The Quest for Faster Join Algorithms

Paris Koutris



Joins are a core relational operator

The Long History of Join Algorithms

- typically computed via a query plan of **binary join operators**
- many different algorithms to compute binary joins
 - nested loop join
 - hash join
 - sort-merge join
 - .. and a lot more!

 $R(A, B) \bowtie S(B, C) \bowtie T(C, D)$



But haven't we already solved everything?

But haven't we already solved everything?



- errors leading to suboptimal plans
- - pattern matching on graphs
 - cyclic joins
- 3. There is a large gap between theory and practice!

No!

1. Query optimizers are notoriously sensitive to cardinality estimation

2. Binary join plans perform well for most joins, but poorly for others:

What is the fastest algorithm to compute a join query?

We still don't know!

What is the fastest algorithm to compute a join query?

But we have made some amazing progress!

In this talk:

the story about the progress and a few key ideas behind it

1980

Yannakakis' algorithm

tree decompositions generalized hypertree width

Yannakakis' algorithm

tree decompositions generalized hypertree width



AGM bound



AGM bound

Roadmap



PANDA algorithm submodular width



Some background

this talk: queries with **natural joins** + **projections**

$\pi_{\emptyset}(R(A,B) \bowtie S(B,C) \bowtie T(C,A))$

Full: no projections

Yannakakis' Algorithm

 $R(A,B) \bowtie S(B,C) \bowtie T(C,D) \bowtie U(C,E)$

join tree

[Yannakakis - '81]

Yannakakis' Algorithm

[Yannakakis - '81]

 $R(A,B) \bowtie S(B,C) \bowtie T(C,D) \bowtie U(C,E)$

Yannakakis' Algorithm

$R(A,B) \bowtie S(B,C) \bowtie T(C,D) \bowtie U(C,E)$

Key property: any intermediate join size is bounded by the output size

[Yannakakis - '81]

Yannakakis' Algorithm [Yannakakis - '81]

Yannakakis' algorithm comes with runtime guarantees for acyclic joins

joins that have a join tree

Boolean

General

O(N)

 $O(N \cdot OUT + OUT)$

Predicate Transfer: Yannakakis in Practice

- Yannakakis can be slow because of the high cost of a semi-join
- We can replace the semi-joins by (approximate) Bloom filters [Yang et al. '24]
- We obtain **speedup** + **robustness** to bad query plans

1 - transfer phase

R forward-backward transfer of Bloom filters U

 $R(A,B) \bowtie S(B,C) \bowtie T(C,D) \bowtie U(C,E)$

Predicate Transfer: Yannakakis in Practice

Robustness Factor := max execution time / min execution time

worst join order best join order

[Zhang et al. - '25]

Predicate Transfer: Yannakakis in Practice

Robustness Factor := max execution time / min execution time

worst join order

RF	TPC-H			JOB			TPC-DS		
	Avg	Min	Max	Avg	Min	Max	Avg	Min	Max
DuckDB	2.7	1.2	9.3	30.4	1.1	371	7.2	1.0	224
RPT	1.3	1.2	1.5	1.2	1.0	1.6	1.1	1.0	1.5

Table 2: Robustness Factors for bushy joins.

RF	TPC-H			JOB			TPC-DS		
	Avg	Min	Max	Avg	Min	Max	Avg	Min	Max
DuckDB	5.1	1.2	13.7	120	1.1	1747	35.0	1.0	1226
RPT	1.8	1.2	3.0	1.6	1.1	7.7	1.8	1.0	4.2

[Zhang et al. - '25]

best join order

Table 1: Robustness Factors for left-deep joins.

1 - Tree Decompositions

- Acyclic joins have plans where the output size.
- What can we do for cyclic joins?

• Acyclic joins have plans where the intermediate size is bounded by the

IDEA: make the join acyclic by computing larger intermediate relations

 $O(N^2 + OUT)$

Tree Decompositions - Definition

- Like a join tree, but each node is an intermediate relation (**bag**)
- Every relation must be included in some bag

Generalized HyperTree Width

integral edge cover w: smallest subset of relations that "covers" all attributes in a bag

Generalized HyperTree Width

integral edge cover w: smallest subset of relations that "covers" all attributes in a bag

$ghw(Q) := \min \max_{\substack{w(X_B) \\ decomp \ T \ bag \ B}} \max_{w(X_B)}$

$$ghw = 2$$

Runtime Result

If Q is a full join, we can evaluate it in time

 $T_1(A, B, C) \bowtie T_2(C, D, E)$

 $O\left(N^{ghw(Q)} + OUT\right)$

$O(N^2 + OUT)$

The algorithm can be expressed as a join-project plan!

2 - Data Partitioning

Join-at-a-time algorithms are **suboptimal!**

We can do better by **partitioning** the input data and applying a different join plan to each partition

best binary join algorithm: $O(N^2)$

we can do better: $O(N^{3/2})$

Data Partitioning: Example

 $R(A, B) \bowtie S(B, C) \bowtie T(C, A)$

 $O(N^{3/2})$

Fractional Hypertree Width

fractional edge cover *f*: assign weights to relations to "cover" all attributes in a bag

Fractional Hypertree Width

fractional edge cover *f*: assign weights to relations to "cover" all attributes in a bag

$fhw(Q) := \min \max_{\substack{a \in A \\ b \in B}} \max_{a \in B} f(X_B)$

$$fhw = 3/2$$

Runtime Result Revisited

If Q is a full join, we can evaluate it in time

 $T_1(A, B, C) \bowtie T_2(C, D, E)$

 $O\left(N^{fhw(Q)} + OUT\right)$

[Ngo, Re, Rudra - '13]

 $O(N^{3/2} + OUT)$

The algorithm can be expressed as a join-project-partition plan!

- WCOJs have a running time matching the worst-case output size • They often work in a attribute-at-a-time manner • Leapfrog Triejoin [Veldhuizen - '14]
- - GenericJoin [Ngo, Re, Rudra - '**13**]
 - FreeJoin [Wang, Willsey, Suciu - '23]

A lot of successful practical implementations!

(Worst-case Optimal Joins)

AGM bound

[Atserias, Grohe, Marx - '08]

Data Partitioning Again

We can partition data to use different tree decompositions!

fhw = 2

Runtime Result (Re)Revisited

If Q is a full join, we can evaluate it in time $\tilde{O}\left(N^{subw(Q)} + OUT\right)$

$subw(Q) := \max_{\substack{\text{submodular } h \text{ decomp } T \text{ bag } B}} \min_{\substack{\text{submodular } h \text{ decomp } T \text{ bag } B}} \max_{\substack{\text{submodular } h \text{ decomp } T \text{ bag } B}} \max_{\substack{\text{submodular } h \text{ decomp } T \text{ bag } B}} \max_{\substack{\text{submodular } h \text{ decomp } T \text{ bag } B}} \max_{\substack{\text{submodular } h \text{ decomp } T \text{ bag } B}} \max_{\substack{\text{submodular } h \text{ decomp } T \text{ bag } B}} \max_{\substack{\text{submodular } h \text{ decomp } T \text{ bag } B}} \max_{\substack{\text{submodular } h \text{ decomp } T \text{ bag } B}} \max_{\substack{\text{submodular } h \text{ decomp } T \text{ bag } B}} \max_{\substack{\text{submodular } h \text{ decomp } T \text{ bag } B}} \max_{\substack{\text{submodular } h \text{ decomp } T \text{ bag } B}} \max_{\substack{\text{submodular } h \text{ decomp } T \text{ bag } B}} \max_{\substack{\text{submodular } h \text{ decomp } T \text{ bag } B}} \max_{\substack{\text{submodular } h \text{ decomp } T \text{ bag } B}} \max_{\substack{\text{submodular } h \text{ decomp } T \text{ bag } B}} \max_{\substack{\text{submodular } h \text{ decomp } T \text{ bag } B}} \max_{\substack{\text{submodular } h \text{ decomp } T \text{ bag } B}} \max_{\substack{\text{submodular } h \text{ decomp } T \text{ bag } B}} \max_{\substack{\text{submodual } h \text{ decomp } T \text{ bag } B}} \max_{\substack{\text{submodual } h \text{ decomp } T \text{ bag } B}} \max_{\substack{\text{submodual } h \text{ decomp } T \text{ bag } B}}} \max_{\substack{\text{submodual } h \text{ decomp } T \text{ bag } B}}} \max_{\substack{\text{submodual } h \text{ decomp } T \text{ bag } B}}} \max_{\substack{\text{submodual } h \text{ decomp } T \text{ bag } B}}} \max_{\substack{\text{submodual } h \text{ decomp } T \text{ bag } B}}} \max_{\substack{\text{submodual } h \text{ decomp } T \text{ bag } B}}} \max_{\substack{\text{submodual } h \text{ decomp } T \text{ bag } B}}} \max_{\substack{\text{submodual } h \text{ decomp } T \text{ bag } B}}} \max_{\substack{\text{submodual } h \text{ decomp } T \text{ bag } B}}} \max_{\substack{\text{submodual } h \text{ decomp } T \text{ bag } B}}} \max_{\substack{\text{submodual } h \text{ decomp } T \text{ bag } B}}} \max_{\substack{\text{submodual } h \text{ decomp } T \text{ bag } B}}} \max_{\substack{\text{submodual } h \text{ decomp } T \text{ bag } B}}} \max_{\substack{\text{submodual } h \text{ decomp } T \text{ bag } B}}} \max_{\substack{\text{submodual } h \text{ decomp } T \text{ bag } B}}} \max_{\substack{\text{submodual } h \text{ decomp } T \text{ bag } B}}} \max_{\substack{\text{submodual } h \text{ decomp } T \text{ bag } B}}} \max_{\substack{\text{submodual } h \text{ decomp } T \text{ bag } B}}} \max_{\substack{\text{submodual } h \text{ decomp } T \text{ bag } B}}} \max_{\substack{\text{submodual } h \text{ decomp } T \text{ bag } B}}} \max_{\substack{\text{submodual } h \text{ de$

[Khamis, Ngo, Suciu - '17]

subw \leq fhw \leq ghw

PANDA algorithm

Output-Sensitive Yannakakis

- $O(N + N \cdot OUT)$
- into account (using **data partitioning** again)

If Q is an acyclic join, we can evaluate it in time $O(N + N \cdot OUT^{1-\epsilon})$

$\pi_{A,C}(R(A,B) \bowtie T(B,C)) \qquad O(N+N \cdot OUT^{1/2})$

• When we have projections in an acyclic join, Yannakakis runs in time

• We can design faster join algorithms by taking the output size OUT

[Deep, Zhao, Fan, Koutris - '25] [Hu - '**25**]

3 - Using Statistical Information

By using more **statistical information**, we can design faster join algorithms

primary key $R(A,B) \bowtie S(B,C) \bowtie T(C,A)$ runs in time only O(N)

3 - Using Statistical Information

Using more detailed statistics • restricts what is the worst-case behaved input closes the gap between theory and practice

By using more statistical information, we can design faster join algorithms

primary key $R(A,B) \bowtie S(B,C) \bowtie T(C,A)$ runs in time only O(N)

3 - Using Statistical Information

Gottlob et al. - '12

Yannakakis

input/output

functional dependencies

cardinalities

AGM bound

We can design join algorithms with stronger **enumeration** guarantees

- useful for LIMIT queries
- constant delay enumeration [Bagan et al. '07]

O(N + OUT)VS

We can design join algorithms with stronger **enumeration** guarantees

- useful for LIMIT queries
- constant delay enumeration [Bagan et al. '07]

O(N + OUT)VS

We can design join algorithms with stronger **enumeration** guarantees

- useful for LIMIT queries
- constant delay enumeration [Bagan et al. '07]

O(N + OUT)VS

We can design join algorithms with stronger **enumeration** guarantees

- useful for LIMIT queries
- constant delay enumeration [Bagan et al. '07]

O(N + OUT)VS

We can design join algorithms with stronger **enumeration** guarantees

- useful for LIMIT queries
- constant delay enumeration [Bagan et al. '07]

O(N + OUT)VS

We can design join algorithms with stronger **enumeration** guarantees

- useful for LIMIT queries
- constant delay enumeration [Bagan et al. '07]

O(N + OUT)VS

Runtime Result Revisited

If Q is a full join, we can do constant delay enumeration with preprocessing time

 $O(N^{subw(Q)})$

Acyclic queries only need linear preprocessing time!

Factorized Databases

preprocessing

The intermediate data structure can be viewed as a compressed (*factorized*) representation of the join output that can be decompressed efficiently

[Olteanu,Zavodny - '16]

constant delay enumeration

Joins with Ranked Order

- 1. Joins where we only want top-ranked results (ORDER BY + LIMIT)
- 2. We can construct algorithms with strong enumeration guarantees
- Works well in practice it avoids intermediate *expensive materializations* 3.

SELECT * FROM R, S For acyclic joins, after O(N) preprocessing time, we can enumerate the output while spending only WHERE R.A = S.B**logarithmic** time between two consecutive outputs ORDER BY S.C + R.D [Deep, Koutris - '21] LIMIT 10 [Tziavelis et al. - '20]

5 - The Algebraic Lens

structures called **semirings**

- aggregation becomes simple sum, count, min, max
- extremely useful for ML pipelines (FAQs, factorized databases)

• join: multiplication (X) • project, union: **addition** (+)

Most existing join algorithms can be *lifted* to work on general algebraic [Green, Karvounarakis, Tannen - '21]

Triangles with Semirings

 $(\mathbb{N}, \min, +, \infty, 0)$

 $(\mathbb{N}, +, \cdot, 0, 1)$ + **arithmetic semiring** = how many triangles?

 $({0,1}, \lor, \land, 0, 1)$ + **Boolean semiring** = is there is a triangle?

+ **tropical semiring** = minimum weight triangle

The Algebraic Lens

If Q is an acyclic full join, we can compute it over any semiring it in time

O(N + OUT)

$R(A,B) \bowtie S(B,C) \bowtie T(C,D)$

and many other results...

FAQs [Khamis, Ngo, Rudra - '16] [Khamis et al. - '19]

What comes next?

Towards Even Faster Joins

We can join faster using Fast Matrix Multiplication (FMM) $\pi_{A,C}(R(A,B) \bowtie S(B,C))$

 $\pi_{\varnothing}(R(A,B) \bowtie S(B,C) \bowtie T(C,A))$

- combinatorial: $O(n^3)$
- FMM: $O(n^{\omega}), \omega = 2.371339$

combinatorial: $O(N^{3/2})$ with FMM: $O(N^{\frac{2\omega}{1+\omega}})$

[Alon, Yuster, Zwick - '97]

Towards Even Faster Joins

- Using MM can also help in practice

Two Path Join 10^{5} Running time in sec 10^{4} 10^{3} [Deep, Hu, Koutris - '22] 10^2 10^{-10} 10^{-1} RoadNet DBLP Words Protein Image Jokes Datasets MMJoin EmptyHeaded **Q1:** Are there other ways to use FMM in query plans? Postgres Non-MMJoin MySQL System X **Q2:** What other non-combinatorial techniques can we use?

• FMM can be incorporated as a new operator in a query plan [Khamis, Hu, Suciu - '25]

Joins for any Semiring - how do we count?

- The PANDA algorithm *cannot count*!
- The key problem is the overlap of output tuples in tree decompositions
- We can count the join size in time $O(N^{\#SubW})$ [Khamis et al. '19]

subw \leq #subw \leq fhw

Q3: What is the best algorithm that counts the size of a join?

(N^{#Subw}) [Khamis et al. - '19]

Lower Bounds for Joins

- Not (currently) feasible to show unconditional bounds
- believe we have an optimal algorithm

Q4: What other conditional lower bounds should we use?

• To show that join algorithms are optimal, we need lower bounds

• We use *conditional lower bounds* by reducing from problems where we

(conditional) lower bounds for all joins

 $\Omega(N^{clemb})$ [Fan, Koutris, Zhao - '23]

clemb \leq subw \leq fhw

Lower Bounds for Joins

- Another approach is to restrict the algorithm
- In particular, we ask that the algorithm structure is encoded via a **circuit**
- We can show matching upper + lower bounds for any join [Fan, Koutris, Zhao '24]

From Theory to Practice

- some have already been successful: WCOJs, predicate transfer
- for some it is still unclear: enumeration?
- ... and others are only theoretical (and will remain)

Q5: Which theoretical developments can be effective in practice?

Thank You!

5 Key Ideas:

- 1. tree decompositions
- 2. data partitioning
- 3. using statistics
- 4. enumeration
- 5. the algebraic lens

