

# Factorized databases I

Jakub Zavodny (University of Oxford, UK)

# DAMOL

DATA ANALYSIS AND MODELING LAB

Palacky University, Olomouc, Czech Republic



europa  
european  
social fund in the  
czech republic



EUROPEAN UNION



MINISTRY OF EDUCATION,  
YOUTH AND SPORTS



OP Education  
for Competitiveness

INVESTMENTS IN EDUCATION DEVELOPMENT

# **Factorised Databases 1.**

## **Factorisation of Relations**

Jakub Závodný, University of Oxford

based on joint work with PhD supervisor Dan Olteanu

# Relational Databases

Data is arranged in a table.

<u>chain</u>	<u>country</u>
Tesco	Czech Republic
Tesco	Slovakia
Tesco	Hungary
Tesco	Poland
Tesco	UK
Tesco	Spain
Lidl	Czech Republic
Lidl	Slovakia
Lidl	Hungary
Lidl	Poland
Lidl	UK
Aldi	Hungary
Aldi	Poland
Aldi	UK
Aldi	Spain

# Relational Databases

A table is modelled by a relation  $HasShop \subseteq Chains \times Countries$

<u>chain , country</u>
{(Tesco , Czech Republic),
(Tesco , Slovakia),
(Tesco , Hungary),
(Tesco , Poland),
(Tesco , UK),
(Tesco , Spain),
(Lidl , Czech Republic),
(Lidl , Slovakia),
(Lidl , Hungary),
(Lidl , Poland),
(Lidl , UK),
(Aldi , Hungary),
(Aldi , Poland),
(Aldi , UK),
(Aldi , Spain)}

Relation = a set of tuples

# Relational Databases

Relation = a set of tuples = union of products of singletons.

$$\begin{aligned} \text{HasShop} = & \langle \text{Tesco} \rangle \times \langle \text{CzechRep} \rangle \cup \\ & \langle \text{Tesco} \rangle \times \langle \text{Slovakia} \rangle \cup \\ & \langle \text{Tesco} \rangle \times \langle \text{Hungary} \rangle \cup \\ & \langle \text{Tesco} \rangle \times \langle \text{Poland} \rangle \cup \\ & \langle \text{Tesco} \rangle \times \langle \text{UK} \rangle \cup \\ & \langle \text{Tesco} \rangle \times \langle \text{Spain} \rangle \cup \\ & \langle \text{Lidl} \rangle \times \langle \text{CzechRep} \rangle \cup \\ & \langle \text{Lidl} \rangle \times \langle \text{Slovakia} \rangle \cup \\ & \langle \text{Lidl} \rangle \times \langle \text{Hungary} \rangle \cup \\ & \langle \text{Lidl} \rangle \times \langle \text{Poland} \rangle \cup \\ & \langle \text{Lidl} \rangle \times \langle \text{UK} \rangle \cup \\ & \langle \text{Aldi} \rangle \times \langle \text{Hungary} \rangle \cup \\ & \langle \text{Aldi} \rangle \times \langle \text{Poland} \rangle \cup \\ & \langle \text{Aldi} \rangle \times \langle \text{UK} \rangle \cup \\ & \langle \text{Aldi} \rangle \times \langle \text{Spain} \rangle \end{aligned}$$

$\langle \text{value} \rangle$  denotes a singleton relation  $\{(\text{value})\}$

# Factorisation in Relational Databases

A union of products can be factorised using distributivity:

$$\begin{aligned} \text{HasShops} = & (\langle \text{Tesco} \rangle \cup \langle \text{Lidl} \rangle) \times (\langle \text{CzechRep} \rangle \cup \langle \text{Slovakia} \rangle) \cup \\ & (\langle \text{Tesco} \rangle \cup \langle \text{Lidl} \rangle \cup \langle \text{Aldi} \rangle) \times (\langle \text{Hungary} \rangle \cup \langle \text{Poland} \rangle \cup \langle \text{UK} \rangle) \cup \\ & (\langle \text{Tesco} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{Spain} \rangle \end{aligned}$$

Factorisation = expression using singletons,  $\cup$  and  $\times$ .

# Factorisation in Relational Databases

A union of products can be factorised using distributivity:

$$\begin{aligned} \text{HasShops} &= (\langle \text{Tesco} \rangle \cup \langle \text{Lidl} \rangle) \times (\langle \text{CzechRep} \rangle \cup \langle \text{Slovakia} \rangle) \cup \\ &\quad (\langle \text{Tesco} \rangle \cup \langle \text{Lidl} \rangle \cup \langle \text{Aldi} \rangle) \times (\langle \text{Hungary} \rangle \cup \langle \text{Poland} \rangle \cup \langle \text{UK} \rangle) \cup \\ &\quad (\langle \text{Tesco} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{Spain} \rangle \\ \\ &= (\langle \text{Tesco} \rangle \cup \langle \text{Lidl} \rangle) \times (\langle \text{CzechRep} \rangle \cup \langle \text{Slovakia} \rangle \cup \langle \text{Hungary} \rangle \cup \langle \text{Poland} \rangle \cup \langle \text{UK} \rangle) \cup \\ &\quad \langle \text{Tesco} \rangle \times \langle \text{Spain} \rangle \cup \\ &\quad \langle \text{Aldi} \rangle \times (\langle \text{Hungary} \rangle \cup \langle \text{Poland} \rangle \cup \langle \text{UK} \rangle \cup \langle \text{Spain} \rangle) \end{aligned}$$

A given relation can have many factorisations, algebraically equivalent by commutativity of  $\cup$  and  $\times$  and by distributivity of  $\times$  over  $\cup$ .

# Connection: Formal Concepts

<u>chain</u>	<u>country</u>
Tesco	Czech Republic
Tesco	Slovakia
Tesco	Hungary
Tesco	Poland
Tesco	UK
Tesco	Spain
Lidl	Czech Republic
Lidl	Slovakia
Lidl	Hungary
Lidl	Poland
Lidl	UK
Aldi	Hungary
Aldi	Poland
Aldi	UK
Aldi	Spain

	CzechRep	Slovakia	Hungary	Poland	Spain	UK
Tesco	✓	✓	✓	✓	✓	✓
Lidl	✓	✓	✓	✓		✓
Aldi			✓	✓	✓	✓



# Connection: Formal Concepts

$$(\langle \text{Tesco} \rangle \cup \langle \text{Lidl} \rangle) \times ((\langle \text{CzechRep} \rangle \cup \langle \text{Slovakia} \rangle) \cup \quad (1)$$

$$(\langle \text{Tesco} \rangle \cup \langle \text{Lidl} \rangle \cup \langle \text{Aldi} \rangle) \times ((\langle \text{Hungary} \rangle \cup \langle \text{Poland} \rangle \cup \langle \text{UK} \rangle) \cup \quad (2)$$

$$(\langle \text{Tesco} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{Spain} \rangle \quad (3)$$

	CzechRep	Slovakia	Hungary	Poland	Spain	UK
Tesco	(1)	(1)	(2)	(2)	(3)	(2)
Lidl	(1)	(1)	(2)	(2)		(2)
Aldi			(2)	(2)	(3)	(2)

- Formal concept, closed  $n$ -set = maximal product of unions

But:

- Factorisations allow non-maximal products of unions.
- Factorisations for  $n$ -ary relations allow deeper nesting.

# Factorisation in Relational Databases

More columns, more complicated factorisations.

$$R_2 \subseteq \text{Chains} \times \text{Products} \times \text{Countries}$$

$$\begin{aligned} & (\langle \text{Tesco} \rangle \times (\langle \text{food} \rangle \cup \langle \text{clothes} \rangle) \cup \langle \text{Lidl} \rangle \times \langle \text{food} \rangle) \times ((\langle \text{CzechRep} \rangle \cup \langle \text{Slovakia} \rangle) \cup \\ & (\langle \text{Tesco} \rangle \cup \langle \text{Lidl} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{food} \rangle \times ((\langle \text{Poland} \rangle \cup \langle \text{UK} \rangle) \cup \\ & (\langle \text{Tesco} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{books} \rangle \times ((\langle \text{UK} \rangle \cup \langle \text{Spain} \rangle)) \end{aligned}$$

chain	product	country
Tesco	food	CzechRep
Tesco	clothes	CzechRep
Lidl	food	CzechRep
Tesco	food	Slovakia
Tesco	clothes	Slovakia
Lidl	food	Slovakia
Tesco	food	Poland

...

# Why Factorisation?

$$\begin{aligned} & (\langle \text{Tesco} \rangle \times (\langle \text{food} \rangle \cup \langle \text{clothes} \rangle) \cup \langle \text{Lidl} \rangle \times \langle \text{food} \rangle) \times (\langle \text{CzechRep} \rangle \cup \langle \text{Slovakia} \rangle) \cup \\ & (\langle \text{Tesco} \rangle \cup \langle \text{Lidl} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{food} \rangle \times (\langle \text{Poland} \rangle \cup \langle \text{UK} \rangle) \cup \\ & (\langle \text{Tesco} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{books} \rangle \times (\langle \text{UK} \rangle \cup \langle \text{Spain} \rangle) \end{aligned}$$

- Data Analysis people say:  
Factorisations reveal deeper structure of the data!
- Database people say:  
Factorisations are useful for data storage!

# Why Factorisation?

$$\begin{aligned} & (\langle \text{Tesco} \rangle \times (\langle \text{food} \rangle \cup \langle \text{clothes} \rangle) \cup \langle \text{Lidl} \rangle \times \langle \text{food} \rangle) \times ((\langle \text{CzechRep} \rangle \cup \langle \text{Slovakia} \rangle) \cup \\ & (\langle \text{Tesco} \rangle \cup \langle \text{Lidl} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{food} \rangle \times ((\langle \text{Poland} \rangle \cup \langle \text{UK} \rangle) \cup \\ & (\langle \text{Tesco} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{books} \rangle \times ((\langle \text{UK} \rangle \cup \langle \text{Spain} \rangle)) \end{aligned}$$

- Factorisations compactly encode combinations of groups of values.

- ▶ Can be exponentially smaller than flat relations.

$$(\langle a \rangle \cup \langle b \rangle) \times (\langle a \rangle \cup \langle b \rangle) \times \cdots \times (\langle a \rangle \cup \langle b \rangle)$$

- ▶ Results of conjunctive queries have succinct factorisations!

# Why Factorisation?

$$\begin{aligned} & \left( \langle \overline{\text{Tesco}} \rangle \times (\langle \overline{\text{food}} \rangle \cup \langle \text{clothes} \rangle) \cup \langle \text{Lidl} \rangle \times \langle \text{food} \rangle \right) \times (\langle \overline{\text{CzechRep}} \rangle \cup \langle \text{Slovakia} \rangle) \cup \\ & \quad (\langle \text{Tesco} \rangle \cup \langle \text{Lidl} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{food} \rangle \times (\langle \text{Poland} \rangle \cup \langle \text{UK} \rangle) \cup \\ & \quad (\langle \text{Tesco} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{books} \rangle \times (\langle \text{UK} \rangle \cup \langle \text{Spain} \rangle) \end{aligned}$$

- Factorisations compactly encode combinations of groups of values.
- Factorisations allow for enumerating tuples with constant delay.

Tesco food CzechRep

# Why Factorisation?

$$\begin{aligned} & \left( \langle \overline{\text{Tesco}} \rangle \times (\langle \overline{\text{food}} \rangle \cup \langle \text{clothes} \rangle) \cup \langle \text{Lidl} \rangle \times \langle \text{food} \rangle \right) \times (\langle \text{CzechRep} \rangle \cup \langle \overline{\text{Slovakia}} \rangle) \cup \\ & \quad (\langle \text{Tesco} \rangle \cup \langle \text{Lidl} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{food} \rangle \times (\langle \text{Poland} \rangle \cup \langle \text{UK} \rangle) \cup \\ & \quad (\langle \text{Tesco} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{books} \rangle \times (\langle \text{UK} \rangle \cup \langle \text{Spain} \rangle) \end{aligned}$$

- Factorisations compactly encode combinations of groups of values.
- Factorisations allow for enumerating tuples with constant delay.

Tesco food CzechRep  
Tesco food Slovakia

# Why Factorisation?

$$\begin{aligned} & (\langle \overline{\text{Tesco}} \rangle \times (\langle \text{food} \rangle \cup \langle \overline{\text{clothes}} \rangle) \cup \langle \text{Lidl} \rangle \times \langle \text{food} \rangle) \times (\langle \overline{\text{CzechRep}} \rangle \cup \langle \text{Slovakia} \rangle) \cup \\ & (\langle \text{Tesco} \rangle \cup \langle \text{Lidl} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{food} \rangle \times (\langle \text{Poland} \rangle \cup \langle \text{UK} \rangle) \cup \\ & (\langle \text{Tesco} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{books} \rangle \times (\langle \text{UK} \rangle \cup \langle \text{Spain} \rangle) \end{aligned}$$

- Factorisations compactly encode combinations of groups of values.
- Factorisations allow for enumerating tuples with constant delay.

Tesco food CzechRep  
Tesco food Slovakia  
Tesco clothes CzechRep

# Why Factorisation?

$$\begin{aligned} & (\langle \overline{\text{Tesco}} \rangle \times (\langle \text{food} \rangle \cup \langle \overline{\text{clothes}} \rangle) \cup \langle \text{Lidl} \rangle \times \langle \text{food} \rangle) \times ((\langle \text{CzechRep} \rangle \cup \langle \overline{\text{Slovakia}} \rangle) \cup \\ & (\langle \text{Tesco} \rangle \cup \langle \text{Lidl} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{food} \rangle \times (\langle \text{Poland} \rangle \cup \langle \text{UK} \rangle) \cup \\ & (\langle \text{Tesco} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{books} \rangle \times (\langle \text{UK} \rangle \cup \langle \text{Spain} \rangle)) \end{aligned}$$

- Factorisations compactly encode combinations of groups of values.
- Factorisations allow for enumerating tuples with constant delay.

Tesco food CzechRep  
Tesco food Slovakia  
Tesco clothes CzechRep  
Tesco clothes Slovakia



# Why Factorisation?

$$\begin{aligned} & (\langle \text{Tesco} \rangle \times (\langle \text{food} \rangle \cup \langle \text{clothes} \rangle) \cup \langle \text{Lidl} \rangle \times \langle \text{food} \rangle) \times (\langle \text{CzechRep} \rangle \cup \langle \text{Slovakia} \rangle) \cup \\ & (\langle \text{Tesco} \rangle \cup \langle \text{Lidl} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{food} \rangle \times (\langle \text{Poland} \rangle \cup \langle \text{UK} \rangle) \cup \\ & (\langle \text{Tesco} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{books} \rangle \times (\langle \text{UK} \rangle \cup \langle \text{Spain} \rangle) \end{aligned}$$

- Factorisations compactly encode combinations of groups of values.
- Factorisations allow for enumerating tuples with constant delay.

Tesco food CzechRep  
Tesco food Slovakia  
Tesco clothes CzechRep  
Tesco clothes Slovakia  
Lidl food CzechRep

# Why Factorisation?

$$\begin{aligned} & (\langle \text{Tesco} \rangle \times (\langle \text{food} \rangle \cup \langle \text{clothes} \rangle) \cup \langle \overline{\text{Lidl}} \rangle \times \langle \overline{\text{food}} \rangle) \times ((\langle \text{CzechRep} \rangle \cup \langle \overline{\text{Slovakia}} \rangle) \cup \\ & (\langle \text{Tesco} \rangle \cup \langle \text{Lidl} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{food} \rangle \times (\langle \text{Poland} \rangle \cup \langle \text{UK} \rangle) \cup \\ & (\langle \text{Tesco} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{books} \rangle \times (\langle \text{UK} \rangle \cup \langle \text{Spain} \rangle)) \end{aligned}$$

- Factorisations compactly encode combinations of groups of values.
- Factorisations allow for enumerating tuples with constant delay.

Tesco food CzechRep  
Tesco food Slovakia  
Tesco clothes CzechRep  
Tesco clothes Slovakia  
Lidl food CzechRep  
Lidl food Slovakia

# Why Factorisation?

$$\begin{aligned} & (\langle \text{Tesco} \rangle \times (\langle \text{food} \rangle \cup \langle \text{clothes} \rangle) \cup \langle \text{Lidl} \rangle \times \langle \text{food} \rangle) \times ((\langle \text{CzechRep} \rangle \cup \langle \text{Slovakia} \rangle) \cup \\ & (\langle \overline{\text{Tesco}} \rangle \cup \langle \text{Lidl} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \overline{\text{food}} \rangle \times (\langle \overline{\text{Poland}} \rangle \cup \langle \text{UK} \rangle) \cup \\ & (\langle \text{Tesco} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{books} \rangle \times (\langle \text{UK} \rangle \cup \langle \text{Spain} \rangle)) \end{aligned}$$

- Factorisations compactly encode combinations of groups of values.
- Factorisations allow for enumerating tuples with constant delay.

Tesco	food	CzechRep
Tesco	food	Slovakia
Tesco	clothes	CzechRep
Tesco	clothes	Slovakia
Lidl	food	CzechRep
Lidl	food	Slovakia
Tesco	food	Poland

...

# Why Factorisation?

$$\begin{aligned} & (\langle \text{Tesco} \rangle \times (\langle \text{food} \rangle \cup \langle \text{clothes} \rangle) \cup \langle \text{Lidl} \rangle \times \langle \text{food} \rangle) \times ((\langle \text{CzechRep} \rangle \cup \langle \text{Slovakia} \rangle) \cup \\ & (\langle \text{Tesco} \rangle \cup \langle \text{Lidl} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{food} \rangle \times ((\langle \text{Poland} \rangle \cup \langle \text{UK} \rangle) \cup \\ & (\langle \text{Tesco} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{books} \rangle \times ((\langle \text{UK} \rangle \cup \langle \text{Spain} \rangle)) \end{aligned}$$

- Factorisations compactly encode combinations of groups of values.
- Factorisations allow for enumerating tuples with constant delay.
  - ▶ Unlike general join decompositions or trivial representation  $(Q, D)$  of query results.

# Why Factorisation?

$$\begin{aligned} & (\langle \text{Tesco} \rangle \times (\langle \text{food} \rangle \cup \langle \text{clothes} \rangle) \cup \langle \text{Lidl} \rangle \times \langle \text{food} \rangle) \times ((\langle \text{CzechRep} \rangle \cup \langle \text{Slovakia} \rangle) \cup \\ & (\langle \text{Tesco} \rangle \cup \langle \text{Lidl} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{food} \rangle \times ((\langle \text{Poland} \rangle \cup \langle \text{UK} \rangle) \cup \\ & (\langle \text{Tesco} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{books} \rangle \times ((\langle \text{UK} \rangle \cup \langle \text{Spain} \rangle)) \end{aligned}$$

- Factorisations compactly encode combinations of groups of values.
- Factorisations allow for enumerating tuples with constant delay.
  - ▶ Unlike general join decompositions or trivial representation  $(Q, D)$  of query results.
- Factorisations allow query evaluation *without expanding* to flat relations
  - ▶ Using factorisations can boost query evaluation (processing a smaller representation).
  - ▶ Unlike general compression methods.

# Why Factorisation?

$$\begin{aligned} & (\langle \text{Tesco} \rangle \times (\langle \text{food} \rangle \cup \langle \text{clothes} \rangle) \cup \langle \text{Lidl} \rangle \times \langle \text{food} \rangle) \times (\langle \text{CzechRep} \rangle \cup \langle \text{Slovakia} \rangle) \cup \\ & (\langle \text{Tesco} \rangle \cup \langle \text{Lidl} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{food} \rangle \times (\langle \text{Poland} \rangle \cup \langle \text{UK} \rangle) \cup \\ & (\langle \text{Tesco} \rangle \cup \langle \text{Aldi} \rangle) \times \langle \text{books} \rangle \times (\langle \text{UK} \rangle \cup \langle \text{Spain} \rangle) \end{aligned}$$

- Factorisations compactly encode combinations of groups of values.
- Factorisations allow for enumerating tuples with constant delay.
  - ▶ Unlike general join decompositions or trivial representation  $(Q, D)$  of query results.
- Factorisations allow query evaluation *without expanding* to flat relations
  - ▶ Using factorisations can boost query evaluation (processing a smaller representation).
  - ▶ Unlike general compression methods.

# Applications of Factorised Representations

Relational data storage and management:

- Storing & querying relationship (many-to-many) data
- Succinct representation of large query results
  - ▶ Storing intermediate results in query evaluation
  - ▶ Knowledge compilation (compile once, speed up all subsequent queries)
  - ▶ Parallel computation (compress intermediate results sent b/w machines)
- Managing a large set of possibilities or choices
  - ▶ Configuration problems (space of valid solutions)
  - ▶ Incomplete information (space of possible worlds)
- Provenance and probabilistic data
  - ▶ Compact encoding for large provenance
  - ▶ Factorisation of provenance is used for query evaluation in probabilistic databases.

# Key Challenges

1. How compact can factorised query results be?

2. Can such factorisations speed up query evaluation?



# Key Results

## 1. How compact can factorised query results be?

- Asymptotic size bounds for factorisations of query results.
- Characterisation of conjunctive queries based on factorisability of their results.

## 2. Can such factorisations speed up query evaluation?

- Compute factorisations directly from query and input data.
- Optimisation and evaluation techniques for queries with selections, projections, joins, aggregates and ordering on factorisations.
- FDB: a main-memory engine using factorised relations.

# Key Challenges

1. How compact can factorised query results be?

# Factorisation of Joins

Products		Locations	
chain	product	chain	country
Tesco	food	Tesco	Slovakia
Tesco	clothes	Tesco	UK
Tesco	books	Aldi	UK
Aldi	food	Aldi	Spain

Join two tables on *chain*:

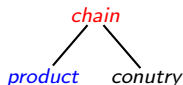
Products		$\bowtie_{chain}$	Locations
shop	product		country
Tesco	food		Slovakia
Tesco	clothes		Slovakia
Tesco	books		Slovakia
Tesco	food		UK
Tesco	clothes		UK
Tesco	books		UK
Aldi	food		UK
Aldi	food		Spain

$$= \langle \text{Tesco} \rangle \times (\langle \text{food} \rangle \cup \langle \text{clothes} \rangle \cup \langle \text{books} \rangle) \times (\langle \text{Slovakia} \rangle \cup \langle \text{UK} \rangle) \cup \\ \langle \text{Aldi} \rangle \times \langle \text{food} \rangle \times (\langle \text{UK} \rangle \cup \langle \text{Spain} \rangle)$$

# Factorisation Trees

Uniform nesting structures  $\longleftrightarrow$  defined by *factorisation trees* (f-trees).

$$\bigcup_{\text{Chains}} \left( \langle \text{chain} \rangle \times \left( \bigcup_{\text{Products}} \langle \text{product} \rangle \right) \times \left( \bigcup_{\text{Countries}} \langle \text{country} \rangle \right) \right) \longleftrightarrow$$

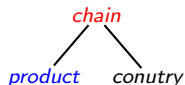


$$\begin{aligned} & \langle \text{Tesco} \rangle \times (\langle \text{food} \rangle \cup \langle \text{clothes} \rangle \cup \langle \text{books} \rangle) \times (\langle \text{Slovakia} \rangle \cup \langle \text{UK} \rangle) \\ \cup & \langle \text{Aldi} \rangle \times \langle \text{food} \rangle \times (\langle \text{UK} \rangle \cup \langle \text{Spain} \rangle) \end{aligned}$$

# Factorisation Trees

Uniform nesting structures  $\longleftrightarrow$  defined by *factorisation trees* (f-trees).

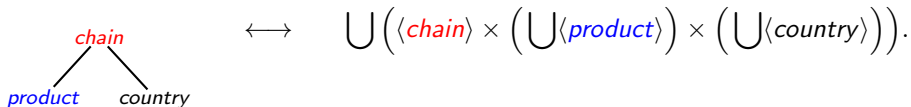
$$\bigcup_{\text{Chains}} \left( \langle \text{chain} \rangle \times \left( \bigcup_{\text{Products}} \langle \text{product} \rangle \right) \times \left( \bigcup_{\text{Countries}} \langle \text{country} \rangle \right) \right) \longleftrightarrow$$



$$\bigcup_{\text{Chains}} \left( \langle \text{chain} \rangle \times \left( \bigcup_{\text{Products}} \langle \text{product} \rangle \times \left( \bigcup_{\text{Countries}} \langle \text{country} \rangle \right) \right) \right) \longleftrightarrow$$



# Factorisation Trees for Relations



Some relations are not factorisable using this f-tree

ch	p	c		
Tesco	food	UK	$\neq$	$\bigcup \left( \langle \text{ch} \rangle \times \left( \bigcup \langle \text{p} \rangle \right) \times \left( \bigcup \langle \text{c} \rangle \right) \right)$
Tesco	clothes	CzechRep		

but the join is *always* factorisable using this f-tree

$$\text{Products}(\text{ch}, p) \bowtie_s \text{Locations}(\text{ch}, c) = \bigcup \left( \langle \text{ch} \rangle \times \left( \bigcup \langle \text{p} \rangle \right) \times \left( \bigcup \langle \text{c} \rangle \right) \right)$$

since  $p$  and  $c$  are *conditionally independent*.

# Factorisation Trees for Query Results

For any join query  $Q$ ,

we characterise **f-trees that always factorise** the result of  $Q$ .

If  $Q$  is an equi-join query and  $\mathcal{T}$  any f-tree, then

the result  $Q(\mathbf{D})$  can be factorised according to  $\mathcal{T}$  for *any* database  $\mathbf{D}$

**iff**

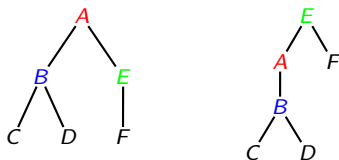
for each relation of  $Q$ , all its attributes are on a root-to-leaf path.

If  $Q$  has projections, a similar but more complicated result holds.

# Factorisation Trees for Query Results

Consider the query:

$$Q(A, B, C, D, E, F) \leftarrow R(A, B, C), S(A, B, D), T(A, E), U(E, F).$$



Left f-tree induces the factorisation structure:

$$\bigcup_{a \in A} \left( \langle a \rangle \times \bigcup_{b \in B} \left( \langle b \rangle \times \left( \bigcup_{c \in C} \langle c \rangle \right) \times \left( \bigcup_{d \in D} \langle d \rangle \right) \right) \times \bigcup_{e \in E} \left( \langle e \rangle \times \left( \bigcup_{f \in F} \langle f \rangle \right) \right) \right)$$



# Key Challenges

1. How compact can factorised query results be?

## Size of Factorisations

The *size* of a factorisation is the number of its singleton data elements.

$$|(\langle a \rangle \langle 1 \rangle \cup \langle a \rangle \langle 2 \rangle \cup \langle b \rangle \langle 1 \rangle \cup \langle b \rangle \langle 2 \rangle \cup \langle c \rangle \langle 1 \rangle \cup \langle c \rangle \langle 2 \rangle)| = 12,$$

$$|(\langle a \rangle \cup \langle b \rangle \cup \langle c \rangle) \times (\langle 1 \rangle \cup \langle 2 \rangle)| = 5.$$

**How much space do we save by factorisation?**

## Size of Factorisations: Main Result

For any conjunctive query  $Q$  there is a number  $s(Q)$  such that

For any database  $\mathbf{D}$ ,  $Q(\mathbf{D})$  admits a factorisation of size  $O(|\mathbf{D}|^{s(Q)})$ .

- **Best possible bound** for a general result of  $Q$  using f-trees.

There exists  $\mathbf{D}$  such that all factorisations over f-trees are  $\Omega(|\mathbf{D}|^{s(Q)})$ .

- Instance-optimal factorisation of relations (i.e., dependent on  $\mathbf{D}$ ) is hard.

## Flat vs. Factorised Size

- For any database  $\mathbf{D}$ ,  $|Q(\mathbf{D})|$  is  $O(|\mathbf{D}|^{\rho^*(Q)})$ . [AGM'08]
- For any database  $\mathbf{D}$ ,  $Q(\mathbf{D})$  admits factorisation of size  $O(|\mathbf{D}|^{s(Q)})$ . [OZ'11]

$$1 \leq s(Q) \leq \rho^*(Q) \leq |Q|$$

There are classes of queries with  $s(Q) \ll \rho^*(Q)$ .

## Appendix: Intuition for Asymptotic Bounds

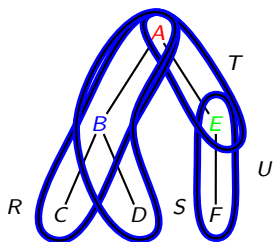
- For any database  $\mathbf{D}$ ,  $|Q(\mathbf{D})|$  is  $O(|\mathbf{D}|^{\rho^*(Q)})$ . [AGM'08]
- For any database  $\mathbf{D}$ ,  $Q(\mathbf{D})$  admits factorisation of size  $O(|\mathbf{D}|^{s(Q)})$ . [OZ'11]

What is  $\rho^*(Q)$  and  $s(Q)$ ?

## Intuition – Edge Covers

$$Q(A, B, C, D, E, F) \leftarrow R(A, B, C), S(A, B, D), T(A, E), U(E, F).$$

The hypergraph of  $Q$ :



First observation:

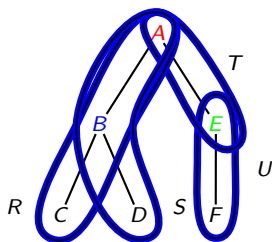
- Cover all attributes by  $k$  relations  $\Rightarrow |Q(\mathbf{D})| \leq |\mathbf{D}|^k$ .
- Set of  $m$  independent attributes  $\Rightarrow$  construct  $\mathbf{D}$  with  $|Q(\mathbf{D})| \sim |\mathbf{D}|^m$ .

$$\max_m = \text{IndependentSet}(Q) \leq \text{EdgeCover}(Q) = \min_k$$

## Intuition – Fractional Edge Covers

$$Q(A, B, C, D, E, F) \leftarrow R(A, B, C), S(A, B, D), T(A, E), U(E, F).$$

The hypergraph of  $Q$ :



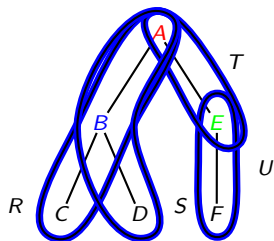
[AGM'08]:

- *Fractional* edge cover of  $Q$  with weight  $k \Rightarrow |Q(\mathbf{D})| \leq |\mathbf{D}|^k$ .
- *Fractional* independent set of weight  $m \Rightarrow$  construct  $\mathbf{D}$  with  $|Q(\mathbf{D})| \sim |\mathbf{D}|^m$ .

By linear programming duality:

$$\max_m = \text{FractionalIndependentSet}(Q) = \text{FractionalEdgeCover}(Q) = \min_k$$

## Intuition – Fractional Edge Covers



For a query  $Q = R_1 \bowtie \dots \bowtie R_n$ , the *fractional edge cover number*  $\rho^*(Q)$  is the cost of an optimal solution to the linear program

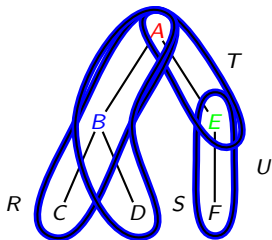
$$\begin{array}{ll} \text{minimising} & \sum_i x_{R_i} \\ \text{subject to} & \sum_{i: R_i \text{ has attribute } A} x_{R_i} \geq 1 \text{ for all attributes } A, \\ & x_{R_i} \geq 0 \text{ for all } R_i. \end{array}$$

- $x_{R_i}$  is the weight of relation  $R_i$ .
- Each attribute has to be covered by relations with sum of weights  $\geq 1$ .
- In the non-weighted edge cover, the variables  $x_{R_i} \in \{0, 1\}$



## Example

$Q(A, B, C, D, E, F) \leftarrow R(A, B, C), S(A, B, D), T(A, E), U(E, F).$

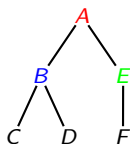


- Relations  $R, S, U$  cover the whole query.  
 $\text{FractionalEdgeCover}(Q) \leq 3$
- Each of the nodes  $C, D,$  and  $F$  must be covered by separate relations.  
 $\text{FractionalIndependentSet}(Q) \geq 3$

$$\Rightarrow \rho^*(Q) = 3$$

$$\Rightarrow |Q(\mathbf{D})| = O(|\mathbf{D}|^3) \text{ and for some inputs } |Q(\mathbf{D})| = \Theta(|\mathbf{D}|^3).$$

## Intuition – Size of Factorisations



$$\bigcup_{a \in A} \left( \langle a \rangle \times \bigcup_{b \in B} \left( \langle b \rangle \times \left( \bigcup_{c \in C} \langle c \rangle \right) \times \left( \bigcup_{d \in D} \langle d \rangle \right) \right) \times \bigcup_{e \in E} \left( \langle e \rangle \times \left( \bigcup_{f \in F} \langle f \rangle \right) \right) \right)$$

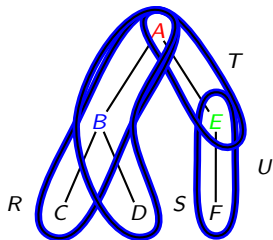
Attributes only depend on ancestor attributes

- $F$  only depends on  $E$  and  $A$ .
- One  $\langle f \rangle$  for each  $(a, e, f) \in Q(\mathbf{D})$ .
- The number of  $F$ -singletons is  $|\pi_{A,E,F}(Q(\mathbf{D}))|$ .

Size of f-representation = sum of sizes of **subqueries along paths**.

## Example

$$Q(A, B, C, D, E, F) \leftarrow R(A, B, C), S(A, B, D), T(A, E), U(E, F).$$



- Path  $A, E, F$  has fractional edge cover 2.  
 $\Rightarrow$  The number of  $F$ -singletons is  $\leq |\mathbf{D}|^2$ , but can be  $\sim |\mathbf{D}|^2$ .
- All other root-to-leaf paths have fractional edge cover 1.  
 $\Rightarrow$  The number of other singletons is  $\leq |\mathbf{D}|$ .

$$s(Q) = 2$$

$$\Rightarrow \text{Factorisation size} \sim |\mathbf{D}|^2$$

$$\text{Recall that } \rho^*(Q) = 3$$

$$\Rightarrow \text{Flat size} \sim |\mathbf{D}|^3$$

## Flat vs. Factorised Size

- For any database  $\mathbf{D}$ ,  $|Q(\mathbf{D})|$  is  $O(|\mathbf{D}|^{\rho^*(Q)})$ . [AGM'08]
- For any database  $\mathbf{D}$ ,  $Q(\mathbf{D})$  admits factorisation of size  $O(|\mathbf{D}|^{s(Q)})$ . [OZ'11]
- $\rho^*(Q)$  = fractional edge cover number of the **entire query**.
- $s(Q)$  = fractional edge cover number of **root-to-leaf paths in best f-tree**.

$$1 \leq s(Q) \leq \rho^*(Q) \leq |Q|$$

There are classes of queries with  $s(Q) \ll \rho^*(Q)$ .

(There are classes of queries with  $s(Q) = 1$  and  $\rho^*(Q) = |Q|$ .)

**Thank you!**