From Hypertree Width to Submodular Width and Data-dependent Structural Decompositions

Francesco Scarcello, DIMES, Università della Calabria and Georg Gottlob, Vienna, Oxford...

Foundational Challenges in Data and Knowledge Management Georg Gottlob special event

# The origins

- I first met Georg Gottlob in Italy, at my first (logic- programming) conference. Georg was an invited speaker
- After a few years I visited Georg and Nicola Leone in Vienna Georg was an invited speaker at a conference on Inductive Logic Programming We worked on PAC learning of <u>acyclic clauses</u>
- We started our investigations on acyclic queries and their generalizations
- It was the time of Georgs's Wittgenstein award
- A beatufil and very stimulating period Kurt Gödel Colloquium and Grohe lecture on the Robber and Cops game
- A never ending friendship and fruitful scientific collaboration
  - Databases, Game Theory, Knowledge Representation and Reasoning

#### The challenge C: a class of conjunctive queries 1): all possible databases A p(d) COMBINED COMPLEXITY C IS AN ISLAND OF TRACTABILITY IF A POLYNOMIAL TIME ALGORITHM A THERE EXISTS 9(b) YgEC and YdED THAT COMPUTES FIND THE LARGEST ISLAND(S) OF TRACTABILITY



### Georg: «The evil is in cycles»

#### **Acyclic queries & join trees**

Join tree: a tree whose nodes are labelled by query hyperedges (or query atoms) such that:

- each hyperedge labels some node, and
- For each query variable V, the tree-nodes containing V span a connected subtree (connectednes condition)



#### Acyclic conjunctive query (ACQ):

A Query whose associated hypergraph is acyclic (more precisely,α-acyclic [Fagin 83])

Query acyclicity was independently defined by

- [Beeri et al. STOC81] acyclic database schemas, and
- [Goodman & Shmueli 1981, TODS'82] tree queries

[Graham; Yu & Özsoyoğlu] GYO reduction

A query is acyclic iff it has a join tree

#### **Good properties:**

- ACQs can be recognized in polynomial (actually linear) time,
- A join-tree for an ACQ can be built in linear time,
- A Bolean ACQ Q can be answered in time  $O(|Q| \times |r_{max}| \times \log |r_{max}|)$  [Yan'81]
- A non-Boolean ACQ can be answered with polynomial delay.

# Generalizing acyclicity

#### **Tree Decomposition**

 $ans \leftarrow a(S, X, X', C, F) \land b(S, Y, Y', C', F') \land c(C, C', Z) \land d(X, Z) \land$  $e(Y,Z) \wedge f(F,F',Z') \wedge g(X',Z') \wedge h(Y',Z') \wedge$  $j(J, X, Y, X', Y') \wedge p(B, X', F) \wedge q(B', X', F)$ S

J,X,Y X' Y'

S,X,X',C,F,Y,Y',C',F'

- Variables of each atom covered by some bag
  - X',Y',F,F',Z' X,YC,C',Z X',Z',F,Z' Y',Z' X.Z Y.Z **Connectedness condition** B' X' F B,X' F **Tree Decomposition of width 8**

#### The power of hyperedges: Generalized Hypertree Decomposition

 $ans \leftarrow a(S, X, X', C, F) \land b(S, Y, Y', C', F') \land c(C, C', Z) \land d(X, Z) \land$  $e(Y, Z) \land f(F, F', Z') \land g(X', Z') \land h(Y', Z') \land$  $j(J, X, Y, X', Y') \land p(B, X', F) \land q(B', X', F)$ 



#### Unfortunately:

GHW is NP-hard to compute, even for small width.

**Theorem** [G., Miklos, Schwentick 07+09]:

Checking whether ghw(Q)=3 is NP complete

Thus, GHDs do not fulfill criterion 2 (efficient recognizability).

→ Slightly restrict GHDs using a *special condition*, yielding HDs

#### Hypertree Decomposition = GHD +Special Condition



#### Hypertree Decomposition = GHD +Special Condition



# Good properties of (G)HDs and of queries of bounded (G)HW



Hybrid decompositions: matching physical DB parameters (G)HDs are query plans. Based on available selectivity and cardinality indexes, we can associate a cost to each such decomposition:



$$cost(D) = \sum_{n \in N} c(n) + \sum_{e \in E} c'(e)$$

Work by Scarcello, Greco, Leone [PODS 04; JCSS 07]:

**Theorem:** Finding a minimum cost HD (or GHD) is NP-hard.

Note: This is just as bad as classical query optimization.

More surprising: Problem is tractable for a slight restriction of HDs:

**Theorem:** Finding a minimum cost *normal-form* HD is tractable.

A k-width HD is in normal form iff  $% \mathcal{A}$  is generated by  $\operatorname{Opt-k-Dekomp}$ 

Algorithm Cost-k-Dekomp



# Many more works implementing HDs and generalizations

- [Afrati, Joglekar, Ré, Salihoglu, Ullman, ICDT 2017]: GYM: A Multiround Join Algorithm In MapReduce.
- [Koch , Ahmad, Kennedy, Nikolic, Nötzli, Lupei, Shaikhha VLDBJ 2014] DBToaster: higher-order delta processing for dynamic, frequently fresh views.
- [Tu,Ré SIGMOD 2015] DunceCap: Query Plans Using Generalized Hypertree Decompositions
- [Abo Khamis, Ngo, Rudra PODS 2016 Best Paper] FAQ: Questions Asked Frequently
- [Aberger, Tu, Olukotun, Ré, SIGMOD 2016] EmptyHeaded: A Relational Engine for Graph Processing.
- [Joglekar, Puttagunta, Ré, PODS'16] AJAR: Aggregations and Joins over Annotated Relations
- [Khamis, Ngo, Suciu, PODS'17] What Do Shannon-type Inequalities, Submodular Width, and Disjunctive Datalog Have to Do with One Another?

# Beyond hypertrees (?)

Unknown recognizability Exponential (but FP)-time query answering





**Theorem** [Grohe & Marx 06] : The answer to a query of fractional cover weight  $\rho^{*}(Q)$  can be computed in time  $|Q| \times rmax^{\rho^{*}(Q)+O(1)}$ 

#### Observe: $\{Q \mid \rho^*(Q) \le k\}$ and $\{Q \mid hw(Q) \le k\}$ are incomparable.

To combine the two notions profitably, Grohe and Marx defined

Fractional Hypertree Decompositions (FHDs) and correspondingly FHW



#### The AGM bound [Atezerias, Grohe, Marx '08]

*Let* q *be a full conjunctive query. For every fractional edge cover*  $\mathbf{u}$  *of* q*, we have:* 

$$|q| \leq \prod_{j=1}^{\ell} N_j^{u_j}$$
  
min  $\sum_j \log_2(N_j) \cdot u_j$   
s.t. $\forall x \in vars(q) : \sum_{j:x \in vars(R_j)} u_j \geq 1$   
 $\forall R_j : u_j \geq 0$ 

The AGM bound is tight

Coloring bound [G. Gottlob, S.T. Lee, G. Valiant, P. Valiant '12]

Given query  $Q = R_0(u_0) \leftarrow R_{i_1}(u_1) \land \ldots \land R_{i_m}(u_m)$ 

This bound is designed to work with output variables, and in the presence of keys and functional dependencies

## Fractional hypertree decompositions

- Assume the fractional hypertree width is k
- Then, at least one vertex has a cover equal to k
- There exists a database that meets the worst-case bound



## Why we can go beyond FHDs



- Two different decompositions with «critical covers»  $C_1$  and  $C_2$
- From the AGM bound: there exists a database  $D_1$  where  $C_1$  meets the worst case bound and a database  $D_2$  where  $C_2$  meets the bound



- Two different decompositions with «critical covers»  $C_1$  and  $C_2$
- From the AGM bound: there exists a database  $D_1$  where  $C_1$  meets the worst case bound and a database  $D_2$  where  $C_2$  meets the bound

## Example: database $D_2$



• From the AGM bound: there exists a database  $D_1$  where  $C_1$  meets the worst case bound and a database  $D_2$  where  $C_2$  meets the bound

## Why we can go beyond FHDs



- In general, we have different decompositions with different «critical covers»
- It is possible that there exists no database that simultaneusly meets the worst-case AGM bound on <u>all</u> hypertree decompositions!

#### What can we do?

- Given a query and a database, choose the best possible decomposition
  - $\rightarrow$  use weighted hypertree decompositions

Hybrid approach

## A more powerful (structural) width?

Think of a measure (width) such that, for every database, there is a hypertree decomposition whose width does not exceed a given threshold

For some classes of queries, such a width can be strictly smaller than the fractional hypertree width

• See [Khamis, Ngo, Suciu 2017]

# A case study: cycles $q(x_1, x_2, x_3) \leftarrow t_1(x_1, x_2) \wedge t_2(x_2, x_3) \wedge t_3(x_1, x_3)$ $q(x_1, x_2, x_3) \leftarrow t_1(x_1, x_2) \wedge t_2(x_2, x_3) \wedge t_3(x_1, x_3)$ Assume $|t_i| = N$

DIFFERENT POSSIBLE RELATIONS

|             | 1 1            |   |
|-------------|----------------|---|
| 1 1 1 1 1 1 | 1 )            |   |
| 2 1 1 2 1 2 | 22             |   |
| 3 2 2 1 1 3 | 3 3            |   |
| 4 3 2 2 1 4 | <del>4</del> 4 | - |

A case study: cycles  $q(x_1, x_2, x_3) \leftarrow t_1(x_1, x_2) \land t_2(x_2, x_3) \land t_3(x_1, x_3)$ Assume  $|\tau_i| = N$ x,•\_\_\_\_ In this case, every variable can take at most TN values, and [g(b)] < N<sup>3</sup>/<sub>2</sub> 2. AGM Bound

A case study: cycles  $q(x_1, x_2, x_3) \leftarrow t_1(x_1, x_2) \land t_2(x_2, x_3) \land t_3(x_1, x_3)$ Assume  $|\tau_i| = N$ In this case,  $X_1$  take just one value and  $1 \ 2 \ |q(b)| \leq |T_{X_1} z_1 \times z_2| = |z_2|$  $1 \ 3 \ 1 \ 4$ 21=23 1





#### The cost of cycles $q \leftarrow t_1(X_1, X_2) \land t_2(X_2, X_3) \land \cdots \land \land t_p(X_p, X_1)$ Assume $|\tau_i| = N$ - × 2 r, (x, x, x) 1 2 (x, x3) 2 ×...• ·×. $\left[ 2_3(x_3, x_4) \wedge 2_4(x_4, x_1) \right] 2''$ EVERY CYCLE CAN BE SOLVED IN NOTE THAT THE WORST-CASE DATABASE $O(N^2)$ MUST HAVE MANY VALUES FOR (FRACTIONAL) HYPERTREE X1 AND X3, AND X2 IS A HEAVY WISTH: 2 HITTER AND SHOULS HAVE FEW VALUES



 $X_{1}$   $X_{2}$   $X_{5}$ 

Xz

Consider the case p=6

 $q \leftarrow t_1(X_1, X_2) \land t_2(X_2, X_3) \land \cdots \land \land t_p(X_p, X_1)$ 

 Consider one variable, say x6, with a few values in its active domain

#### The cost of cycles $q \leftarrow t_1(X_1, X_2) \land t_2(X_2, X_3) \land \cdots \land \land t_p(X_p, X_q)$ 2, (x, x2) 1 2, (x6, X1) 2 XI WITH OUR ASSUMPTION 22(X2,X3)~ 26(X6, -) 2" $\times_{z}$ eN ×6 THE COST 13 23(×3,×4)×26(×6, \_) 2" ALMOST LINEAR • X5 24(×4,×5) × 25(×5,×6) 2" ×4

 $q \leftarrow t_1(X_1, X_2) \land t_2(X_2, X_3) \land \cdots \land \land \uparrow (X_p, X_1)$ 



- Consider the case *p=6*
- Assume now that all variables have many values. In this case, they should be "almost keys", and a different kind of decomposition should be considered

 $q \leftarrow t_1(X_1, X_2) \land t_2(X_2, X_3) \land \cdots \land \land t_p(X_p, X_1)$ 



2, (×1,×2) ~ 22 (×2,×3) ~ 23 (×3,×4) 24(x4, ×5) × 25(×5, ×6) × 26(×6,×1)

WE HALVE THE CYCLE IN TWO SUBPROBLEMS HAVING AT MOST [] ATOMS Z

If we have quasi-keys of maximum degree d, the cost of each node is  $Nd^3$ 

(d=1 for actual keys)

XJ

Xu

X1

 $X_2$ 

Xz

 $q \leftarrow t_1(X_1, X_2) \land t_2(X_2, X_3) \land$ 

We can show that the following «width» bound holds:





The cost of cycles  $\Lambda \stackrel{\sim}{}_{\mathcal{P}} (X_{\mathcal{P}}, X_{\mathcal{I}})$  $q \leftarrow t_1(X_1, X_2) \land t_2(X_2, X_3) \land$ Assume  $|\tau_i| = N$ ×2 ×5 SUBQUADRATIC FOR EVERY CLASS Cp of cycles (having length at most p) To evaluate such a query, we get a cost of O(p<sup>1-1/[p/27</sup>. N<sup>2-1/[p/27</sup>)

 $\begin{array}{c} \times_{4} \\ \times_{2} \\ \times_{3} \\ \times_{3} \\ \end{array}$ 

 $q \leftarrow t_1(X_1, X_2) \land t_2(X_2, X_3) \land \cdots \land \cdots \land$  $X_1 \qquad Assume |t_i| = N$ 

THE SUBRUADRATIC BOUND COMES FROM MARX'S SUBMOBULAR WIDTH:

) · Horizontal gragments ? · Different decompositions

To evaluate such a quary, we get a cost of O(p1-1/17/27. N2-1/17/27)

#### Surprising results

- The submodular width [Marx 2013] is based on data-dependent decompositions, but
  - It provides a structural worst-case bound
- Our first feeling was that (hypertree) decompositions techniques were practically useful for long and complex queries, but
  - It turns out that we can get important results even for simple queries, such as the short cycles

#### Some hot issues for future research

- The challenge: when conjunctive queries are tractable?
- Can minimum cost decompositions be well approximated?
- Better integrate classical query optimization methods with HDs
- Provide efficient FPT algorithms that match the submodular width
  - N.B.: the PANDA algorithm [Khamis, Ngo, Suciu 2017] is not far from that, but the  $\tilde{O}$  notation hides a  $\log(n)^{|f(q)|}$  factor (it is not FPT)

# Appendix

#### Can HDs be applied outside CQs?

#### Example of CSP: Crossword Puzzle

| 2  | 3             | 4                           | 5                                |  | 6   |
|----|---------------|-----------------------------|----------------------------------|--|---|
|    | +             |                             | 8                                | 9  | 10  |
| 12 | 13            |                             | 14                               |  | 15  |
|    | 17            |                             | 18                               |  | 19  |
| 21 | 22            | 23                          | 24                               | 25   | 26  |
|    | 2<br>12<br>21 | 2 3<br>12 13<br>17<br>21 22 | 2 3 4<br>12 13<br>17<br>21 22 23 | 2       3       4       5         12       13       14         17       18         21       22       23       24 | 2       3       4       5         12       13       14         17       18         21       22       23       24       25 |



1h: |

P A R I S P A N D A L A U R A A N I T A 1v: LIMBO

LINGO PETRA PAMPA PETER

# Constraint satisfaction problems: Renault example<sub>(1/2)</sub>

• Renault Megane configuration [Amilhastre, Fargier, Marquis AIJ, 2002] Used in CSP competitions and as a benchmark problem



- Variables encode type of engine, country, options like air cooling, etc.
- The considered instances consist of about 150 atoms/constraints and 110 variables, with database instances where attributes have at most 42 distinct values, and the largest constraint relation contains 48721 tuples.

# Constraint satisfaction problems: Renault example<sub>(2/2)</sub>



- We discovered that the generalized hypertree width is 3 for most instances (with a maximum of 4).
- The total number of solutions is about 2 · 10<sup>12</sup>, however this information is not very meaningful, because of the many auxiliary variables occurring in the problem.
- Rather, by using the algorithms based on generalized hypertree decompositions, it is possible to compute the solutions of these instances (or just their number) over the actual variables of interest.
- None of the other available engines that we know, either in the database or in the CSP community, were able to compute such a result for those large instances.

For more information, we refer to the Hypertree Decomposition web-page: <u>http://www.dimes.unical.it/scarcello/Hypertrees/</u>

# **Combinatorial Auctions**

Bidders can place bids on
Packages of items.
<u>Winner determination:</u> Choose a set of
compatible bids of maximum revenue
or minimum cost.
For classical auctions, winner determination
is obviously tractable. Not so for CAs.

Interesting tractable results based on the hypertree width of the dual hypergraphs [Gottlob & Greco, JACM 2013]



## Applications in different domains





London Regional Transports: Combinatorial auctions of bus routes. Private bus companies bid on bundles of routes.

#### Airport slot auction

