

Properties	Semantics	Insensitivity	Stability	Conclusions
••	0	00	0	0
Insensitivi	ty			

Insensitivity to Query Rewrite

• Equivalent queries have the same provenance

•
$$Q \equiv Q' \Rightarrow \mathcal{P}(Q, I, t) = \mathcal{P}(Q', I, t)$$



Properties	Semantics	Insensitivity	Stability	Conclusions
●O	0	00	0	0
Insensitivit	tv			

Insensitivity to Query Rewrite

• Equivalent queries have the same provenance

•
$$Q \equiv Q' \Rightarrow \mathcal{P}(Q, I, t) = \mathcal{P}(Q', I, t)$$

- Caveat: Which queries are equivalent?
 - Set vs. Bag semantics
 - Query language / Operators





Holy Grails

Properties	Semantics	Insensitivity	Stability	Conclusions
0.	0	00	0	0
Stability				

Stability with Respect to Query Language Extension

• Extend query language with new operators \Rightarrow no change to provenance of queries that do not use new operators



Properties	Semantics	Insensitivity	Stability	Conclusions
00	•	00	0	0
Where				

Where [Buneman et al., 2003]

Captures which attribute values in the result of a query have been copied from which attribute values in the instance. Representation: $\mathbb{P}(Attr(I))$

- Where: Operator-level syntax-based annotation propagation
- IWhere: Insensitive variant: Union of Where for all Q' with $Q' \equiv Q$



Properties	Semantics	Insensitivity	Stability	Conclusions
00	0	•0	0	0
Where				

Where

- Sensitive, traditionally attributed to being based on query syntax
- Depends on the internal data-flow inside the query
 - How values are routed through the query

IWhere

- Insensitive by combining Where for all equivalent queries
- Counterintuitive effect that if (R, t, A) is in the provenance then all (R, t', A) with $t \cdot A = t' \cdot A$ are in the provenance too.
 - Reason: Can construct equivalent query adding self-join on A



Properties	Semantics	Insensitivity	Stability	Conclusions
00	0	•O	0	0
Where				

$$Q_{a} = R$$

$$Q_{b} = \pi_{A,B}(R \bowtie_{A=C} \pi_{A \to C,B \to D}(R))$$

$$r$$

R			Q_a	& Q	b
	Α	В		Α	
r_1	1	2	a ₁	1	
<i>r</i> ₂	1	3	a ₂	1	
r ₃	2	3	a ₃	2	
<i>r</i> 4	2	5	a4	2	

B 2 3

3 5

$$Where(Q_a, a_1, A) = \{(r_1, A)\}$$

$$Where(Q_b, a_1, A) = \{(r_1, A), (r_2, A)\}$$

$$IWhere(Q_a, a_1, A) = IWhere(Q_b, a_1, A) = \{(r_1, A), (r_2, A)\}$$



Properties	Semantics	Insensitivity	Stability	Conclusions
00	0	00	0	0
Arguments f	for Insensit	ivity		

- Traditionally observed as advantageous in database research
 Tradition not a solid argument
- External implementation of *sensitive* semantics. Computing provenance for a query different from the one that will be executed by DBMS

 \Rightarrow No way to solve this, but provenance based on user query seems to be reasonable

- Implementation of sensitive semantics in DB-engine limits optimizer search space
 - \Rightarrow Insensitive semantics may be harder to compute
 - \Rightarrow Lack of practical experience
 - \Rightarrow "Realistic" sensitivity example?



Properties	Semantics	Insensitivity	Stability	Conclusions
00	0	00	•	0
Instability	of IWhere			

• *IWhere* is union of *Where* for all equivalent queries



Properties	Semantics	Insensitivity	Stability	Conclusions
00	0	00	•	0
Instability	of IWhere			

- *IWhere* is union of *Where* for all equivalent queries
- e.g., SPJ and USPJ equivalences are different



Properties	Semantics	Insensitivity	Stability	Conclusions
00	0	00	•	0
Instability	of IWhere			

- *IWhere* is union of *Where* for all equivalent queries
- e.g., SPJ and USPJ equivalences are different
- e.g., union Q with a join of Q with some other relation



Properties	Semantics	Insensitivity	Stability	Conclusions
00	0	00	•	0
Instability	of IWhere			

- *IWhere* is union of *Where* for all equivalent queries
- e.g., SPJ and USPJ equivalences are different
- e.g., union Q with a join of Q with some other relation
- Let UWhere be IWhere for USPJ queries

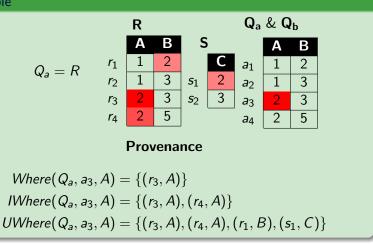


Properties	Semantics	Insensitivity	Stability	Conclusions
00	0	00	•	0
Instability	of IWhere			

- *IWhere* is union of *Where* for all equivalent queries
- e.g., SPJ and USPJ equivalences are different
- e.g., union Q with a join of Q with some other relation
- Let UWhere be IWhere for USPJ queries
- → UWhere an attribute value is annotated with all annotations from attribute positions that have the same value



Properties	Semantics	Insensitivity	Stability	Conclusions
00	0	00	•	0
Instability	of IWhere			





Properties	Semantics	Insensitivity	Stability	Conclusions
	O	00	⊙	•
Conclusions				

Take Away Messages

- Be careful how to achieve a property
- Insensitivity less applicable to semantics that address internal data-flow
 - Queries with the same external but possibly different internal behaviour have the same provenance

Some Things I'd Like to See

- "Declarative" Semantics \Rightarrow derive operator-level construction
- Semantics model processing, but have a insensitive "core"
- The never-ending quest: Deal with Negation
- Other data-models (order)



Questions	Overview	Properties	Semantics	Insensitivity
•	O	○		00000
Questions				

Semantic	cs	Sound	Complete	Responsible	Insensitive (set)	Insensitive (bag)	Stable
	Wit	-	Х	-	Х	Х	Х
Why	Why	-	Х	-	-	?	Х
	IWhy	-	Х	Х	Х	Х	Х
Where	Where	-	-	-	-	?	Х
vvnere	IWhere	-	-	-	Х	Х	-
How		-	Х	-	-	Х	Х
	Lineage	Х	Х	-	-	-	Х
Lineage-based	PI-CS	Х	Х	-	-	-	Х
	C-CS	Х	-	-	-	-	Х
Causality		-	Х	Х	Х	Х	Х



Questions	Overview	Properties	Semantics	Insensitivity
0	•	0	0000	00000
Semantic	s Summary			

Representation	
	Used by
$\mathbb{P}(Attr(I))$	Where,
	IWhere
$\mathbb{P}(\mathbb{P}(Tuple(I)))$	Wit,
	Why,
	IWhy
$\mathbb{N}[Tuple(I)]$	How
$\{ < R_1^*, \ldots, R_n^* > \mid R_i^* \subseteq R_i(Q) \}$	Lineage
$\mathbb{P}(\{\langle t_1,\ldots,t_n\rangle \ t_i\in R_i(Q)\lor t_i=\bot\})$	PI-CS, C-
	CS
$\mathbb{P}(Tuple(I))$	Causality



Questions	Overview	Properties	Semantics	Insensitivity
0	0	•	0000	00000
Sound Co	omplete, Res	ponsible		

• **Sound:** Provenance of *t* produces nothing different from *t*.

• $t' \neq t \Rightarrow t' \notin Q(\mathcal{P}(Q, I, t))$



Questions	Overview	Properties	Semantics	Insensitivity
0	0	•	0000	00000
Sound C	omplete, Res	nonsible		
Sound, C	Jumpiele, ives			

- **Sound:** Provenance of *t* produces nothing different from *t*.
 - $t' \neq t \Rightarrow t' \notin Q(\mathcal{P}(Q, I, t))$
- Caveat: Semantics of evaluating query over provenance





Questions	Overview	Properties	Semantics	Insensitivity
0	0	•	0000	00000
Sound Co	omplete, Res	ponsible		

• **Sound:** Provenance of *t* produces nothing different from *t*.

•
$$t' \neq t \Rightarrow t' \notin Q(\mathcal{P}(Q, I, t))$$

- **Complete:** Provenance of *t* produces at least *t*
 - $t \in Q(\mathcal{P}(Q, I, t))$



Questions	Overview	Properties	Semantics	Insensitivity
0	0	•	0000	00000
Sound Co	omplete, Res	ponsible		

• **Sound:** Provenance of *t* produces nothing different from *t*.

•
$$t' \neq t \Rightarrow t' \notin Q(\mathcal{P}(Q, I, t))$$

• **Complete:** Provenance of *t* produces at least *t*

•
$$t \in Q(\mathcal{P}(Q, I, t))$$

• Caveat: Semantics of evaluating query over provenance





Reexamining some Holy Grails of Provenance

Questions	Overview	Properties	Semantics	Insensitivity
0	0	•	0000	00000
Sound Co	omplete, Res	ponsible		

• **Sound:** Provenance of *t* produces nothing different from *t*.

•
$$t' \neq t \Rightarrow t' \notin Q(\mathcal{P}(Q, I, t))$$

• **Complete:** Provenance of *t* produces at least *t*

• $t \in Q(\mathcal{P}(Q, I, t))$

• **Responsible:** Every tuple in the provenance of *t* is necessary to derive *t*





• **Sound:** Provenance of *t* produces nothing different from *t*.

•
$$t' \neq t \Rightarrow t' \notin Q(\mathcal{P}(Q, I, t))$$

• **Complete:** Provenance of *t* produces at least *t*

• $t \in Q(\mathcal{P}(Q, I, t))$

- **Responsible:** Every tuple in the provenance of *t* is necessary to derive *t*
- Caveat: ... from every alternative derivation in the provenance ...
 - $\bullet \ \Rightarrow$ factor provenance into alterative derivations
- Caveat: Different ways to model that.

• E.g., $\forall t' \in \mathcal{P}(Q, I, t) : t \notin Q(\mathcal{P}(Q, I, t) - \{t'\})$

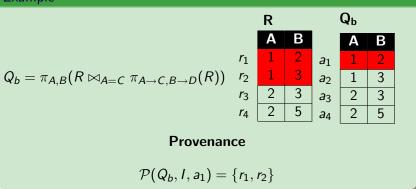








Example





Questions	Overview	Properties	Semantics	Insensitivity
0	0	0	0000	00000
1. A.	1. Sec.			
Lineage-	based			

Lineage-based [Cui et al., 2000]

Operator-level declarative semantics similar to *Why*. Provenance is modeled as a list of subsets of the relations accessed by the query (leafs of the algebra tree of Q) Representation: $\{ < R_1^*, \ldots, R_n^* > | R_i^* \subseteq R_i(Q) \}$

• Lineage: List of subsets of the algebra-tree nodes



Questions	Overview	Properties	Semantics	Insensitivity
0	0	0	0000	00000
Lineage-b	based			

Lineage-based [Glavic et al., 2009]

Provenance is modeled as a set of **witness lists**. A witness list is a list of tuples - one from each relation accessed by the query. Representation: $\mathbb{P}(\{ < t_1, ..., t_n > | t_i \in R_i(Q) \lor t_i = \bot\})$

- **PI-CS**: *Lineage* with different representation and broader query language coverage
- C-CS: Similar to Where but with tuple granularity



Questions	Overview	Properties	Semantics	Insensitivity
0	0	0	0000	00000
1. A.	100 Aug. 100			
Lineage-b	based			

Example

$$Q_{c} = \pi_{A}(R) \cup \pi_{B}(R) \quad \begin{matrix} R \\ r_{1} \\ r_{2} \\ q_{d} = \pi_{A}(R \bowtie_{B=C} S) \end{matrix} \stackrel{r_{1}}{r_{2}} \begin{matrix} \frac{1}{2} \\ \frac{1}{2} \\ r_{3} \\ \frac{2}{2} \\ \frac{5}{5} \end{matrix} \stackrel{s_{1}}{s_{2}} \begin{matrix} 2 \\ c_{2} \\ c_{2} \\ c_{3} \\ c_{4} \end{matrix} \stackrel{s_{1}}{s_{2}} \begin{matrix} Q_{d} \& Q_{e} \\ A \\ B \\ c_{1} \\ 1 \\ c_{2} \\ c_{2} \\ c_{3} \\ c_{4} \end{matrix} \stackrel{s_{1}}{s_{2}} \begin{matrix} Q_{d} \& Q_{e} \\ A \\ c_{1} \\ 1 \\ c_{2} \\ c_{2} \\ c_{3} \\ c_{4} \\ c_{5} \\ c_{4} \\ c_{5} \\ c_{4} \\ c_{5} \\ c_{5} \\ c_{6} \\ c_{7} \\ c_$$

$$C - CS(Q_d, d_1) = \{ < r_1, \bot >, < r_2, \bot > \}$$

$$Lineage(Q_d, d_1) = \{ < r_1, \bot >, < r_2, \bot > \}$$

Questions	Overview	Properties	Semantics	Insensitivity
0	0	0	0000	00000
Why				

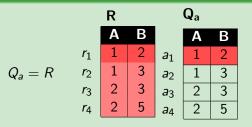
Why [Buneman et al., 2003]

Why-provenance models provenance as a set of **witnesses**. A witness w for a tuple t is a subset of the instance I where $t \in Q(w)$. Representation: $\mathbb{P}(\mathbb{P}(Tuple(I)))$

- Wit: Set of all witnesses
- Why: Query-syntax based "proof-witnesses"
- IWhy: Minimal elements from Wit resp. Why



Questions O	Overview O	Properties ○	Semantics	Insensitivity
Why				



$$Wit(Q_a, a_1) = \{J \mid J \subset R \land r_1 \in J\}$$
$$Why(Q_a, a_1) = \{\{r_1\}\}$$
$$IWhy(Q_a, a_1) = \{\{r_1\}\}$$



Questions	Overview	Properties	Semantics	Insensitivity
				00000
How				

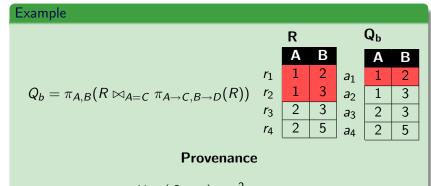
Provenance Semirings [Green et al., 2007]

Tuples of relations annotated with elements from a semiring. Annotation propagation defined for positive relational algebra as operations of the semiring (set difference and aggregration later). Representation: $\mathbb{N}(Tuple(I))$

• How: Most general form of annotations: polynomials over variable representing the instance tuples. Addition indicates alternative use of tuples; multiplication conjunctive use.



Questions	Overview	Properties	Semantics	Insensitivity
0	0	0	0000	00000
How				
11000				



$$How(Q_b,a_1)=r_1^2+r_1\times r_2$$



Questions	Overview	Properties	Semantics	Insensitivity
0	0	0	000	00000
Causality				

Causality [Meliou et al., 2010]

Provenance is modeled as a set of causes. A **cause** $c \in I$ for a tuple *t* is defined as follows:

$$t \notin Q(I - \{c\})$$

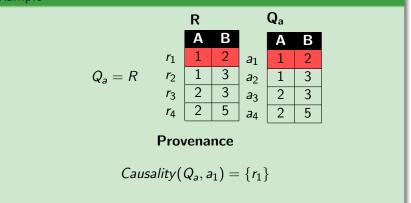
② there exists a set $C \subset I$ called contingency so that $t \in Q(I - C)$ and $t \notin Q(I - C - \{c\})$

Representation: $\mathbb{P}(Tuple(I))$

• Causality: Set of all causes



Questions	Overview	Properties	Semantics	Insensitivity
0	0	0	0000	00000
Causality				





Questions	Overview	Properties	Semantics	Insensitivity •0000
and the second		1. S.		

Insensitivity to Query Rewrite

Insensitivity to Query Rewrite

• Equivalent queries have the same provenance

•
$$Q \equiv Q' \Rightarrow \mathcal{P}(Q, I, t) = \mathcal{P}(Q', I, t)$$

• Set resp. Bag semantics

Overview

- Insensitive (Set, Bag) Wit, IWhy, IWhere, Causality
- Insensitive (Bag) How
- Sensitive: Lineage, PI-CS, C-CS, Where



Questions	Overview	Properties	Semantics	Insensitivity
O	○	O		○●○○○
Why				

Wit

• Defined over black-box behaviour of query \Rightarrow trivially insensitive

Why

- Sensitive, traditionally attributed to being based on query syntax
- Why may contain tuples that do not contribute to t
- $\bullet \Rightarrow$ Equivalent queries that apply redundant computations may contain larger provenance
- Caveat: But why does this argument not apply to Wit?
- Positive queries: super-set of a witness is also a witness ⇒ tuples used by redundant computations are in *Wit*



IWhy

Questions	Overview	Properties	Semantics	Insensitivity
○	O	○		○●○○○
Why				

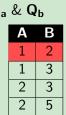
$$Q_{a} = R$$
$$Q_{b} = \pi_{A,B}(R \bowtie_{A=C} \pi_{A \to C,B \to D}(R))$$

	R		Q_a
	Α	В	
1	1	2	a ₁
2	1	3	a ₂
3	2	3	a ₃
4	2	5	а ₂ а ₃ а ₄

r r

r

r



$$Wit(Q_a, a_1) = Wit(Q_b, a_1) = \{J \mid J \subset R \land r_1 \in J\}$$
$$Why(Q_a, a_1) = \{\{r_1\}\}$$
$$Why(Q_b, a_1) = \{\{r_1\}, \{r_1, r_2\}\}$$
$$Why(Q_a, a_1) = IWhy(Q_b, a_1) = \{\{r_1\}\}$$



Questions	Overview	Properties	Semantics	Insensitivity
0	0	0	0000	00000
	and the second			
Lineage-b	based			
Encage r	Juseu			

- Provenance representation based on query syntax
- Trivial examples for sensitivity based on reordering of the arguments of commutative operators



Questions	Overview	Properties	Semantics	Insensitivity
0	0	0	0000	00000
Lineage-	pased			

O

$$\begin{aligned} &Lineage(Q_d, d_1) = <\{r_1, r_2\}, \{s_1, s_2\} > \\ &Lineage(Q_e, d_1) = <\{s_1, s_2\}, \{r_1, r_2\} > \\ &PI - CS(Q_d, d_1) = \{, \} \\ &PI - CS(Q_e, d_1) = \{, \} \end{aligned}$$



Questions	Overview	Properties	Semantics	Insensitivity
O	○	O	0000	○○○●○
How				

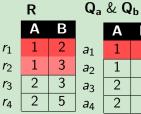
- Sensitive (Set):
- Insensitive (Bag): Operator semantics defined to take bag semantics into account



Questions O	Overview ○	Properties ○	Semantics	Insensitivity
How				

$$Q_{a} = R$$

$$Q_{b} = \pi_{A,B}(R \bowtie_{A=C} \pi_{A \to C,B \to D}(R))$$



$$How(Q_a, a_1) = r_1$$
$$How(Q_b, a_1) = r_1^2 + r_1 \times r_2$$

Questions	Overview	Properties	Semantics	Insensitivity
0	0	0	0000	00000
Causality				
Causanty				

• Trivially insensitive: Defined over the black-box behaviour of a query



Questions	Overview ○	Properties O	Semantics	Insensitivity ○○○○●
Causality				

