



# Provenance for database transformations or an algebraic view of annotated data

**Val Tannen et al : "Provenance Semirings"**  
[GreenKarvounarakis&T PODS 07]

# MOTIVATION



## ■ Annotations capture:

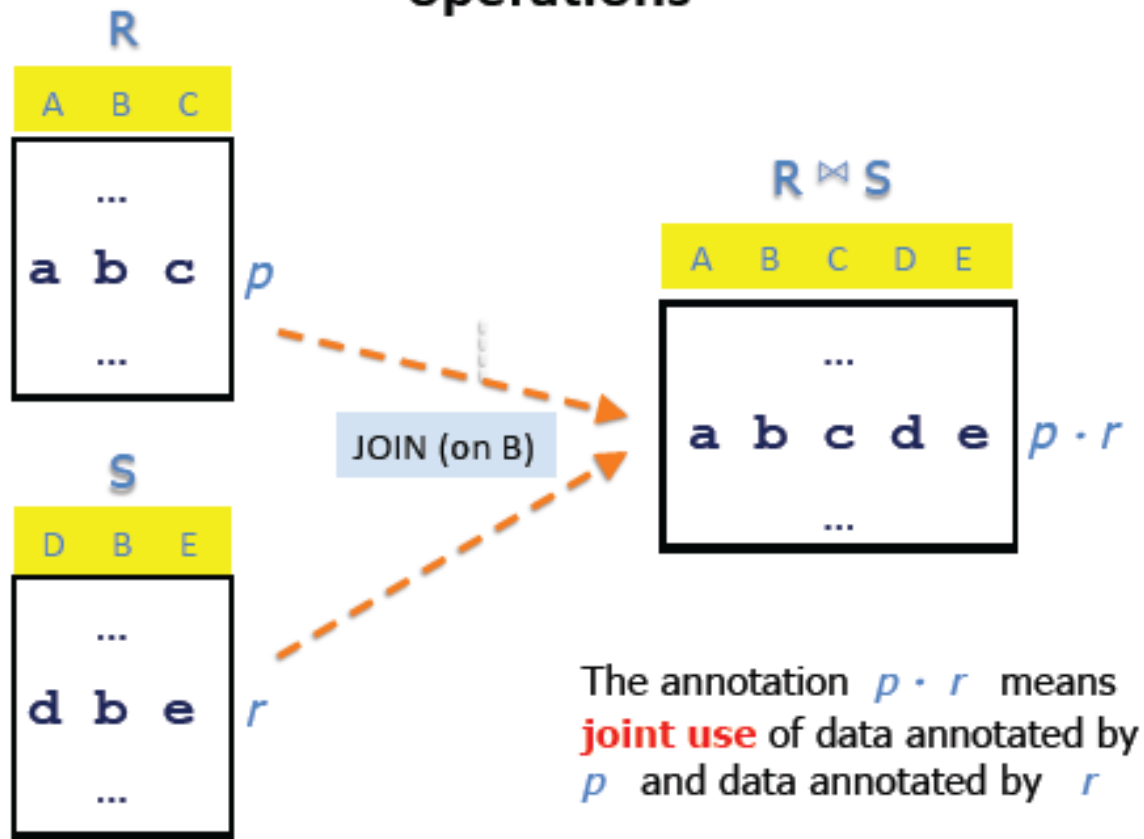
- PROVENANCE
- UNCERTAINTY
- TRUST
- SECURITY
- MULTIPLICITY

## ■ Semirings bring:

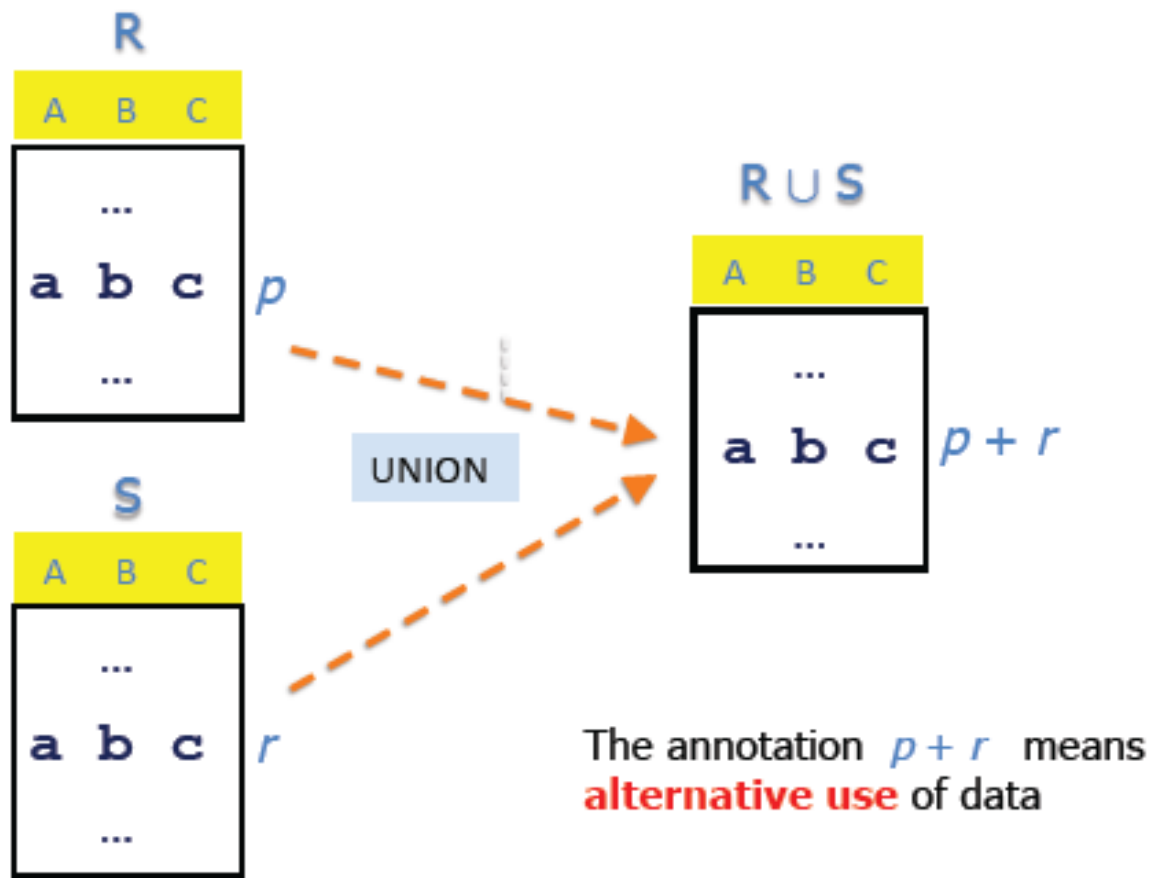
- Operations for annotation propagation in a uniform view and with good properties



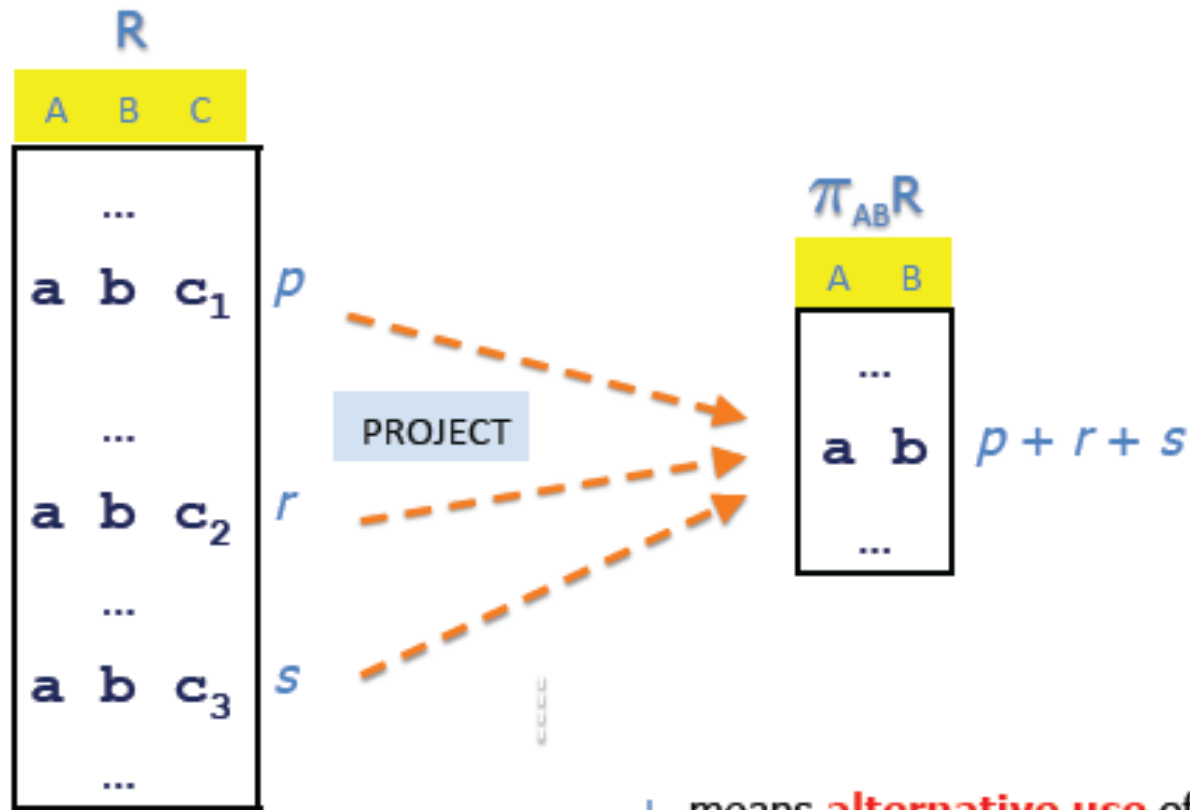
# Propagating annotations through database operations



## Another way to propagate annotations



# Another use of +



+ means **alternative use** of data

## An example in positive relational algebra (SPJU)

$$R \quad Q = \sigma_{C=e} \pi_{AC} ( \pi_{AC} R \bowtie \pi_{BC} R \cup \pi_{AB} R \bowtie \pi_{BC} R )$$

A	B	C	
a	b	c	<i>p</i>
d	b	e	<i>r</i>
f	g	e	<i>s</i>

## An example in positive relational algebra (SPJU)

$$Q = \sigma_{C=e} \pi_{AC} (\pi_{AC} R \bowtie \pi_{BC} R \cup \pi_{AB} R \bowtie \pi_{BC} R)$$

A	B	C
a	b	c
d	b	e
f	g	e

*p*  
*r*  
*s*

A	C
a	c
a	e
d	c
d	e
f	e

$(p \cdot p + p \cdot p) \cdot 0$   
 $p \cdot r \cdot 1$   
 $r \cdot p \cdot 0$   
 $(r \cdot r + r \cdot s + r \cdot r) \cdot 1$   
 $(s \cdot s + s \cdot r + s \cdot s) \cdot 1$

For selection we multiply  
with two special annotations, 0 and 1

# FORMALIZATION

A space of annotations,  $K$

**$K$ -relations**: every tuple annotated with some element from  $K$ .

Binary operations on  $K$ :  $\cdot$  corresponds to joint use (join),  
and  $+$  corresponds to alternative use (union and projection).

We assume  $K$  contains special annotations  $0$  and  $1$ .

“Absent” tuples are annotated with  $0$ !

$1$  is a “neutral” annotation (no restrictions).

**Algebra of annotations?** What are the **laws** of  $(K, +, \cdot, 0, 1)$  ?



# REQUIRED LAWS ON $K$

## ■ Equivalent queries should produce the same annotations

- Union is associative and commutative
- Join is associative, commutative and distributes over union
- Projection and selection commute with each other and with union and join (when applicable)

## ■ Theorem:

- Above identities holds for queries on  $K$ -relations iff  $(K, +, \cdot, 1, 0)$  is a commutative semiring



## What is a commutative semiring?

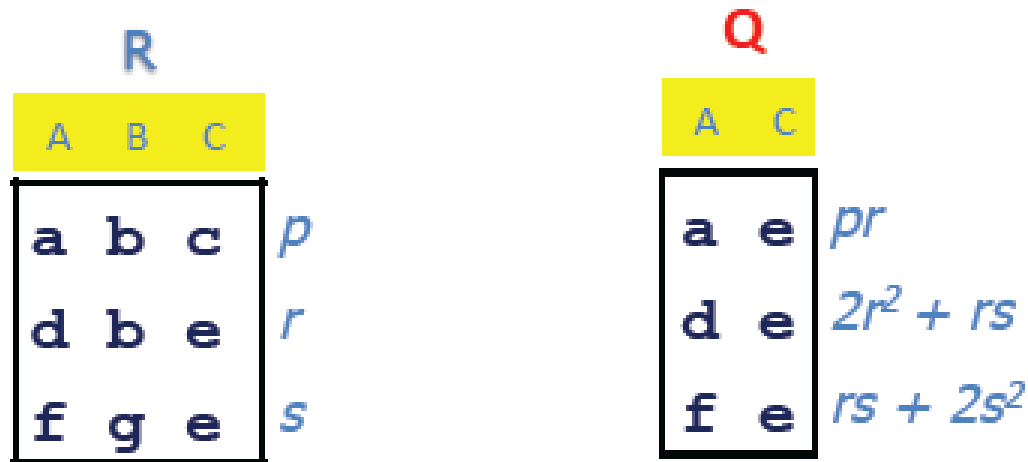
An algebraic structure  $(K, +, \cdot, 0, 1)$  where:

- $K$  is the domain
  - $+$  is associative, commutative, with  $0$  identity
  - $\cdot$  is associative, with  $1$  identity
  - $\cdot$  distributes over  $+$
  - $a \cdot 0 = 0 \cdot a = 0$
- } **semiring**
- $\cdot$  is also **commutative**

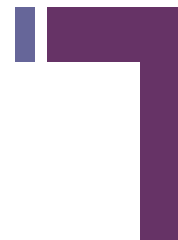
Unlike ring, no requirement for inverses to  $+$



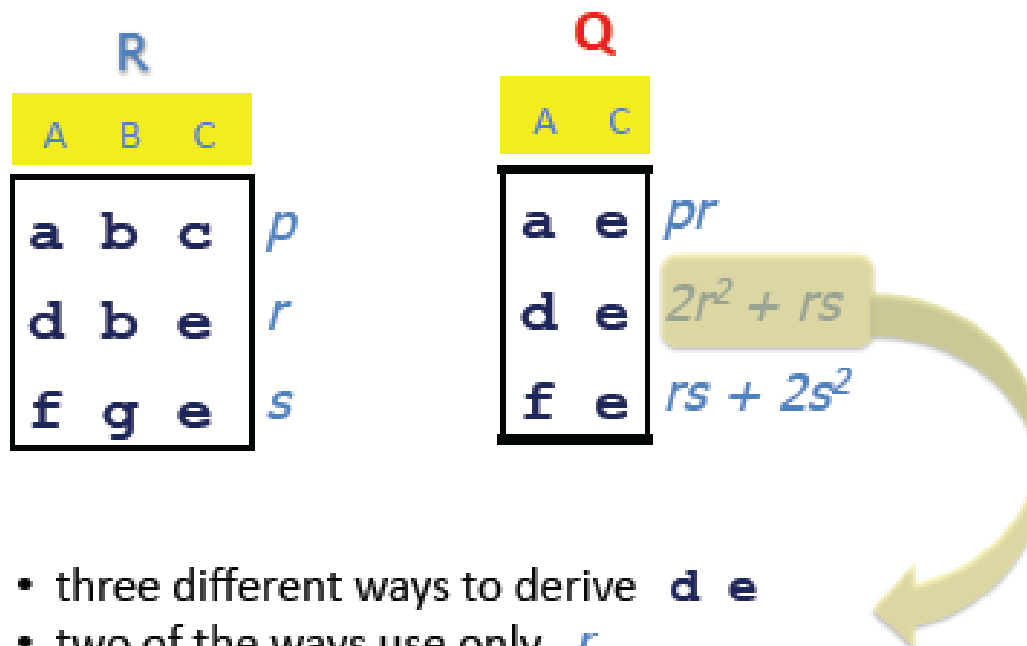
## Using the laws: **polynomials**



Polynomials with coefficients in  $\mathbb{N}$  and **annotation tokens** as indeterminates  $p, r, s$  capture a very general form of **provenance**



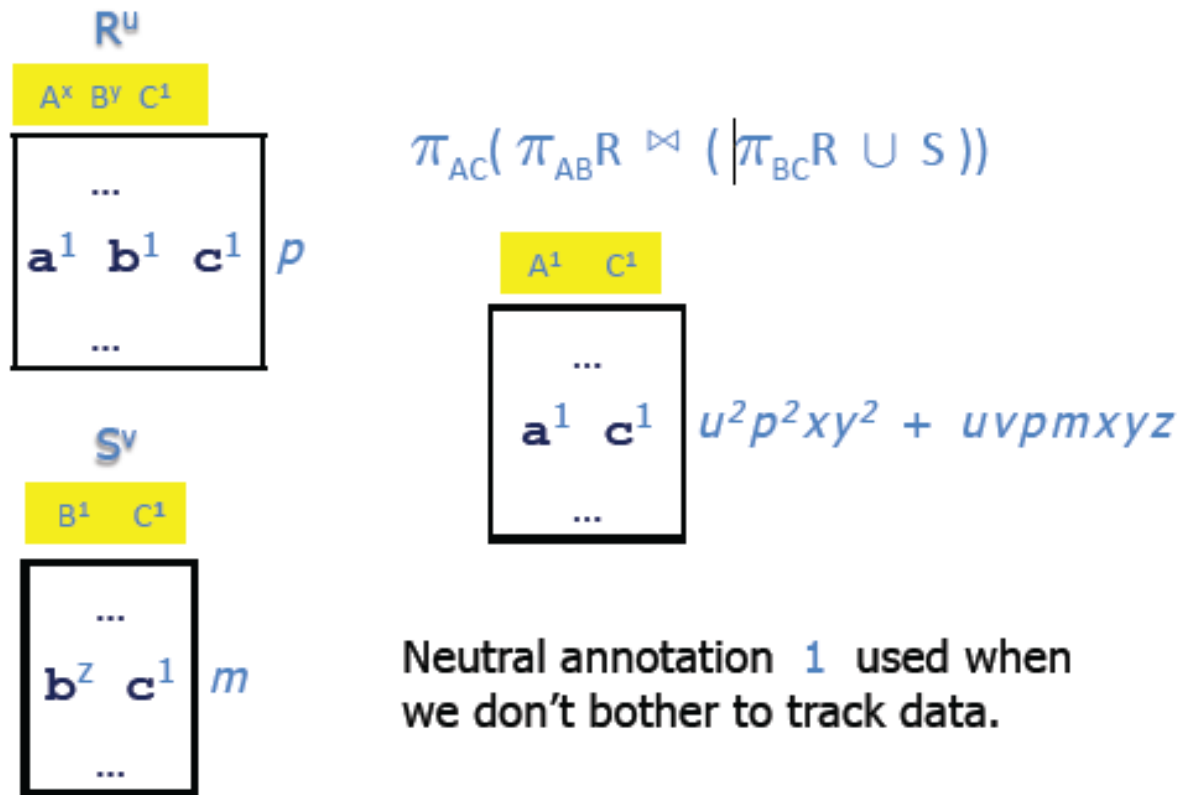
## Provenance reading of the polynomials



- three different ways to derive **d e**
- two of the ways use only *r*
- but they use it twice
- the third way uses *r* once and *s* once

# BEYOND TUPLE ANNOTATIONS

## Relation, attribute and field annotation (1)



# HOMOMORPHISMS

## ■ Useful to evaluate the polynomials

- by mapping  $X$  to  $K$  and extending it to an homomorphism from  $N[X]$  to  $K$  :

- $h(p_1 + p_2) = h(p_1) + h(p_2)$                        $h(p_1 \cdot p_2) = h(p_1) \cdot h(p_2)$
- $h(1) = 1$      $h(0) = 0$

- to obtain:

- trust scores
- multiplicity
- uncertainty values
- access control levels

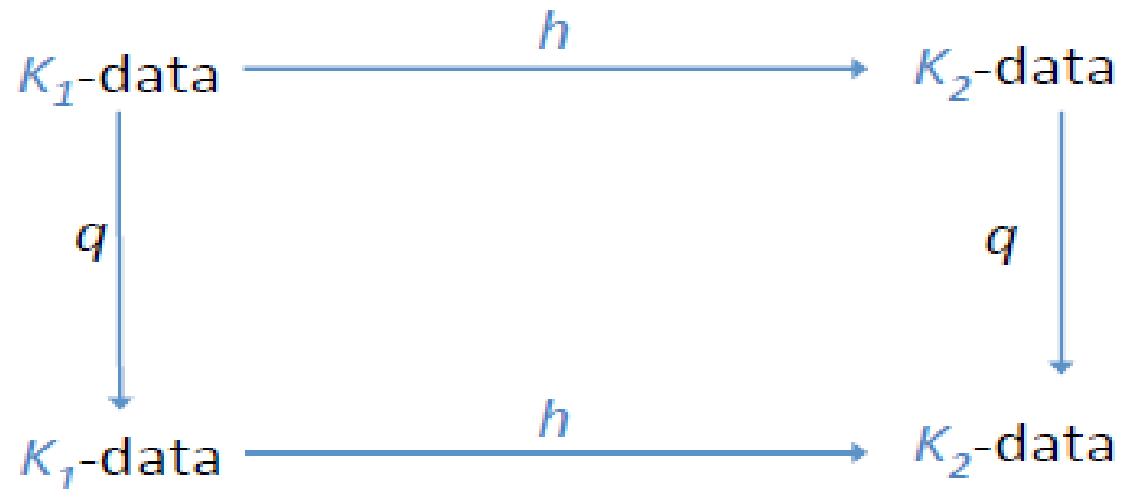
## ■ Useful to hide detail and increase abstraction

- by mapping provenance tokens, many to few
- Stop tracking tokens by mapping them to  $\mathbf{1}$  (neutral)

Doesn't always work, eg. difference.

## Fundamental property

For every query  $q$  and every homomorphism of commutative semirings  $h : K_1 \rightarrow K_2$  the following “commutes”:



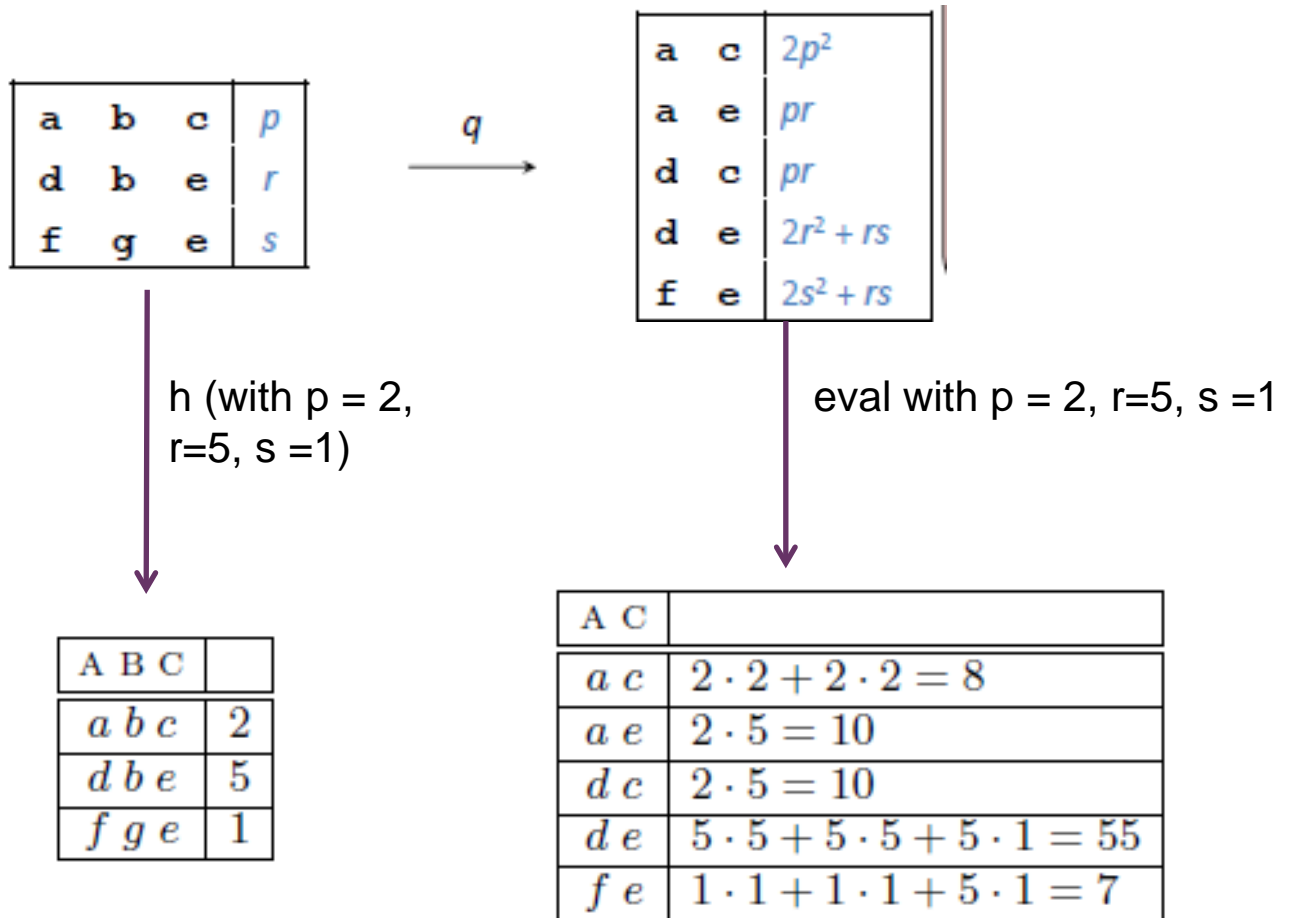
03/24/10

EDBT Keynote, Lausanne

42

$$q(h(R)) = h(q(R))$$

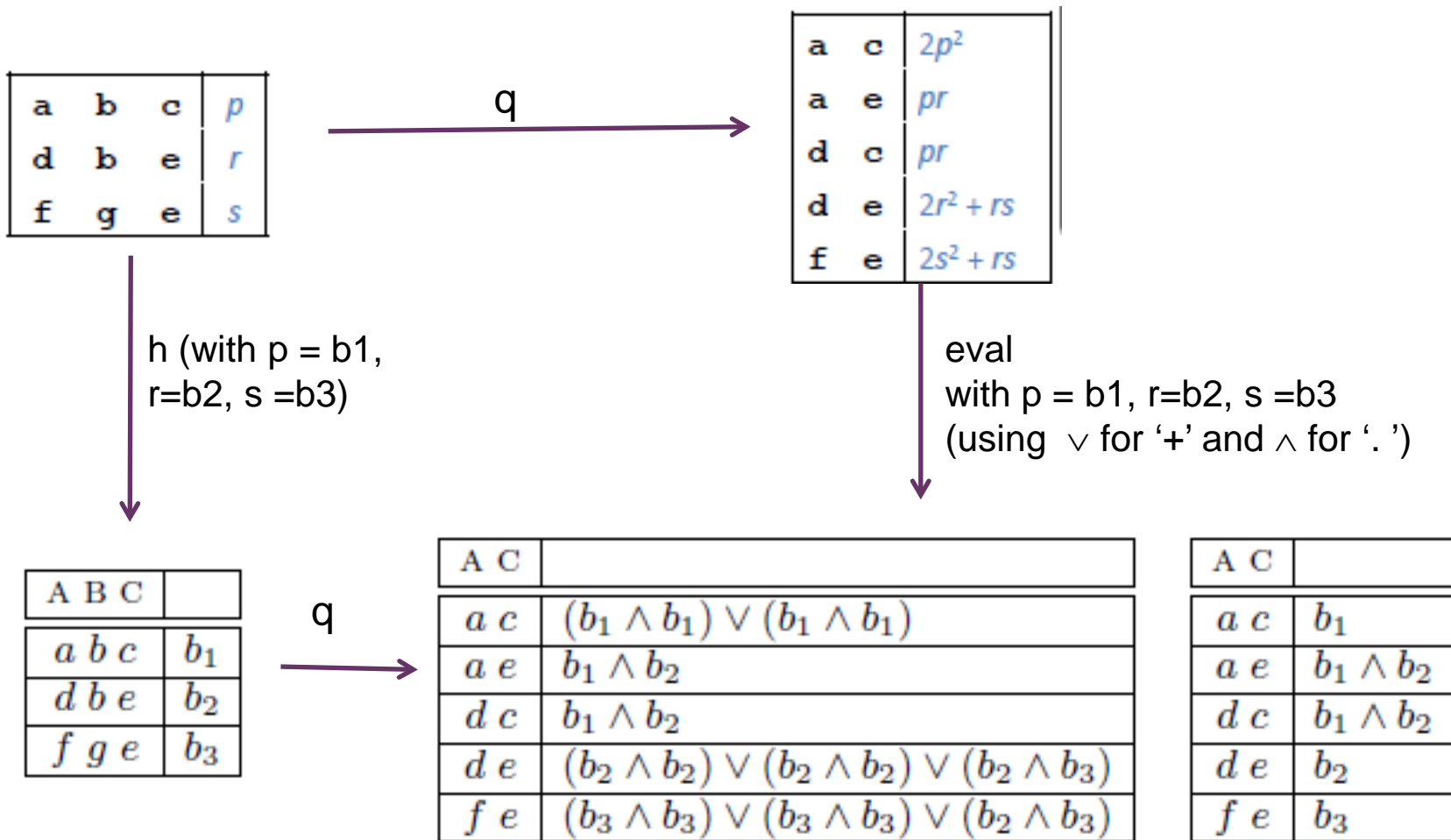
# (Direct) application to multiplicity: bag semantics





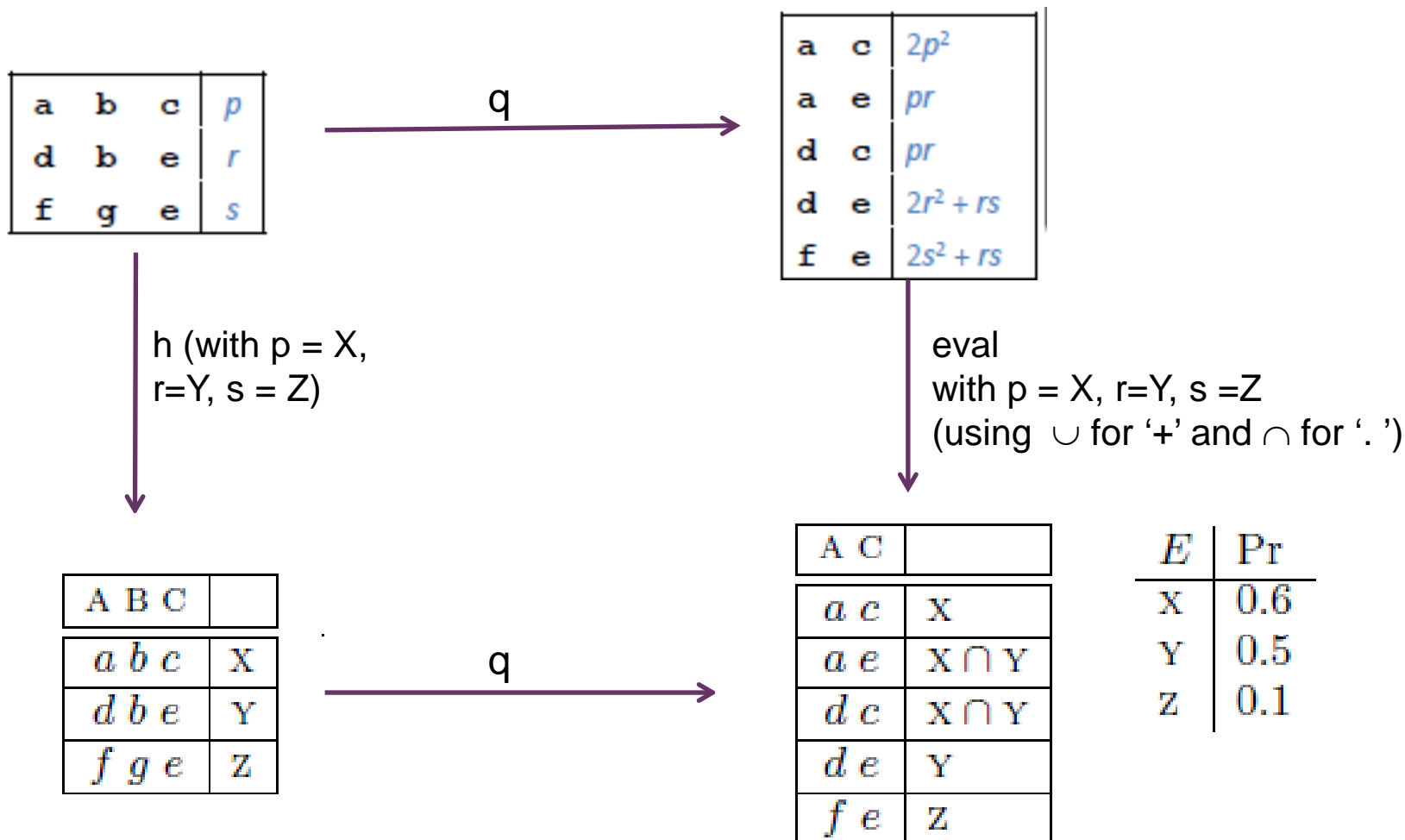
# Application to c-tables: boolean conditions used as annotations for modeling incomplete databases

$(\text{BoolExp}(X), \wedge, \vee, \top, \perp)$



# Application to event tables: probabilistic databases

$$(\mathcal{P}(\Omega), \cup, \cap, \emptyset, \Omega)$$



# Application to access control

$(\mathbb{A}, \min, \max, 0, P)$  where  $\mathbb{A} = P < C < S < T < 0$

Suppose  $p$  is public,  $r$  is secret,  $s$  is top secret

a	b	c	$p$
d	b	e	$r$
f	g	e	$s$

$q$

a	c	$2p^2$
a	e	$pr$
d	c	$pr$
d	e	$2r^2 + rs$
f	e	$2s^2 + rs$

Fundamental property implies that applying the clearance to the database or to the query answer yields the **same result**. (But only the second is actually feasible!)

with  $p = P, r = S, s = T$

↓

a	b	c	P
d	b	e	S
f	g	e	T

$q$

a	c	P
a	e	S
d	c	S
d	e	S
<del>f</del>	<del>e</del>	<del>T</del>

↓ eval with  $p = P, r = S, s = T$   
(using min for "+", max for ".")

"User with secret clearance"

# CONCLUSION



- Provenance is the basis for :
  - UNCERTAINTY
  - TRUST
  - SECURITY
  - MULTIPLICITY
- The semiring  $(\mathbb{N}[X], +, \cdot, 0, 1)$  is the algebraic foundation for computing provenance over queries
- Homomorphisms are the tools for evaluating provenance in different settings

