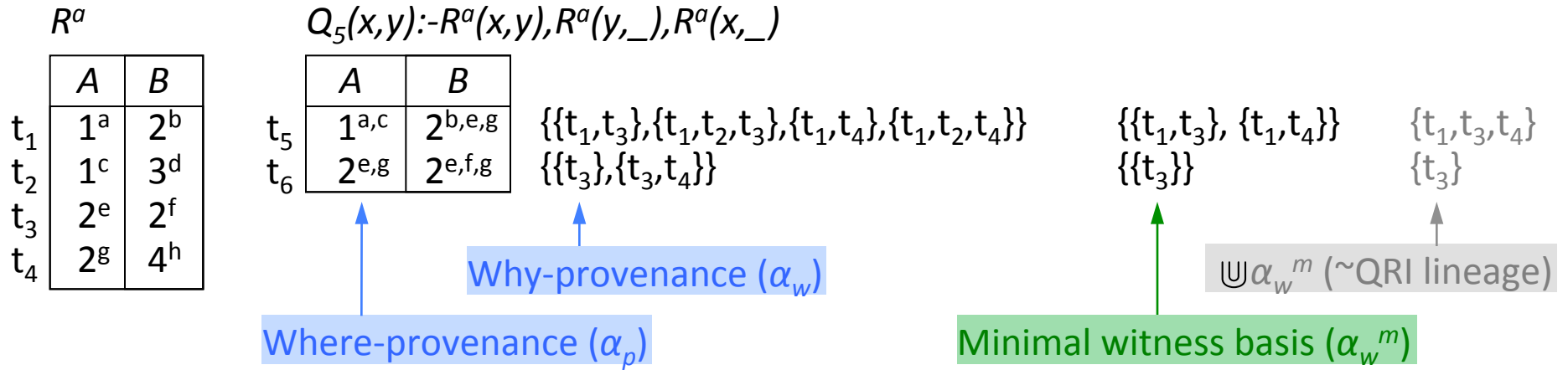
Annotation table S^a :

Date	annotation
Dec 25	This is a holiday.

Counter-Argument: But then these annotations can be modeled in a separate table as normalized tables.

* Alternative interpretation suggested by Wang-Chiew Tan (example created after conversation at Sigmod 2011)

Backup: Detailed Example 2



Default-all propagation (α_p^d)

A	B
$1^{a,c}$	$2^{b,e,f,g}$
$2^{e,g}$	$2^{b,e,f}$

$\alpha_p^d(t_4, B, Q_5) = \alpha_p(t_4, B, Q_6)$ with
 $Q_6(x,y):-R^a(x,y),R^a(y,_),R^a(x,_),S^a(_ ,y)$

Minimal propagation (α_p^m)

	A	B
t_4	1^a	$2^{b,e,g}$
t_5	2^e	$2^{e,f}$

$\alpha_p^m(t_4, A, Q_5) = \bigcup_{t' \in \{t_1, t_3, t_4\}, A'} \alpha_p(t', A')$
 $= \alpha_p(t_1, A) = \{a\}$

$\alpha_p^m(t_5, B, Q_5) = \bigcup_{t' \in \{t_3\}, A'} \alpha_p(t', A')$
 $= \alpha_p(t_3, B) \cup \alpha_p(t_3, A) = \{e, f\}$

Note minimal propagation is not equivalent to just evaluating the where-provenance for the query:
 $Q_7(x,y):-R^a(x,y),R^a(y,_)$. E.g. $\alpha_p(t_5, B, Q_7) = \{e, f, g\}$