



# Data Quality Problems (DQPs) at the Instance Level

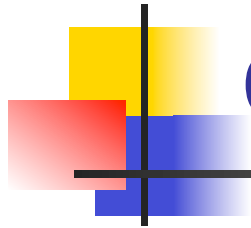
---

Paulo Oliveira



Toward Efficient Portuguese and Brazilian Electricity Markets Workshop

September 25, 2013

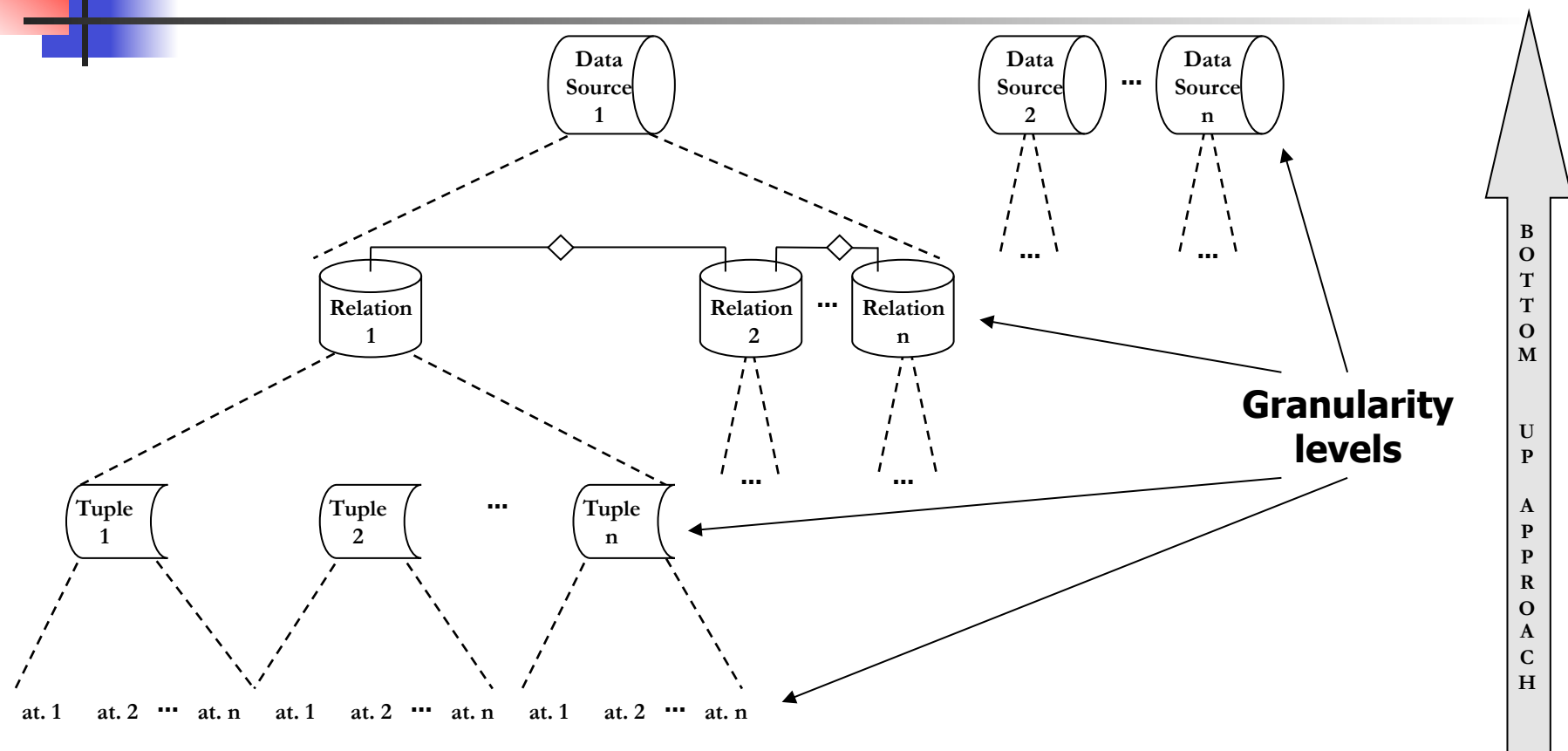


# Concept of DQP

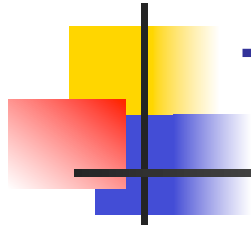
---

- DQP is restricted to problems related with the quality of the data values or instances
- Data values are affected by different kinds of quality problems (errors, anomalies, or *dirty*)
- DQPs arise in:
  - single data source
  - data migration
  - integration of multiple data sources
  - data-based projects
    - data warehouses
    - data mining
- Important due to the *Garbage In Garbage Out* principle

# Approach to Identify DQPs



- Based on this model of data organization
- Identification of all the DQPs that can be found at each granularity level

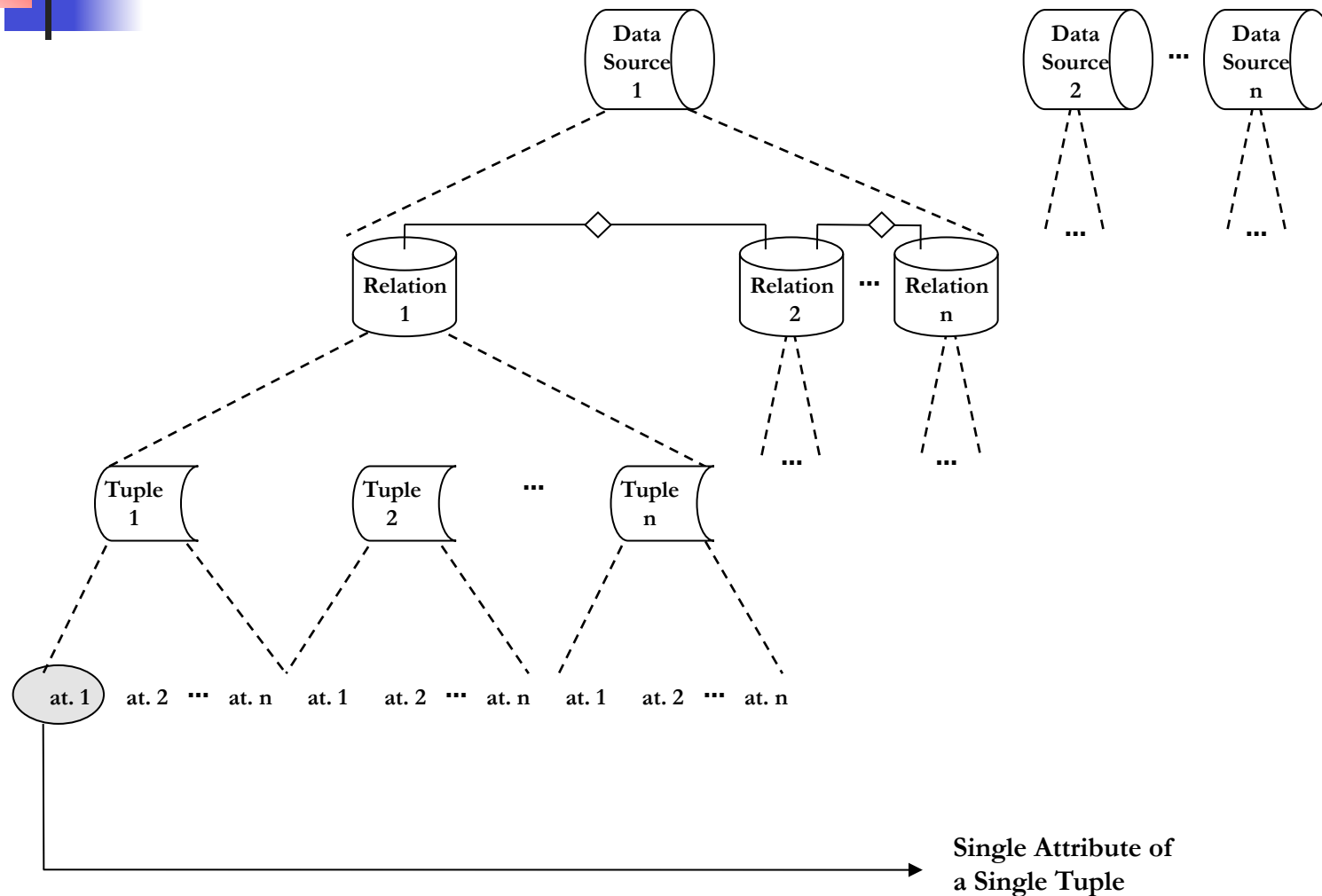


# Taxonomy of DQPs

---

- Covers
  - DQPs of tabular data
  - All data types, except multimedia
- Is complete
- Is limited to DQPs related with the attribute values

# DQPs in a Single Attribute of a Single Tuple





# Missing Value

---

	$at_1$	$at_2$	<b>name</b>	...	$at_n$
$t_1$	xxx	xxx	Carl Louis	...	xxx
$t_2$	xxx	xxx		...	xxx
⋮	⋮	⋮	⋮	⋮	⋮
$t_n$	xxx	xxx	Ben Johnson	...	xxx



# Syntax Violation

---

	at <sub>1</sub>	at <sub>2</sub>	<b>order_date</b>	...	at <sub>n</sub>
t <sub>1</sub>	xxx	xxx	24/10/2012	...	xxx
t <sub>2</sub>	xxx	xxx	2012/10/26	...	xxx
.	.	.	.	.	.
.	.	.	.	...	.
.	.	.	.	.	.
t <sub>n</sub>	xxx	xxx	27/10/2012	...	xxx



# Domain Violation

---

	at <sub>1</sub>	at <sub>2</sub>	<b>ordered_quantity</b>	...	at <sub>n</sub>
t <sub>1</sub>	xxx	xxx	4	...	xxx
t <sub>2</sub>	xxx	xxx	-1	...	xxx
.	.	.	.	.	.
.	.	.	.	...	.
.	.	.	.	.	.
t <sub>n</sub>	xxx	xxx	1	...	xxx





# Misspelling Error

---

	$at_1$	$at_2$	<b>city</b>	...	$at_n$
$t_1$	xxx	xxx	New York	...	xxx
$t_2$	xxx	xxx	Bostom	...	xxx
.	.	.	.	.	.
.	.	.	.	...	.
.	.	.	.	.	.
$t_n$	xxx	xxx	Washington	...	xxx



# Overloaded Attribute

---

	at <sub>1</sub>	at <sub>2</sub>	name	...	at <sub>n</sub>
t <sub>1</sub>	xxx	xxx	Bill Clinton	...	xxx
t <sub>2</sub>	xxx	xxx	Dr. Barack Obama	...	xxx
.	.	.	.	.	.
.	.	.	.	...	.
.	.	.	.	.	.
t <sub>n</sub>	xxx	xxx	George Bush	...	xxx



# Incomplete Value

---

	$at_1$	$at_2$	<b>address</b>	...	$at_n$
$t_1$	xxx	xxx	Sun Street, 123	...	xxx
$t_2$	xxx	xxx	Flowers Street	...	xxx
⋮	⋮	⋮	⋮	⋮	⋮
$t_n$	xxx	xxx	Baker Street, 321	...	xxx

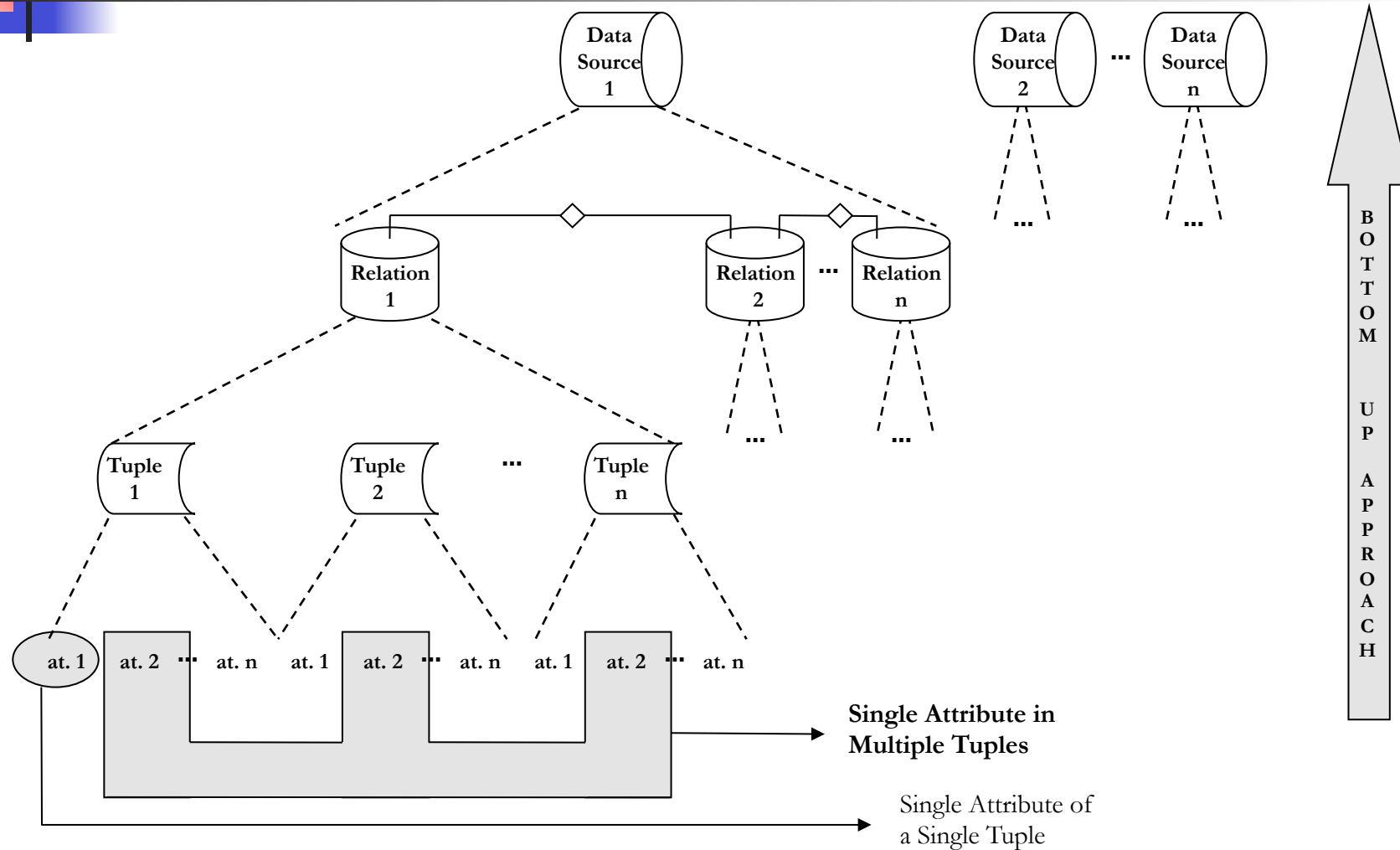


# Wrong Value

---

	$at_1$	$at_2$	<b>Marital Status</b>	...	$at_n$
$t_1$	xxx	xxx	Married	...	xxx
$t_2$	xxx	xxx	Single	...	xxx
⋮	⋮	⋮	⋮	⋮	⋮
$t_n$	xxx	xxx	Divorced	...	xxx

# DQPs in a Single Attribute in Multiple Tuples

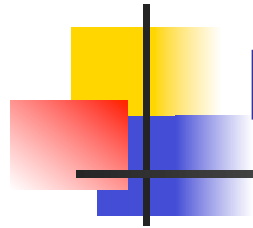




# Uniqueness Violation

---

	name	at <sub>2</sub>	<b>taxpayer_nr</b>	...	at <sub>n</sub>
t <sub>1</sub>	George Clooney	xxx	196 567 931	...	xxx
t <sub>2</sub>	Harrison Ford	xxx	187 323 436	...	xxx
t <sub>3</sub>	Brad Pitt	xxx	205 239 894	...	xxx
⋮	⋮	⋮	⋮	⋮	⋮
t <sub>n</sub>	Judy Foster	xxx	187 323 436	...	xxx



# Existence of Synonyms

	$at_1$	$at_2$	<b>job</b>	...	$at_n$
$t_1$	xxx	xxx	Researcher	...	xxx
$t_2$	xxx	xxx	Schoolmaster	...	xxx
$t_3$	xxx	xxx	Electrician	...	xxx
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$t_n$	xxx	xxx	Teacher	...	xxx



# Violation of Business Rule

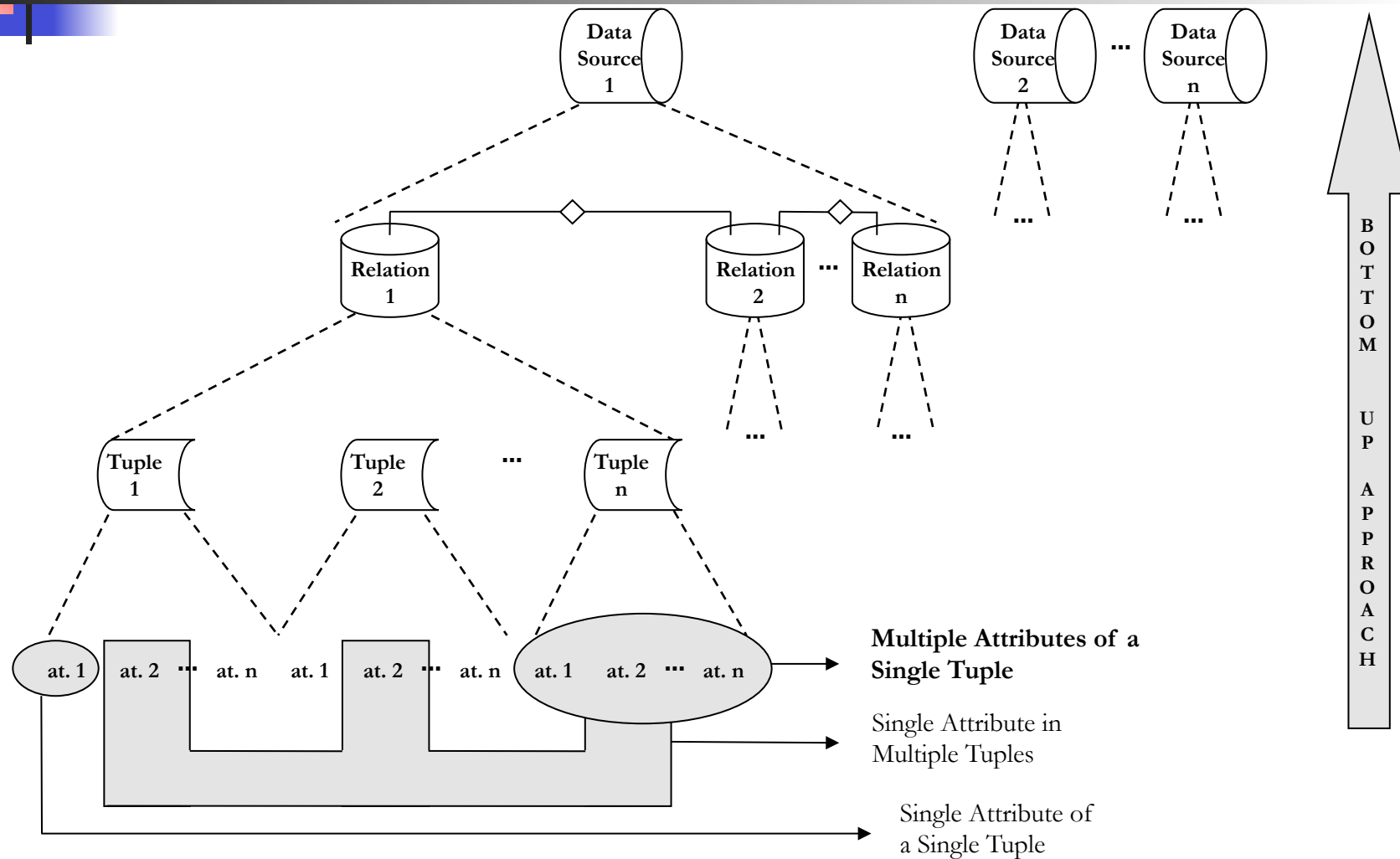
---

	invoice_nr	at <sub>2</sub>	invoice_date	...	at <sub>n</sub>
t <sub>1</sub>	20121100	xxx	25/10/2012	...	xxx
t <sub>2</sub>	20121101	xxx	24/10/2012	...	xxx
t <sub>3</sub>	20121102	xxx	25/10/2012	...	xxx
⋮	⋮	⋮	⋮	⋮	⋮
t <sub>n</sub>	20121178	xxx	05/11/2012	...	xxx

The values of *invoice\_date* must appear by ascending order !



# DQPs in Multiple Attributes of a Single Tuple





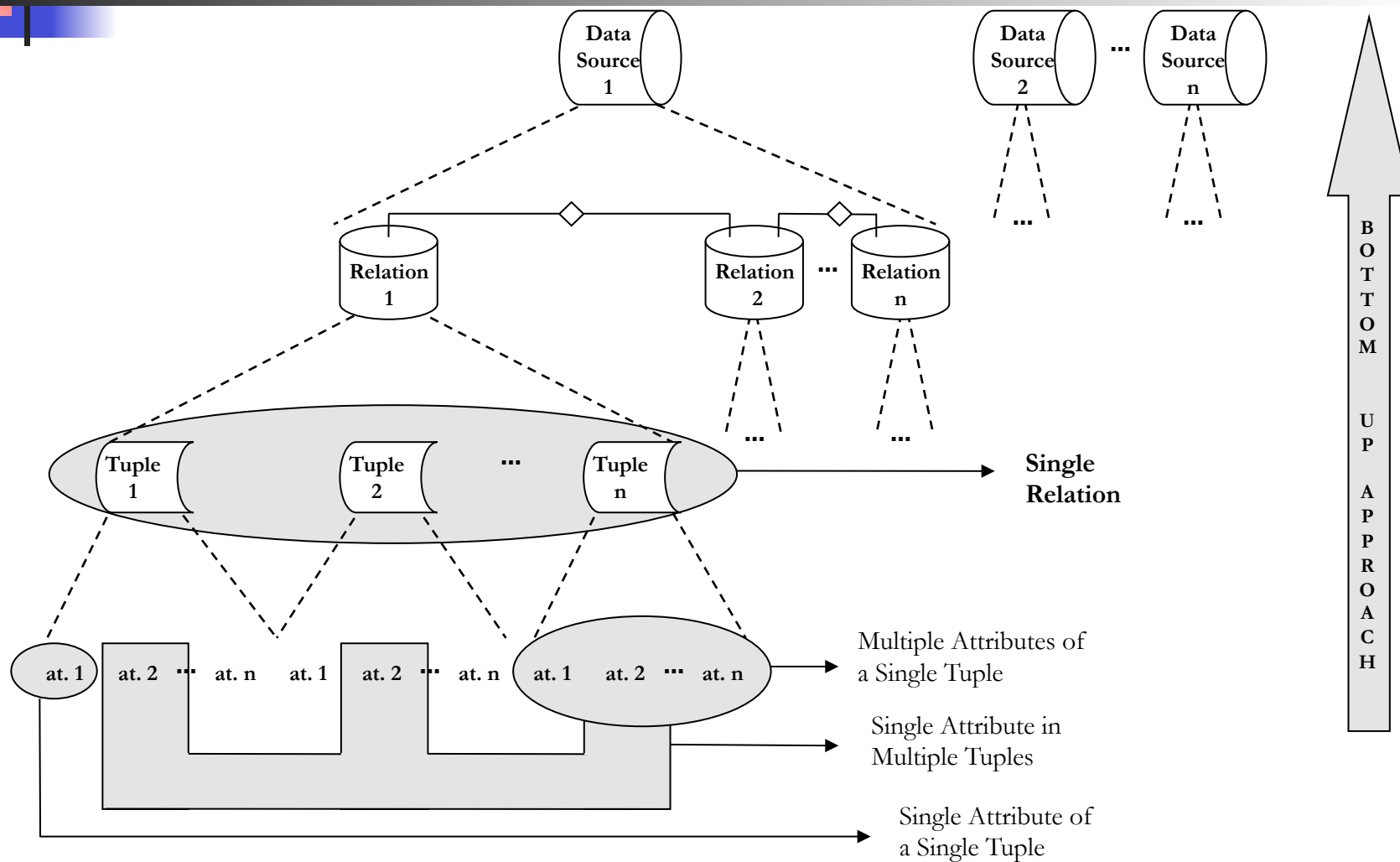
# Violation of Business Rule

---

	$at_1$	<b>quantity</b>	<b>unit_price</b>	<b>total_prod</b>	...	$at_n$
$t_1$	xxx	2	3	6	...	xxx
$t_2$	xxx	2	5	5	...	xxx
$t_3$	xxx	3	4	12	...	xxx
⋮	⋮	⋮	⋮		⋮	⋮
$t_n$	xxx	2	1	2	...	xxx

*total\_prod* must be equal  
to *quantity* \* *unit\_price* !

# DQPs at the Single Relation Level





# Violation of Functional Dependency

---

	$at_1$	<b>zip_code</b>	<b>city</b>	...	$at_n$
$t_1$	xxx	4000	Oporto	...	xxx
$t_2$	xxx	4000	Lisbon	...	xxx
$t_3$	xxx	1000	Lisbon	...	xxx
⋮	⋮	⋮	⋮	⋮	⋮
$t_n$	xxx	4000	Oporto	...	xxx



# Duplicate Tuples (Equal)

---

	id	name	address	taxpayer_nr	...	at <sub>n</sub>
t <sub>1</sub>	xxx	xxx	xxx	xxx	...	xxx
t <sub>2</sub>	10	<b>Cliff Barnes</b>	<b>Flowers Street, 123</b>	<b>205 239 894</b>	...	xxx
⋮	⋮	⋮	⋮	⋮	⋮	⋮
t <sub>n</sub>	72	<b>Cliff Barnes</b>	<b>Flowers Street, 123</b>	<b>205 239 894</b>	...	xxx



# Duplicate Tuples (Approximate)

---

	id	name	address	taxpayer_nr	...	at <sub>n</sub>
t <sub>1</sub>	xxx	xxx	xxx	xxx	...	xxx
t <sub>2</sub>	10	<b>Cliff Barnes</b>	<b>Flowers Street, 123</b>	<b>205 239 894</b>	...	xxx
⋮	⋮	⋮	⋮	⋮	⋮	⋮
t <sub>n</sub>	72	<b>C. Barnes</b>	<b>Flowers St., 123</b>	<b>205 239 894</b>	...	xxx



# Duplicate Tuples (Inconsistent)

---

	id	name	address	taxpayer_nr	...	at <sub>n</sub>
t <sub>1</sub>	xxx	xxx	xxx	xxx	...	xxx
t <sub>2</sub>	10	Cliff Barnes	<b>Flowers Street, 123</b>	205 239 894	...	xxx
⋮	⋮	⋮	⋮		⋮	⋮
t <sub>n</sub>	72	Cliff Barnes	<b>Sun Street, 321</b>	205 239 894	...	xxx



# Violation of Business Rule

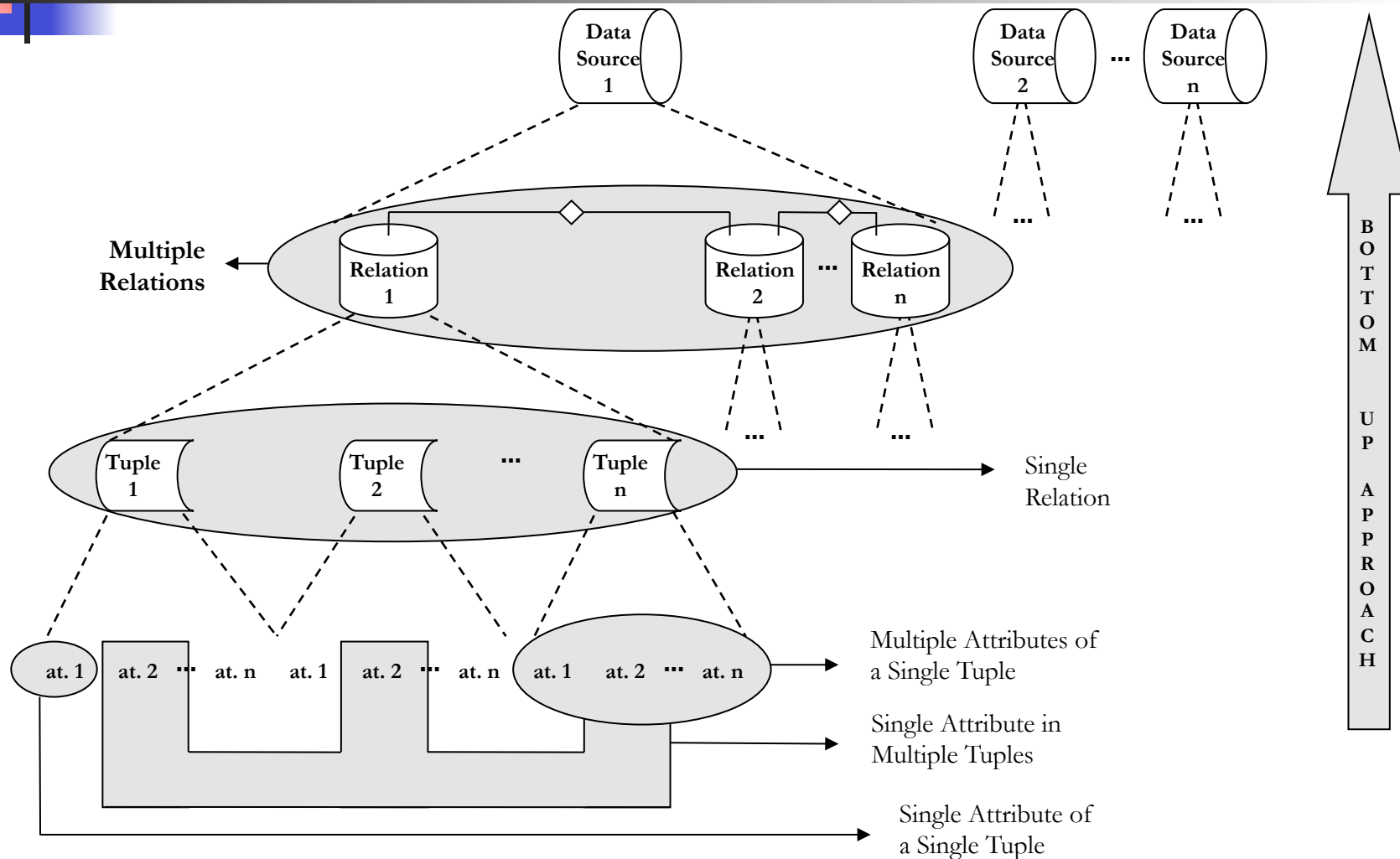
---

	$at_1$	$at_2$	...	$at_n$
$t_1$	xxx	xxx	...	xxx
$t_2$	xxx	xxx	...	xxx
⋮	⋮	⋮	⋮	⋮
$t_{12}$	xxx	xxx	...	xxx

The number of product families  
(tuples) must not be superior to 10!



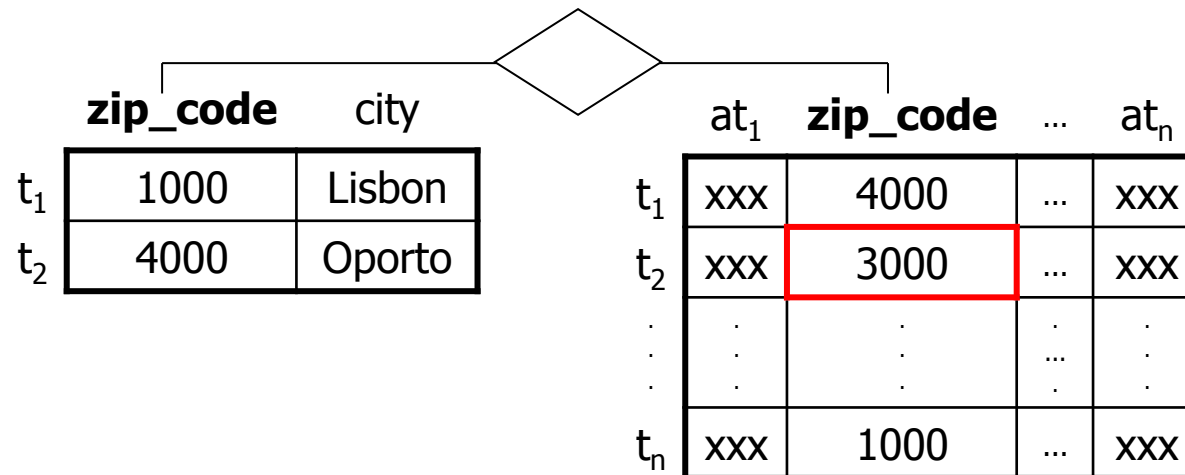
# DQPs at the Level of Multiple Relations





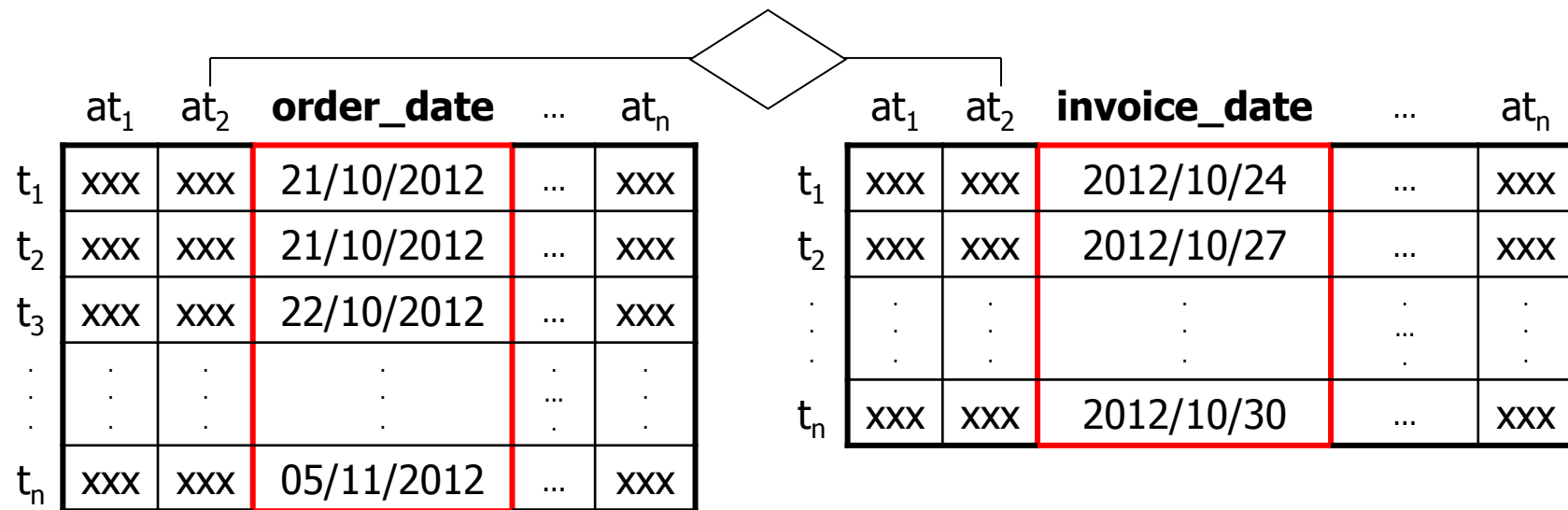
# Referential Integrity Violation

---

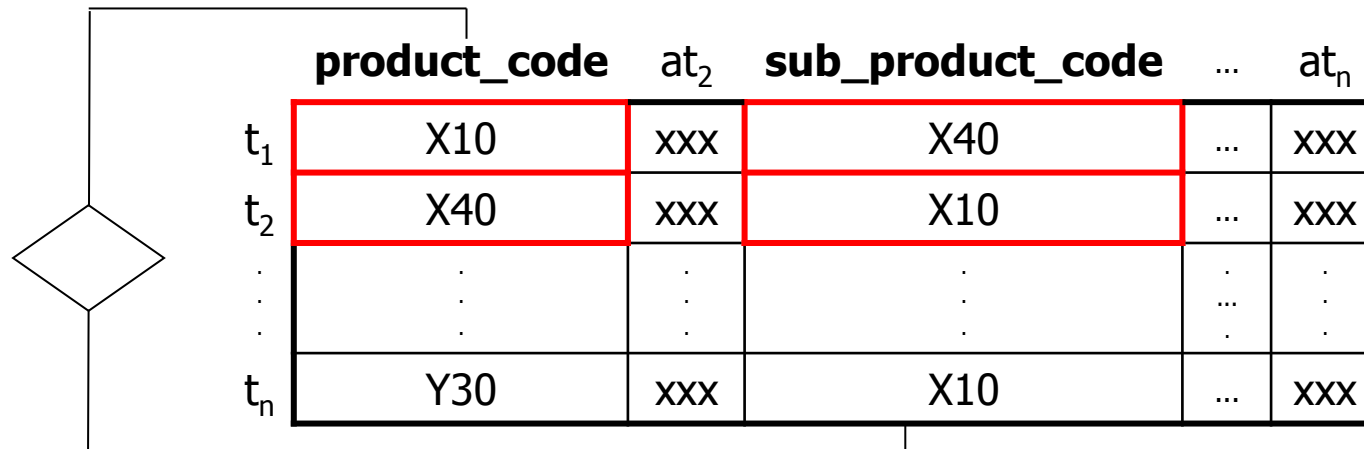




# Heterogeneity of Syntaxes

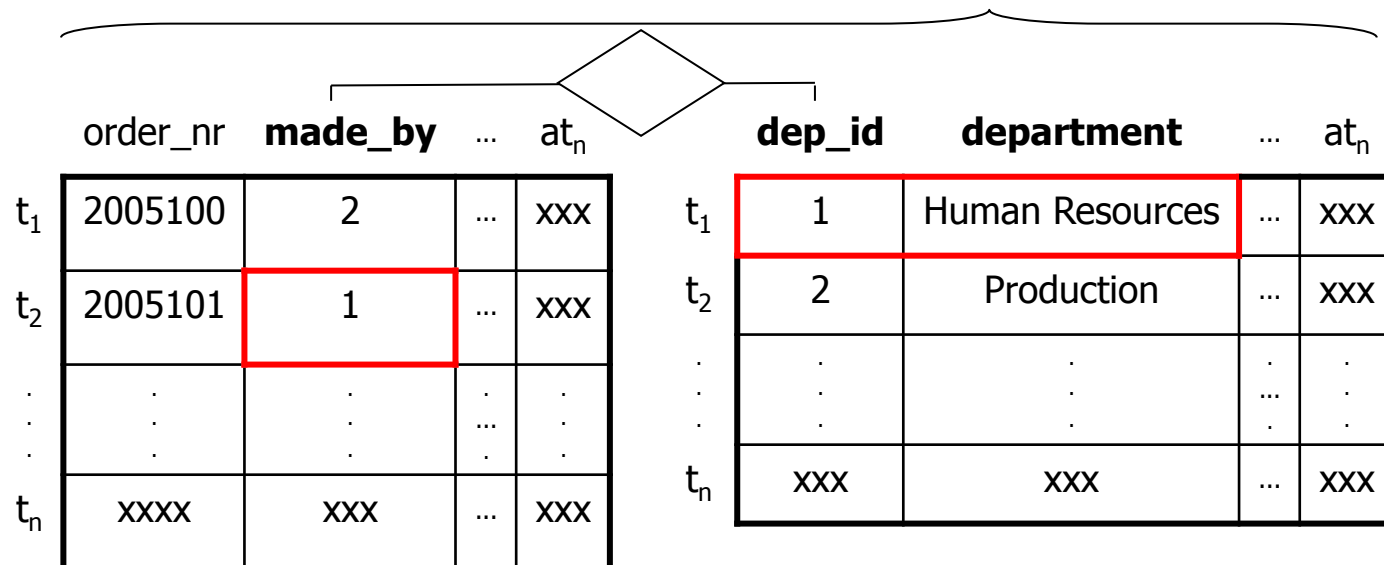


# Circularity among Tuples in a Self-Relationship

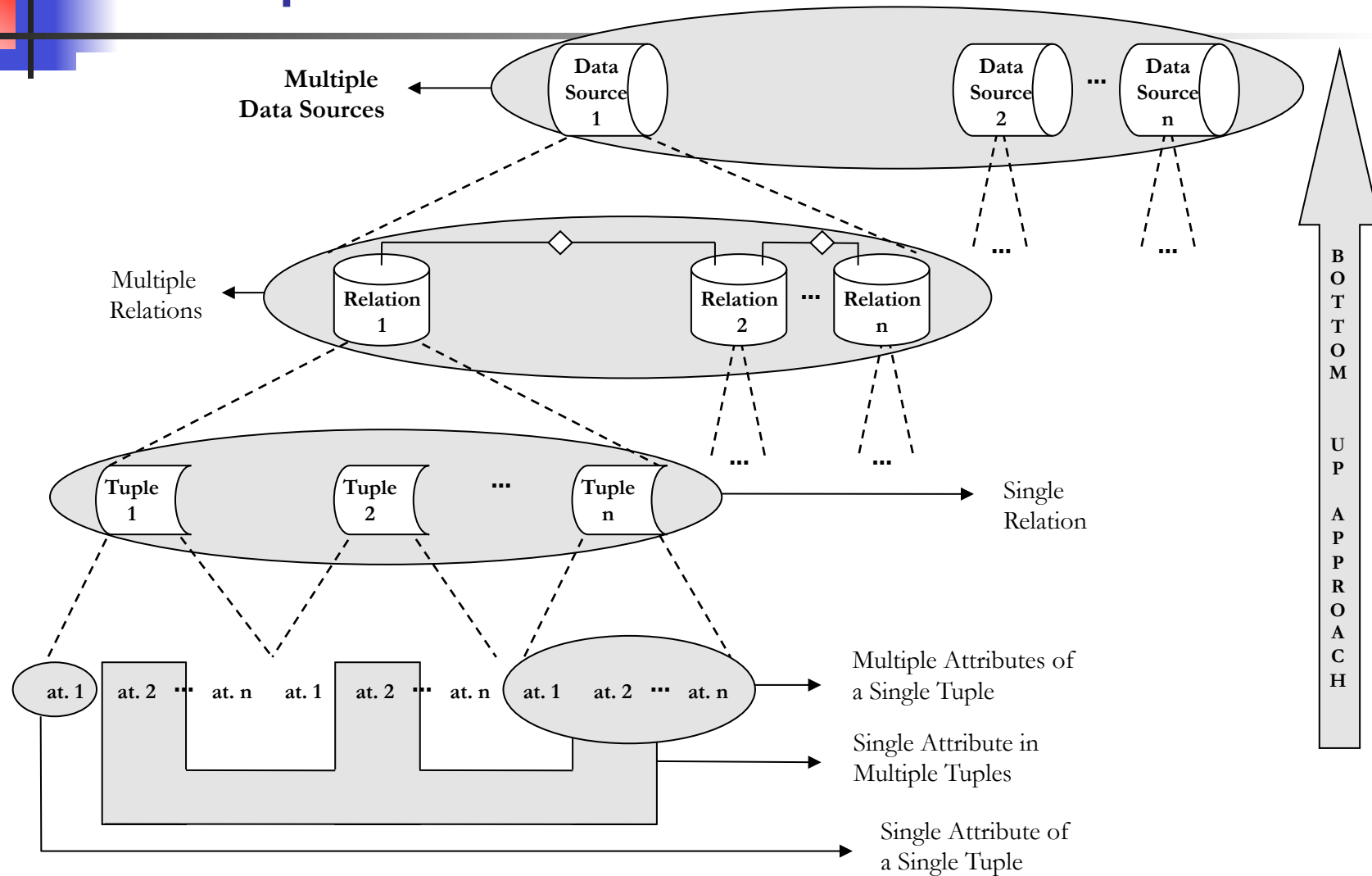


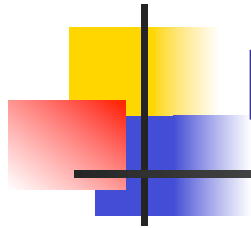
# Violation of Business Rule

Orders can only be made by the provisions department or production department !



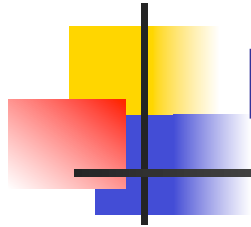
# DQPs at the Level of Multiple Data Sources





# Heterogeneity of Syntaxes

					at <sub>1</sub>	order_date	...	at <sub>n</sub>	
<b>DS<sub>1</sub></b>	t <sub>1</sub>	xxx	21/10/2012	...	xxx				<b>DS<sub>2</sub></b>
	t <sub>2</sub>	xxx	21/10/2012	...	xxx				
	⋮	⋮	⋮	⋮	⋮				
	⋮	⋮	⋮	⋮	⋮				
	t <sub>n</sub>	xxx	05/11/2012	...	xxx				
					at <sub>1</sub>	order_date	...	at <sub>n</sub>	
	t <sub>1</sub>	xxx	2012/10/21	...	xxx				
	t <sub>2</sub>	xxx	2012/10/21	...	xxx				
	⋮	⋮	⋮	⋮	⋮				
	⋮	⋮	⋮	⋮	⋮				
	t <sub>n</sub>	xxx	2012/11/05	...	xxx				



# Heterogeneity of Measure Units

**DS<sub>1</sub>**

	prod_id	unit_price	...	at <sub>n</sub>
t <sub>1</sub>	xpto	5	...	xxx
t <sub>2</sub>	ypto	12	...	xxx
⋮	⋮	⋮	⋮	⋮
t <sub>n</sub>	zpto	4	...	xxx

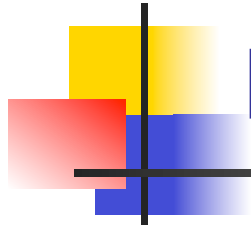
→ Dollars

**DS<sub>2</sub>**

	prod_id	unit_price	...	at <sub>n</sub>
t <sub>1</sub>	xpto	6.05	...	xxx
t <sub>2</sub>	ypto	13.20	...	xxx
⋮	⋮	⋮	⋮	⋮
t <sub>n</sub>	zpto	4.40	...	xxx

Euros ←





# Heterogeneity of Domains

		at <sub>1</sub>	gender	... at <sub>n</sub>
DS <sub>1</sub>	t <sub>1</sub>	xxx	M	... xxx
	t <sub>2</sub>	xxx	F	... xxx
	⋮	⋮	⋮	⋮
	⋮	⋮	⋮	⋮
	t <sub>n</sub>	xxx	M	... xxx

		at <sub>1</sub>	gender	... at <sub>n</sub>
DS <sub>2</sub>	t <sub>1</sub>	xxx	1	... xxx
	t <sub>2</sub>	xxx	2	... xxx
	⋮	⋮	⋮	⋮
	⋮	⋮	⋮	⋮
	t <sub>n</sub>	xxx	1	... xxx



# Existence of Synonyms

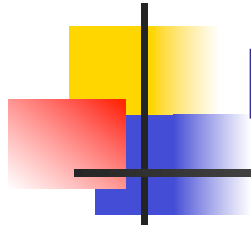
---

**DS<sub>1</sub>**

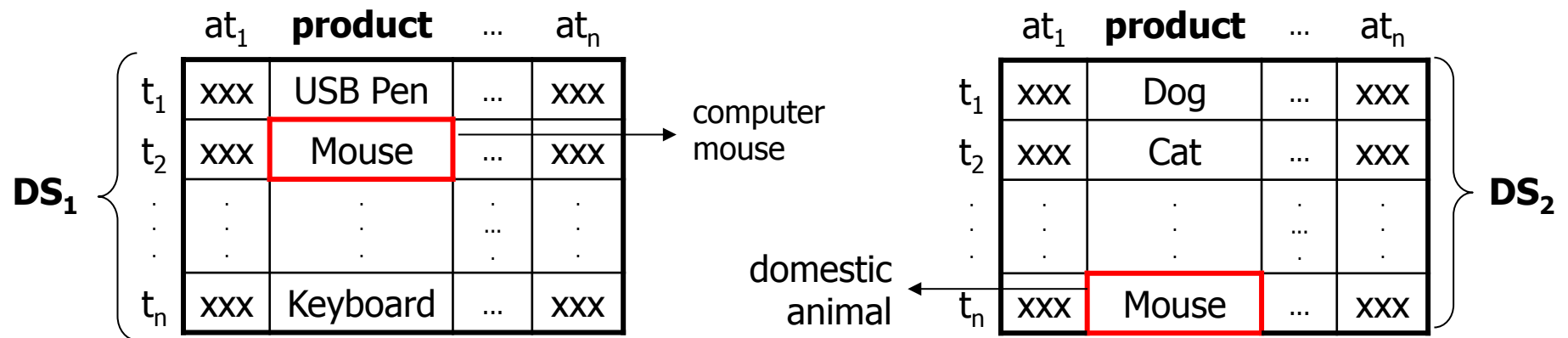
	at <sub>1</sub>	job	...	at <sub>n</sub>
t <sub>1</sub>	xxx	Policeman	...	xxx
t <sub>2</sub>	xxx	Teacher	...	xxx
⋮	⋮	⋮	⋮	⋮
t <sub>n</sub>	xxx	Researcher	...	xxx

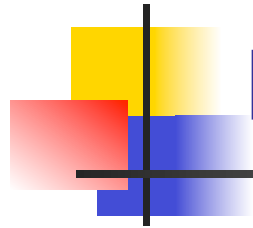
**DS<sub>2</sub>**

	at <sub>1</sub>	job	...	at <sub>n</sub>
t <sub>1</sub>	xxx	Electrician	...	xxx
t <sub>2</sub>	xxx	Plumber	...	xxx
⋮	⋮	⋮	⋮	⋮
t <sub>n</sub>	xxx	Schoolmaster	...	xxx



# Existence of Homonyms





# Duplicate Tuples (Equal)

	name	address	taxpayer_nr	...	at <sub>n</sub>
t <sub>1</sub>	xxx	xxx	xxx	...	xxx
t <sub>2</sub>	<b>C. Barnes</b>	<b>Flowers St. 123</b>	205 239 894	...	xxx
⋮	⋮	⋮	⋮	⋮	⋮
t <sub>n</sub>	xxx	xxx	xxx	...	xxx

DS<sub>1</sub>

	name	address	taxpayer_nr	...	at <sub>n</sub>
t <sub>1</sub>	xxx	xxx	xxx	...	xxx
t <sub>2</sub>	xxx	xxx	xxx	...	xxx
⋮	⋮	⋮	⋮	⋮	⋮
t <sub>n</sub>	<b>C. Barnes</b>	<b>Flowers St., 123</b>	205 239 894	...	xxx

DS<sub>2</sub>

# Duplicate Tuples (Approximate)

	name	address	taxpayer_nr	...	at <sub>n</sub>		name	address	taxpayer_nr	...	at <sub>n</sub>
t <sub>1</sub>	xxx	xxx	xxx	...	xxx	t <sub>1</sub>	xxx	xxx	xxx	...	xxx
t <sub>2</sub>	<b>Cliff Barnes</b>	<b>Flowers Street, 123</b>	205 239 894	...	xxx	t <sub>2</sub>	xxx	xxx	xxx	...	xxx
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
t <sub>n</sub>	xxx	xxx	xxx	...	xxx	t <sub>n</sub>	<b>C. Barnes</b>	<b>Flowers St., 123</b>	205 239 894	...	xxx

DS<sub>1</sub>

DS<sub>2</sub>



# Duplicate Tuples (Inconsistent)

	name	address	taxpayer_nr	...	at <sub>n</sub>
t <sub>1</sub>	xxx	xxx	xxx	...	xxx
t <sub>2</sub>	Cliff Barnes	<b>Flowers Street, 123</b>	205 239 894	...	xxx
⋮	⋮	⋮	⋮	⋮	⋮
t <sub>n</sub>	xxx	xxx	xxx	...	xxx

DS<sub>1</sub>

	name	address	taxpayer_nr	...	at <sub>n</sub>
t <sub>1</sub>	xxx	xxx	xxx	...	xxx
t <sub>2</sub>	xxx	xxx	xxx	...	xxx
⋮	⋮	⋮	⋮	⋮	⋮
t <sub>n</sub>	Cliff Barnes	<b>Sun Street, 123</b>	205 239 894	...	xxx

DS<sub>2</sub>

# Violation of Business Rule

The maximum number of projects for a manager is two !

<b>id_proj</b> <b>manager_name</b> ... <b>at<sub>n</sub></b>					<b>id_proj</b> <b>manager_name</b> ... <b>at<sub>n</sub></b>				
<b>DS<sub>1</sub></b>	t <sub>1</sub>	XY100	A. Schwarzenegger	... xxx	t <sub>1</sub>	AB900	Steven Segal	... xxx	<b>DS<sub>2</sub></b>
	t <sub>2</sub>	YW200	Sylvester Stallone	... xxx	t <sub>2</sub>	BC800	Sylvester Stallone	... xxx	
	.	.	.	... .	.	.	.	... .	
	.	.	.	... .	.	.	.	... .	
	t <sub>n</sub>	WZ300	Sylvester Stallone	... xxx	t <sub>n</sub>	CB700	Steven Segal	... xxx	



## Conclusion

---

- Taxonomy of DQPs is useful to:
  - Alert for the existence of many DQPs
  - Assess the coverage of existing DQ tools





# Data Quality Problems (DQPs) at the Instance Level

---

Paulo Oliveira



Toward Efficient Portuguese and Brazilian Electricity Markets Workshop

September 25, 2013



# Data Quality Problems (DQPs) at the Instance Level

---

Paulo Oliveira, Fátima Rodrigues and Pedro Henriques – “A Formal Definition of Data Quality Problems”. In *Proceedings of the 10th International Conference on Information Quality*, MIT, Boston, EUA, November of 2005. p. 13-26.