

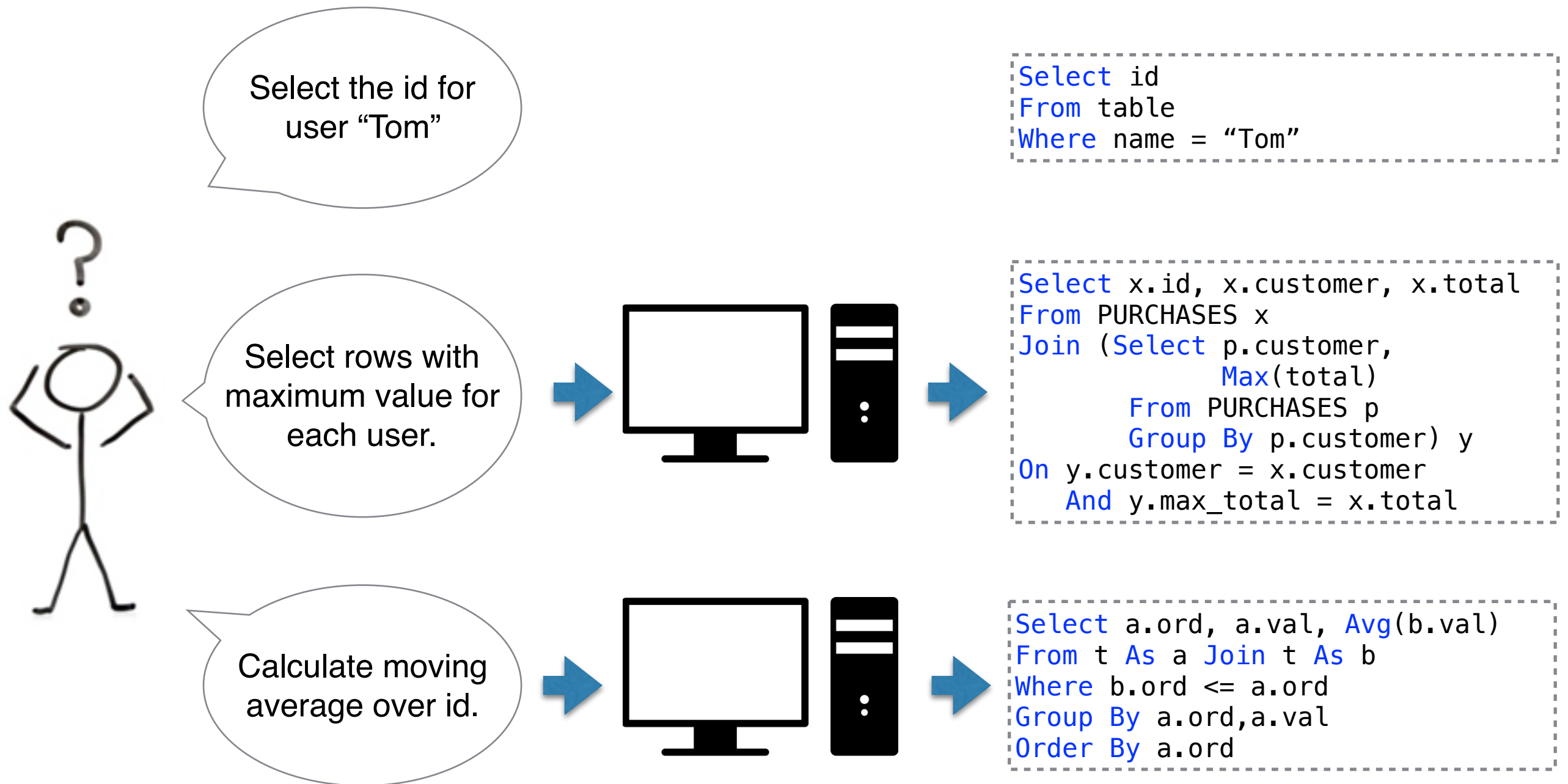
Synthesizing Highly Expressive SQL Queries From Input-Output Examples

<http://scythe.cs.washington.edu>

Chenglong Wang, Alvin Cheung, Ras Bodík
University of Washington

Tasks

SQL Query



Problem: Advanced SQL operators make SQL powerful but hard to master.

How to select the first N rows of each group?

I have two SQLite tables like this:

AuthorId	AuthorName
1	Alice
2	Bob
3	Carol

BookId	AuthorId	Title
1	1	aaa1
2	1	aaa2
3	1	aaa3
4	2	ddd1
5	2	ddd2
19	3	fff1
20	3	fff2
21	3	fff3
22	3	fff4

I want to make a SELECT query that will return the first N (e.g. two) rows for each AuthorId, ordering by Title ("Select the first two books of each author").

Sample output:

BookId	AuthorId	AuthorName	Title
1	1	Alice	aaa1
2	1	Alice	aaa1
4	2	Bob	ddd1
5	2	Bob	ddd2
19	3	Carol	fff1
20	3	Carol	fff2

How can I build this query?

(Yes, I found a similar topic, and I know how to return only one row (first or top). The problem is with the two).

[sql](#) [sqlite](#) [greatest-n-per-group](#) [limit-per-group](#)

Synthesize queries from ...?

Input Example

AuthorId	AuthorName
1	Alice
2	Bob
3	Carol

BookId	AuthorId	Title
1	1	aaa1
2	1	aaa2
3	1	aaa3
4	2	ddd1
5	2	ddd2
19	3	fff1
20	3	fff2
21	3	fff3
22	3	fff4

Output Example

BookId	AuthorId	AuthorName	Title
1	1	Alice	aaa1
2	1	Alice	aaa2
4	2	Bob	ddd1
5	2	Bob	ddd2
19	3	Carol	fff1
20	3	Carol	fff2

Constants

{2}

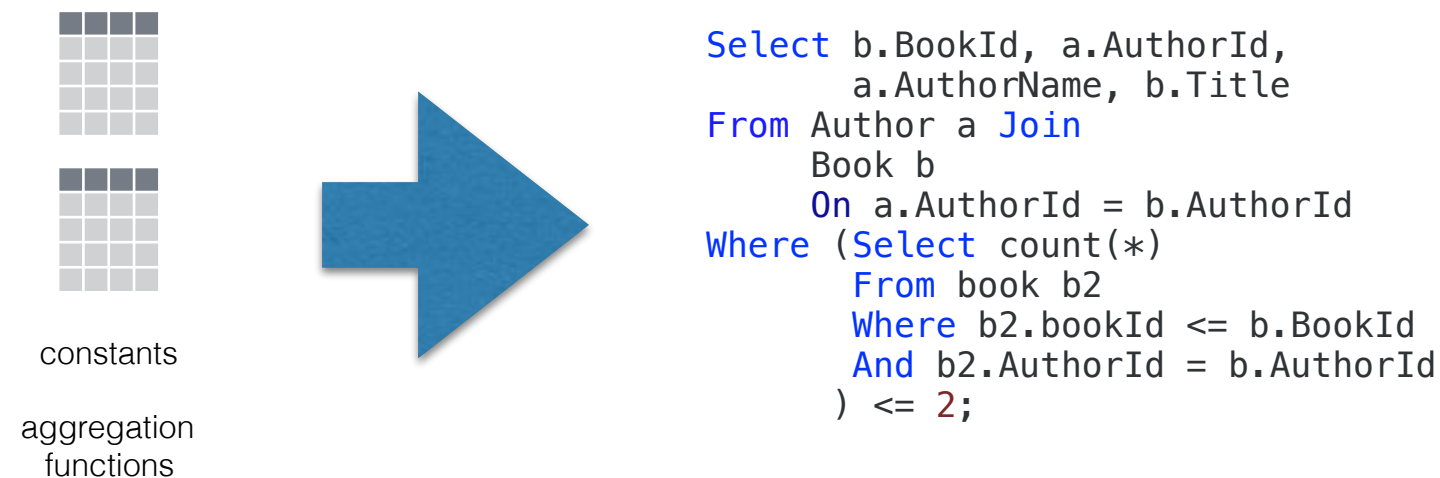
Aggregation Functions (Optional)

{ Count, Max, Min, Sum, Avg ... }

Key: The synthesizer takes inputs that users can provide online.

Talk Outline

- Motivation & Problem Definition
- Synthesis Algorithm



- Evaluation on Stack Overflow Posts

Running Example

Task: Collect the max vals below 50 for all `oid` groups in T2 and join them with T1.

T1

id	date	uid
1	12/25	1
2	11/21	3
4	12/24	2

T2

oid	val
1	30
1	10
1	10
2	50
2	10

Out

oid	date	uid	oid	MaxVal
1	12/25	1	1	30
4	12/24	2	2	10

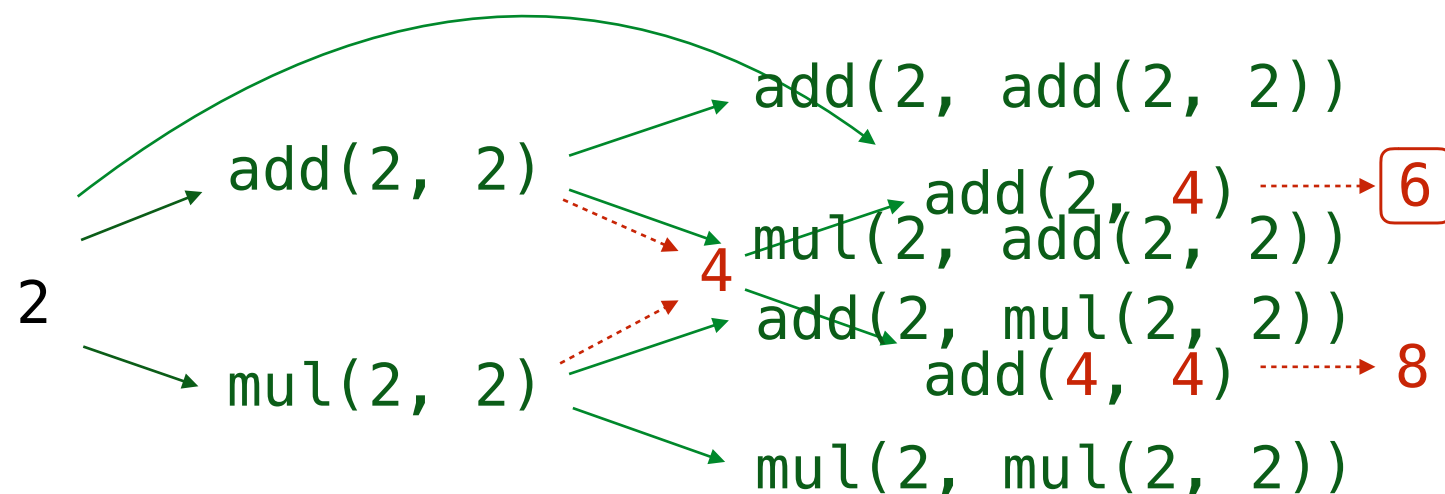
```
Select *  
From (Select oid, Max(val)  
      From T2  
      Where val < 50  
      Group By oid) T3  
Join T1  
On T3.oid = T1.uid
```

Constants = { 50 }

AggrFunc = { Max, Min }

Basic Algorithm: Enumerative Search

Input: 2 Output: 6 Operators: add, mul



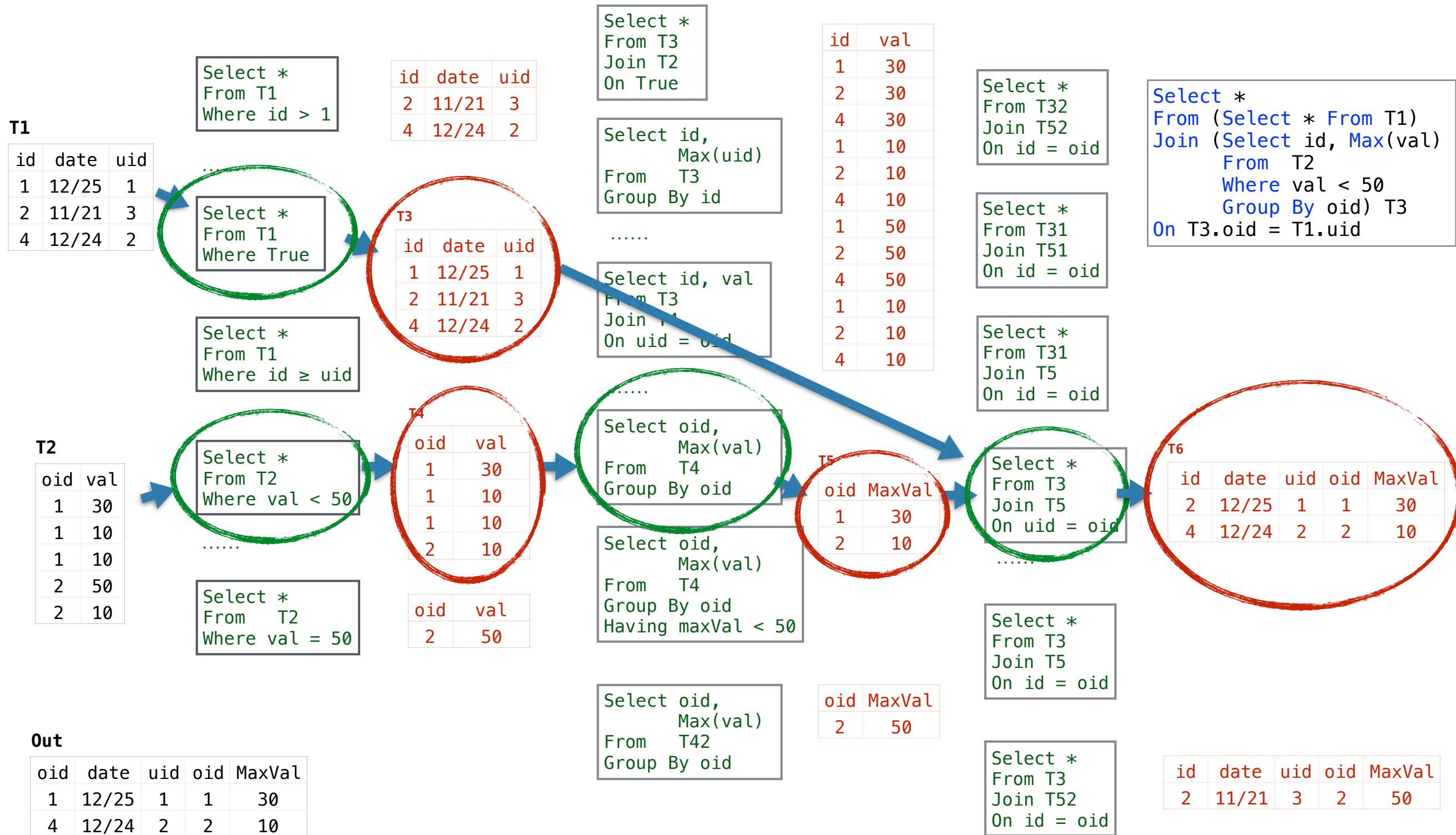
Synthesize Distributed protocols,
Super-optimization

$$\text{add}(2, \text{add}(2, 2)) = 6$$

$$\text{add}(2, \text{mul}(2, 2)) = 6$$

Key: Compressing the search space by memoizing values.

Input: T1, T2 Output: T_{out} Operators: Select, Join, Aggr



Input: T1, T2 Output: T_{out} Operators: Select, Join, Aggr

T1

id	date	uid
1	12/25	1
2	11/21	3
4	12/24	2

```
Select *
From T1
Where id > 1
```

.....

```
Select *
From T1
Where True
```

```
Select *
From T1
Where id ≥ uid
```

T2

oid	val
1	30
1	10
1	10
2	50
2	10

```
Select *
From T2
Where val < 50
```

.....

```
Select *
From T2
Where val = 50
```

Out

oid	date	uid	oid	MaxVal
1	12/25	1	1	30
4	12/24	2	2	10

id	date	uid
2	11/21	3
4	12/24	2

T3

id	date	uid
1	12/25	1
2	11/21	3
4	12/24	2

```
Select *
From T3
Join T2
On True
```

```
Select id,
        Max(uid)
From T3
Group By id
```

.....

```
Select id, val
From T3
Join T4
On uid = oid
```

.....

```
Select oid,
        Max(val)
From T4
Group By oid
```

```
Select oid,
        Max(val)
From T4
Group By oid
Having maxVal < 50
```

```
Select oid,
        MaxVal
From T4
Group By
```

id	date	uid	oid	MaxVal
1	12/25	1	1	30
2	11/21	1	1	30
4	12/24	2	1	30
1	12/25	1	1	10
2	11/21	1	1	10
4	12/24	2	1	10

id	date	uid	oid	MaxVal
1	12/25	1	1	30
2	11/21	1	1	30
4	12/24	2	1	30
1	12/25	1	1	10
2	11/21	1	1	10
4	12/24	2	1	10
1	12/25	1	2	50
2	11/21	1	2	50
4	12/24	2	2	50
1	12/25	1	2	10
2	11/21	1	2	10
4	12/24	2	2	10

oid	MaxVal
1	30
2	10

```
From T3
Join T5
On id = oid
```

```
Select *
From (Select * From T1)
Join (Select id, Max(val)
      From T2
      Where val < 50
      Group By oid) T3
On T3.oid = T1.uid
```

T6

id	date	uid	oid	MaxVal
2	11/21	3	2	50

**Challenge 2: Big tables
1,889 --> 42,600 cells**

**Challenge 1: Large number
of queries per-stage.
~500,000 in the last stage.**

id	date	uid	oid	MaxVal
2	11/21	3	2	50

Problem: Value-based compression is inefficient & ineffective.

Insight: Decompose Search Process

T1

id	date	uid
1	12/25	1
2	11/21	3
4	12/24	2

T2

oid	val
1	30
1	10
1	10
2	50
2	10

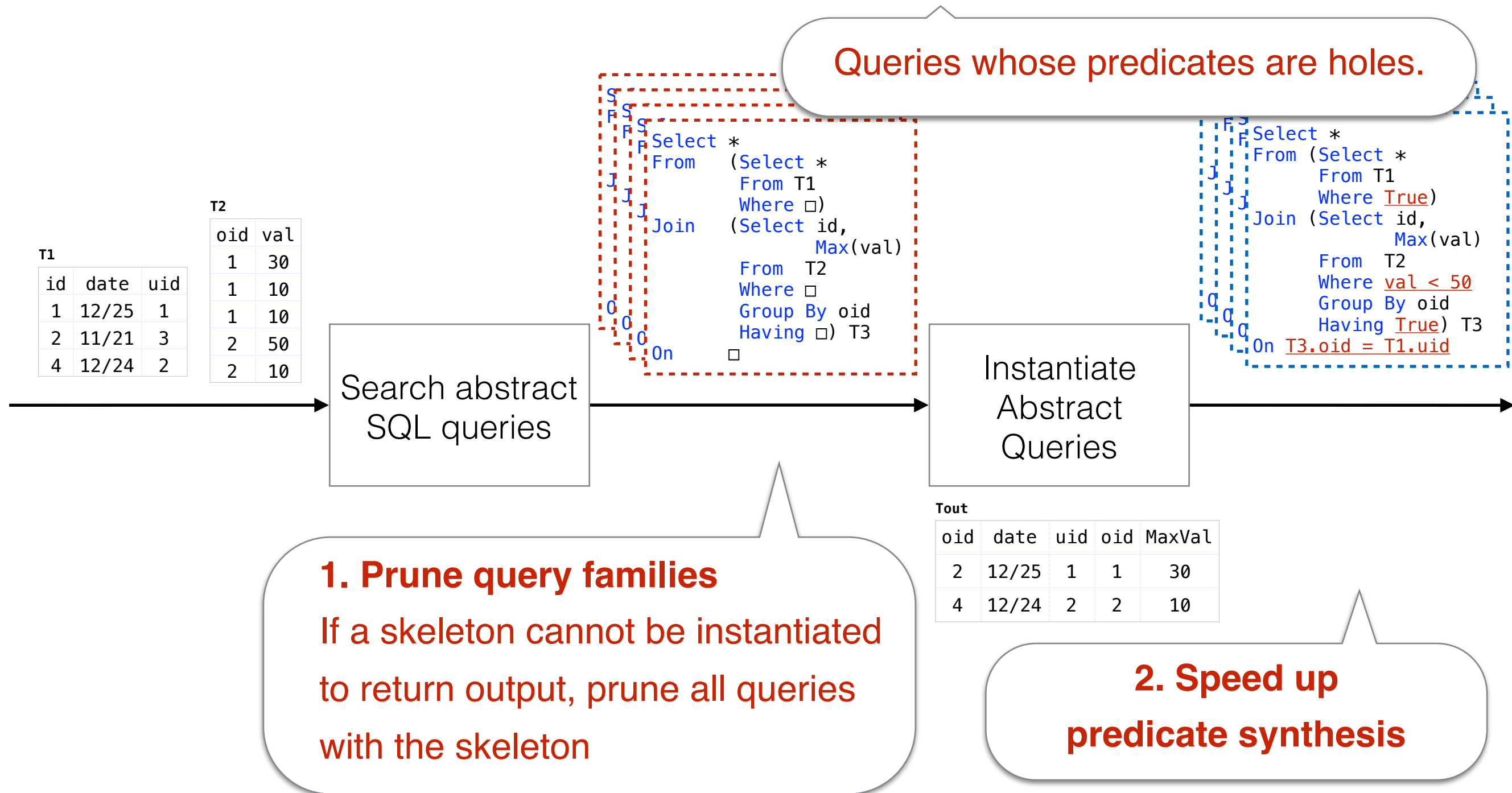
Tout

oid	date	uid	oid	MaxVal
2	12/25	1	1	30
4	12/24	2	2	10

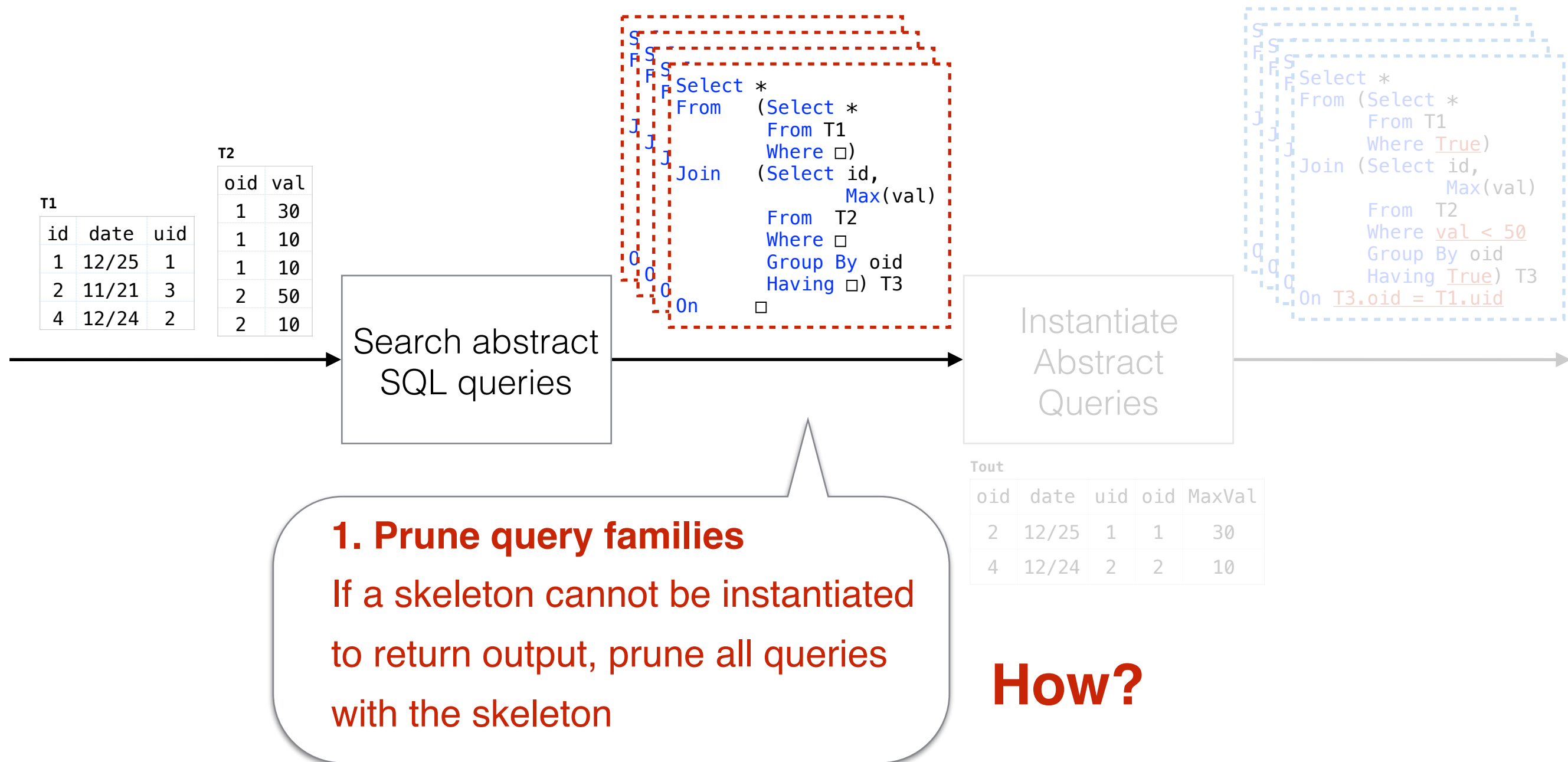
Search SQL queries

```
Select *
From (Select *
      From T1
      Where True)
Join (Select id,
            Max(val)
      From T2
      Where val < 50
      Group By oid
      Having True) T3
On T3.oid = T1.uid
```

Insight: Decompose Search Process With Abstract Queries



Insight: Decompose Search Process With Abstract Queries



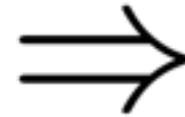
Evaluating Abstract Queries with Over-Approximation

Inductively defined
over abstract SQL
operators

T1

id	date	uid
1	12/25	1
2	11/21	3
4	12/24	2

Select id, date
From T1
Where \square

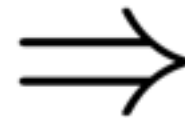


Summary

id	date
1	12/25
2	11/21
4	12/24

\cup

Select id, date
From T1
Where id \leq 2

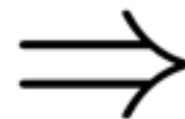


id	date	uid
1	12/25	1
2	11/21	3

T2

oid	MaxVal
1	30
2	10

T1 Join T2
On \square



Summary2

id	date	uid	oid	MaxVal
1	12/25	1	1	30
2	11/21	1	1	30
4	12/24	2	1	30
1	12/25	1	1	10
2	11/21	1	1	10
4	12/24	2	1	10

Key: Evaluating abstract queries into over-approximations of concrete query results.

Pruning with Abstract Queries

Input: T1, T2, **Output:** T_{out}, **Operators:** Select, Aggr, Join

T1

id	date	uid
1	12/25	1
2	11/21	3
4	12/24	2

Select *
From T1
Where □

T3

id	date	uid
1	12/25	1
2	11/21	3
4	12/24	2

Select id,
Max(date)
From T3
Group By id
Having □

id	date
1	12/25
2	11/21
4	12/24

Select *
From T1
Join T3
On □

T2

oid	val
1	30
1	10
1	10
2	50
2	10

Select *
From T2
Where □

T4

oid	uid
1	30
1	10
1	10
2	50
2	10

Select id,
Max(uid)
From T3
Group By id
Having □

id	uid
1	1
2	3
4	2

Select *
From T2
Join T4
On □

Select oid,
Max(val)
From T4
Group By oid
Having □

T5

oid	MaxVal
1	30
1	10
2	50
2	10

Select *
From T3
Join T5
On □

Out

oid	date	uid	oid	MaxVal
1	12/25	1	1	30
4	12/24	2	2	10

On average, number of
tables generated is **7x** less
v.s. concrete case.

oid	date	uid	oid	MaxVal
1	12/25	1	1	30
4	12/24	2	2	10

Select *
From (Select *
From T1
Where □)
Join (Select id,
Max(val)
From T2
Where □
Group By oid
Having □) T3
On □

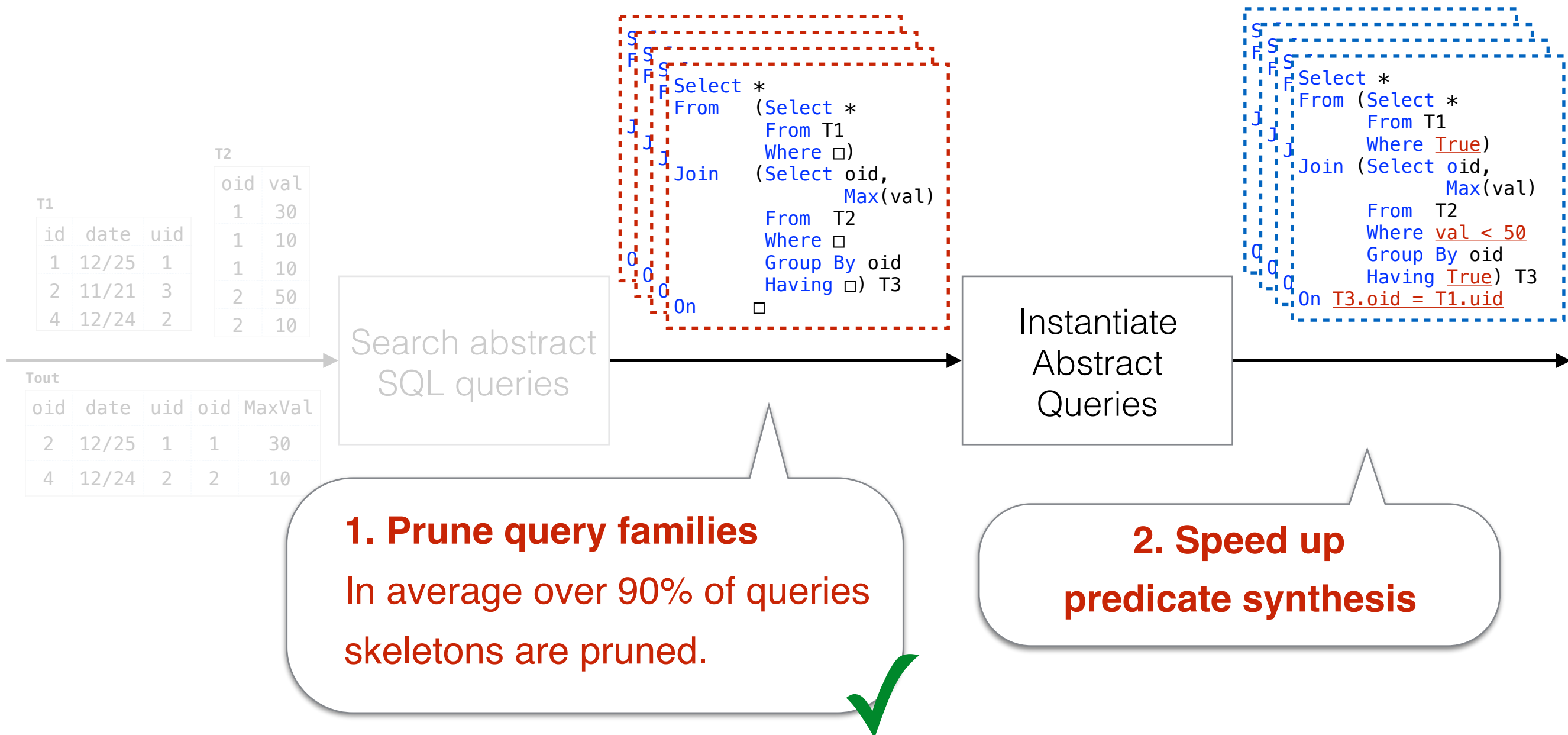
T6

id	date	uid	oid	MaxVal
1	12/25	1	1	30
2	11/21	3	1	30
4	12/24	2	1	30
1	12/25	1	1	10
2	11/21	3	1	10
4	12/24	2	1	10
1	12/25	1	2	50
2	11/21	3	2	50
4	12/24	2	2	50
1	12/25	1	2	10
2	11/21	3	2	10
4	12/24	2	2	10

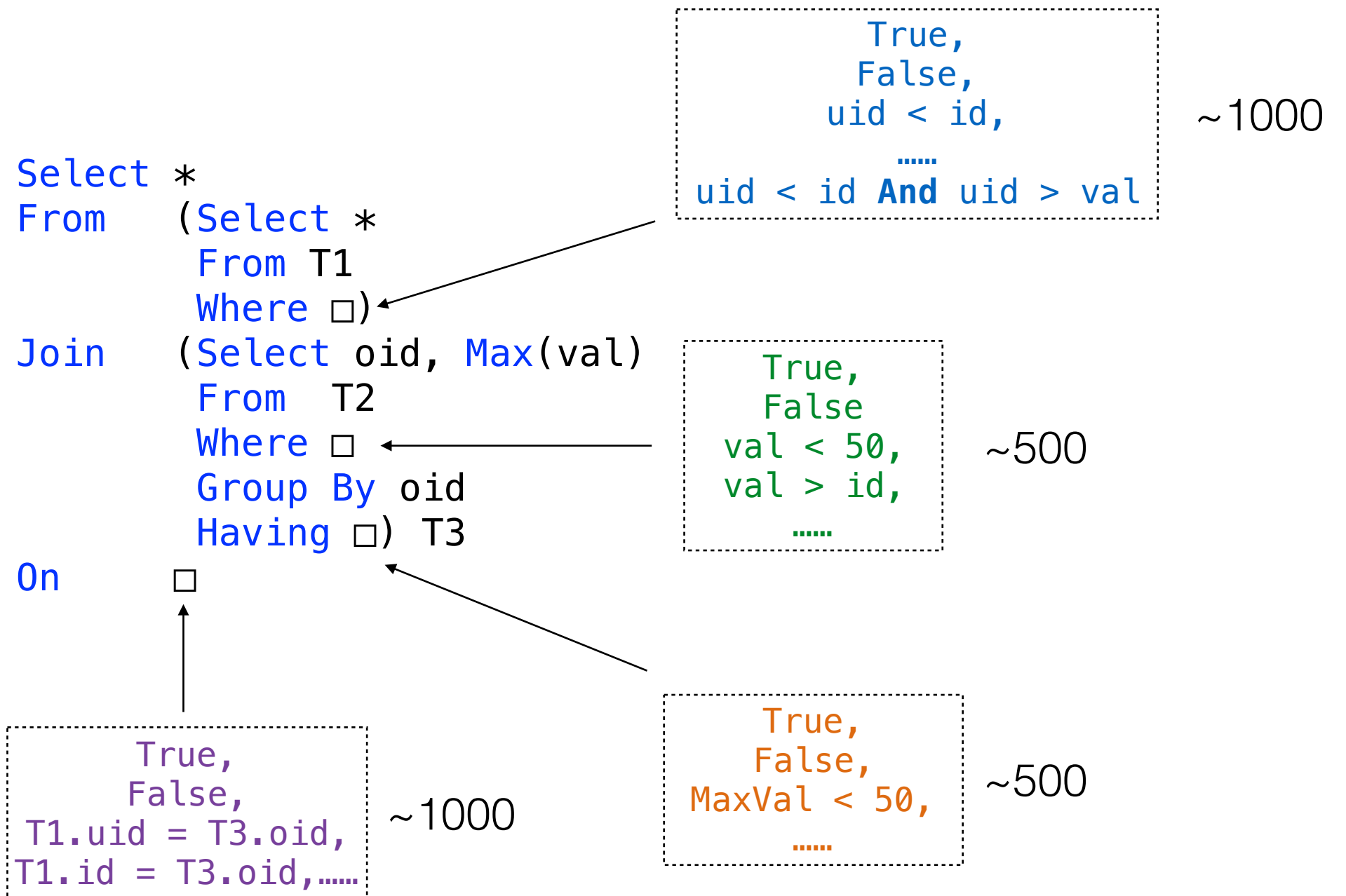
UI

oid	date	uid	oid	MaxVal
1	12/25	1	1	30
4	12/24	2	2	10

Search with Abstract Queries



Predicate Search Space



Challenge: Large number of predicate combinations to search.

Enumerative Predicate Synthesis

```

Select *
From (Select *
      From T1
      Where □) As T3
Join (Select id, Max(val)
      From T2
      Where □
      Group By oid
      Having □) T5
On □
    
```

T1

id	date	uid
1	12/25	1
2	11/21	3
4	12/24	2

```

Select *
From T1
Where True
Where id < uid
Where id = uid
    
```

T3

id	date	uid
1	12/25	1
2	11/21	3
4	12/24	2

```

Select *
From T1
Where □
    
```

```

Select *
From T2
Where □
    
```

```

Select oid,
MAX(val)
From T4
Group By oid
Having □
    
```

```

Select *
From T3
Join T5
On □
    
```

T6

oid	date	uid	oid	MaxVal
1	12/25	1	1	30
4	12/24	2	2	10

T2

oid	val
1	30
1	10
1	10
2	50
2	10

```

Select *
From T2
Where val < 50
Where val ≤ 50
Where oid < val
Where True
    
```

T4

oid	val
1	30
1	10
1	10
2	10
2	10

```

Select oid,
MAX(val)
From T4
Group By oid
Having True
Having oid < mVal
Having mVal < 50
    
```

T5

oid	val
1	30
2	10
2	10
2	10
2	10

```

Select *
From T3
Join T5
On uid = oid
On id < oid
On True
    
```

Out

oid	date	uid	oid	MaxVal
1	12/25	1	1	30
4	12/24	2	2	10

Enumerative Predicate Synthesis

```

Select *
From (Select *
      From T1
      Where □) As T3
Join (Select id, Max(val)
      From T2
      Where □
      Group By oid
      Having □) T5
On □
    
```

T1

id	date	uid
1	12/25	1
2	11/21	3
4	12/24	2

```

Select *
From T1
Where True
    
```

```

Where id < uid
Where id = uid
    
```

T3

id	date	uid
1	12/25	1
2	11/21	3
4	12/24	2

```

Select *
From T1
Where □
    
```

```

Select *
From T2
Where □
    
```

```

Select oid,
MAX(val)
From T4
Group By oid
Having □
    
```

```

Select *
From T3
Join T5
On □
    
```

T6

oid	date	uid	oid	MaxVal
1	12/25	1	1	30
4	12/24	2	2	10

T2

oid	val
1	30
1	10
1	10
2	50
2	10

```

Select *
From T2
Where val < 50
    
```

```

Where val ≤ 50
    
```

```

Where oid < val
    
```

```

Where True
    
```

T4

oid	val
1	30
1	10
1	10
2	10
2	10

```

Select oid,
MAX(val)
From T4
Group By oid
Having True
    
```

```

Having oid < mVal
    
```

```

Having mVal < 50
    
```

T5

oid	val
1	30
2	10
2	10
2	10
2	10

```

Select *
From T3
Join T5
On uid = oid
    
```

```

On id < oid
    
```

```

On True
    
```

Out

oid	date	uid	oid	MaxVal
1	12/25	1	1	30
4	12/24	2	2	10

Enumerative Predicate Synthesis

```

Select *
From (Select *
      From T1
      Where □) As T3
Join (Select id, Max(val)
      From T2
      Where □
      Group By oid
      Having □) T5
On □
    
```

T1

id	date	uid
1	12/25	1
2	11/21	3
4	12/24	2

```

Select *
From T1
Where True
Where id < uid
Where id = uid
    
```

T3

id	date	uid
1	12/25	1
2	11/21	3
4	12/24	2

```

Select *
From T1
Where □
    
```

Inefficient representation

```

Select *
From T2
Where □
    
```

```

Select oid,
MAX(val)
From T4
Group By oid
Having □
    
```

```

Select *
From T3
Join T5
On □
    
```

T6

oid	date	uid	oid	MaxVal
1	12/25	1	1	30
4	12/24	2	2	10

T2

oid	val
1	30
1	10
1	10
2	50
2	10

```

Select *
From T2
Where val < 50
Where val ≤ 50
Where oid < val
Where True
    
```

T4

oid	val
1	30
1	10
1	10
2	10
2	10

```

Select oid,
MAX(val)
From T4
Group By oid
Having True
Having oid < mVal
Having mVal < 50
    
```

T5

oid	val
1	30
2	10
2	10
2	10
2	10

```

Select *
From T3
Join T5
On uid = oid
On id < oid
On True
    
```

Computation overhead

Out

oid	date	uid	oid	MaxVal
1	12/25	1	1	30
4	12/24	2	2	10

Enumerative Predicate Syn

“Evaluating abstract queries into over-approximations of concrete query results.”

```
Select *  
From (Select *  
      From T1  
      Where  $\square$ ) As T3  
Join (Select id, Max(val)  
      From T2  
      Where  $\square$   
      Group By oid  
      Having  $\square$ ) T5  
On  $\square$ 
```

T1

id	date	uid
1	12/25	1
2	11/21	3
4	12/24	2

```
Select *  
From T1  
Where True
```

```
Where id < uid  
Where id = uid
```

```
Select *  
From T1  
Where  $\square$ 
```

id	date	uid
1	12/25	1
2	11/21	3
4	12/24	2

```
Select *  
From T2  
Where  $\square$ 
```

oid	val
1	30
1	10
1	10
2	50
2	10

```
Select oid,  
MAX(val)  
From T4  
Group By oid  
Having  $\square$ 
```

oid	val
1	30
1	10
2	50
2	10

```
Select *  
From T3  
Join T5  
On  $\square$ 
```

id	date	uid	oid	MaxVal
1	12/25	1	1	30
2	11/21	1	1	30
.....				
1	12/25	1	1	10
2	11/21	1	1	10
4	12/24	2	1	10

T2

oid	val
1	30
1	10
1	10
2	50
2	10

```
Select *  
From T2  
Where val < 50
```

```
Where val  $\leq$  50
```

```
Where oid < val
```

```
Where True
```

oid	val
1	30
1	10
1	10
2	10

```
Select oid,  
MAX(val)  
From T4  
Group By oid  
Having True
```

```
Having oid < mVal
```

```
Having mVal < 50
```

T5

oid	val
1	30
2	10
2	10
2	10
2	10

```
Select *  
From T3  
Join T5  
On uid = oid
```

```
On id < oid
```

```
On True
```

Out

oid	date	uid	oid	MaxVal
1	12/25	1	1	30
4	12/24	2	2	10

Enumerative Predicate Synthesis

```

Select *
From (Select *
      From T1
      Where □) As T3
Join (Select id, Max(val)
      From T2
      Where □
      Group By oid
      Having □) T5
On □
    
```

T1

id	date	uid
1	12/25	1
2	11/21	3
4	12/24	2

```

Select *
From T1
Where True
Where id < uid
Where id = uid
    
```

T3

id	date	uid
1	12/25	1
2	11/21	3
4	12/24	2

```

Select *
From T1
Where □
    
```

```

Select *
From T2
Where □
    
```

```

Select oid,
MAX(val)
From T4
Group By oid
Having □
    
```

```

Select *
From T3
Join T5
On □
    
```

T6

oid	date	uid	oid	MaxVal
1	12/25	1	1	30
4	12/24	2	2	10

T2

oid	val
1	30
1	10
1	10
2	50
2	10

```

Select *
From T2
Where val < 50
Where val ≤ 50
Where oid < val
Where True
    
```

T4

oid	val
1	30
1	10
1	10
2	10
2	10

```

Select oid,
MAX(val)
From T4
Group By oid
Having True
Having oid < mVal
Having mVal < 50
    
```

T5

oid	val
1	30
2	10
2	10
2	10
2	10

```

Select *
From T3
Join T5
On uid = oid
On id < oid
On True
    
```

Out

oid	date	uid	oid	MaxVal
1	12/25	1	1	30
4	12/24	2	2	10

Encoding Tables using Bit-vectors

```

Select *
From (Select *
      From T1
      Where □) As T3
Join (Select id, Max(val)
      From T2
      Where □
      Group By oid
      Having □) T5
On □
    
```

T1

id	date	uid
1	12/25	1
2	11/21	3
4	12/24	2

```

Select *
From T1
Where True
Where id < uid
Where id = uid
    
```

1
1
1
1
0

```

Select *
From T1
Where □
    
```

id	date	uid
1	12/25	1
2	11/21	3
4	12/24	2

```

Select *
From T2
Where □
    
```

oid	val
1	30
1	10
1	10
2	50
2	10

```

Select oid,
MAX(val)
From T4
Group By oid
Having □
    
```

oid	val
1	30
1	10
2	50
2	10

```

Select *
From T3
Join T5
On □
    
```

id	date	uid	oid	MaxVal
1	12/25	1	1	30
2	11/21	1	1	30
.....				
1	12/25	1	1	10
2	11/21	1	1	10
4	12/24	2	1	10

T2

oid	val
1	30
1	10
1	10
2	50
2	10

```

Select *
From T2
Where val < 50
Where val ≤ 50
Where oid < val
Where True
    
```

1
1
1
0
1
1

```

Select oid,
MAX(val)
From T4
Group By oid
Having True
Having oid < mVal
Having mVal < 50
    
```

1
0
0
0
1
1

```

Select *
From T3
Join T5
On uid = oid
On id < oid
On True
    
```

1
0
.
.
.
1
1

Computation overhead

Out

oid	date	uid	oid	MaxVal
1	12/25	1	1	30
4	12/24	2	2	10

Optimize computation: Grouping Predicates

oid	val
1	30
1	10
1	10
2	50
2	10

oid	mVal
1	30
1	10
2	50
2	10

Select oid,
MAX(val)
From T4
Group By oid
Having □

Alternative inputs
from its subquery

1
1
1
0
1

1
1
1
0
1

Select oid,
MAX(val)
From T4
Group By oid
Having True
Having oid < mVal
Having mVal <= 50

Having True
Having oid < mVal

All possible outputs
of this query

1
0
0
1

1
0
1
0

1
0
0
0

Problem: need to
perform 2×3 operations
to get only 3 results.

Discovery: Grouping
predicates on the
summary table.

Number of predicates
reduced by 40,000×

Enumerative Predicate Synthesis

```
Select *  
From (Select *  
      From T1  
      Where  $\square$ ) As T3  
Join (Select id, Max(val)  
      From T2  
      Where  $\square$   
      Group By oid  
      Having  $\square$ ) T5  
On  $\square$ 
```

T1

id	date	uid
1	12/25	1
2	11/21	3
4	12/24	2

```
Select *  
From T1  
Where True  
Where id < uid  
Where id = uid
```

T3

id	date	uid
1	12/25	1
2	11/21	3
4	12/24	2
4	12/24	2
4	12/24	2

```
Select *  
From T1  
Where  $\square$ 
```

```
Select *  
From T2  
Where  $\square$ 
```

```
Select oid,  
MAX(val)  
From T4  
Group By oid  
Having  $\square$ 
```

```
Select *  
From T3  
Join T5  
On  $\square$ 
```

T6

oid	date	uid	oid	MaxVal
1	12/25	1	1	30
4	12/24	2	2	10
4	12/24	2	2	10
4	12/24	2	2	10

T2

oid	val
1	30
1	10
1	10
2	50
2	10

```
Select *  
From T2  
Where val < 50  
Where val ≤ 50  
Where oid < val  
Where True
```

T4

oid	val
1	30
1	10
1	10
2	10
2	10

```
Select oid,  
MAX(val)  
From T4  
Group By oid  
Having True  
Having oid < mVal  
Having mVal < 50
```

T5

oid	val
1	30
2	10
2	10
2	10
2	10
2	10

```
Select *  
From T3  
Join T5  
On uid = oid  
On id < oid  
On True
```

Out

oid	date	uid	oid	MaxVal
1	12/25	1	1	30
4	12/24	2	2	10

Grouping Predicates + Bit-vector Representation

```

Select *
From (Select *
      From T1
      Where □) As T3
Join (Select id, Max(val)
      From T2
      Where □
      Group By oid
      Having □) T5
On □
    
```

T1

id	date	uid
1	12/25	1
2	11/21	3
4	12/24	2

```

Select *
From T1
Where True
Where id < uid
Where id = uid
    
```

1
1
1
1
0

```

Select *
From T1
Where □
    
```

id	date	uid
1	12/25	1
2	11/21	3
4	12/24	2

```

Select *
From T2
Where □
    
```

oid	val
1	30
1	10
1	10
2	50
2	10

```

Select oid,
MAX(val)
From T4
Group By oid
Having □
    
```

oid	val
1	30
1	10
2	50
2	10

```

Select *
From T3
Join T5
On □
    
```

id	date	uid	oid	MaxVal
1	12/25	1	1	30
2	11/21	1	1	30
.....				
1	12/25	1	1	10
2	11/21	1	1	10
4	12/24	2	1	10

T2

oid	val
1	30
1	10
1	10
2	50
2	10

```

Select *
From T2
Where val < 50
Where True
    
```

1
1
1
0
1
1

```

Select oid,
MAX(val)
From T4
Group By oid
Having True
Having mVal < 50
    
```

1
0
0
0
1
1

```

Select *
From T3
Join T5
On uid = oid
On id < oid
On True
    
```

1
0
.
.
.
1
1

Out

oid	date	uid	oid	MaxVal
1	12/25	1	1	30
4	12/24	2	2	10

As a Programming-by-Example System

- **Synthesis process**
 - Iterating over the search depth for abstract queries
 - Instantiate abstract queries in the current depth and check results
- **Dealing with ambiguity**
 - Ranking programs by heuristic
 - complexity, naturalness, constant coverage
 - Provide a new example / restrict aggregation functions.

Implementation — Scythe

<http://scythe.cs.washington.edu>

- Supported features:
 - Select, Join, Group By, Aggregation,
 - Subqueries, Outer Join, Exists, Union
- Unsupported
 - Arithmetics, Pivot, Window functions, Limit, Insert

Evaluation

- Benchmarks from Stack Overflow:

- 57 used in development
- 57 top-voted posts
- 51 recent posts

- Benchmarks from prior work:

[Zhang et al. ASE'13]

- 23 textbook questions.
- 5 forum posts.

In total 193 benchmarks.
Avg. Example Size: 34 cells

- Algorithms

- Enumerative Search

[Udupa et al. PLDI'13]

- SqlSynthesizer

(Decision tree algorithm)

[Zhang et al. ASE'13]

- Scythe

- Evaluation Condition

- 4G memory, 600s timeout

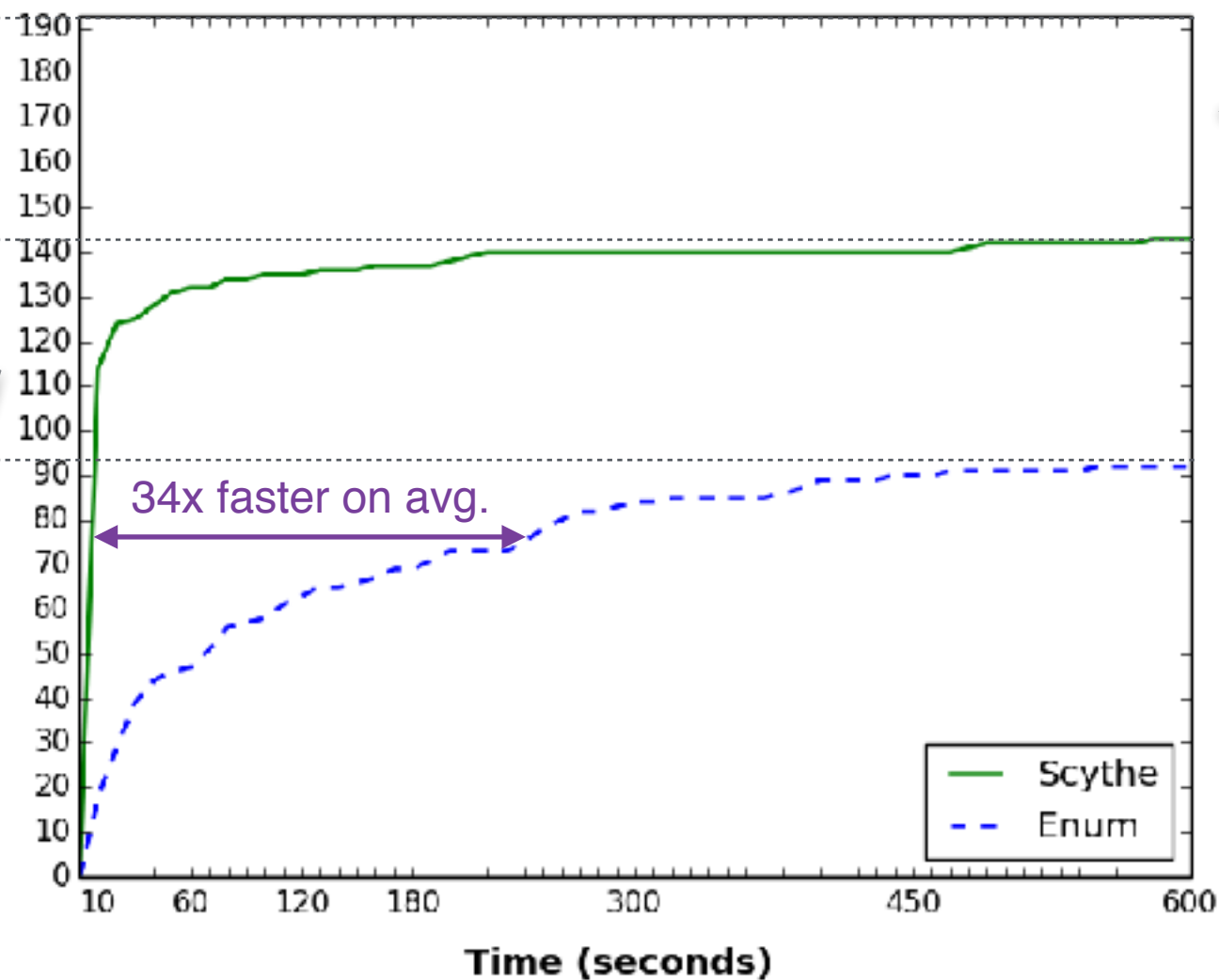
Evaluation

Benchmark: 193

Scythe: 143

Enum: 92

59% can be answered
within 10 seconds



Reasons for failures

34: missing features
15: timeout
1: failed to disambiguate

Comparing with SQLSynthesizer

Scythe:
18/28 in 120s
SQLSynthesizer:
15/28 in 120s

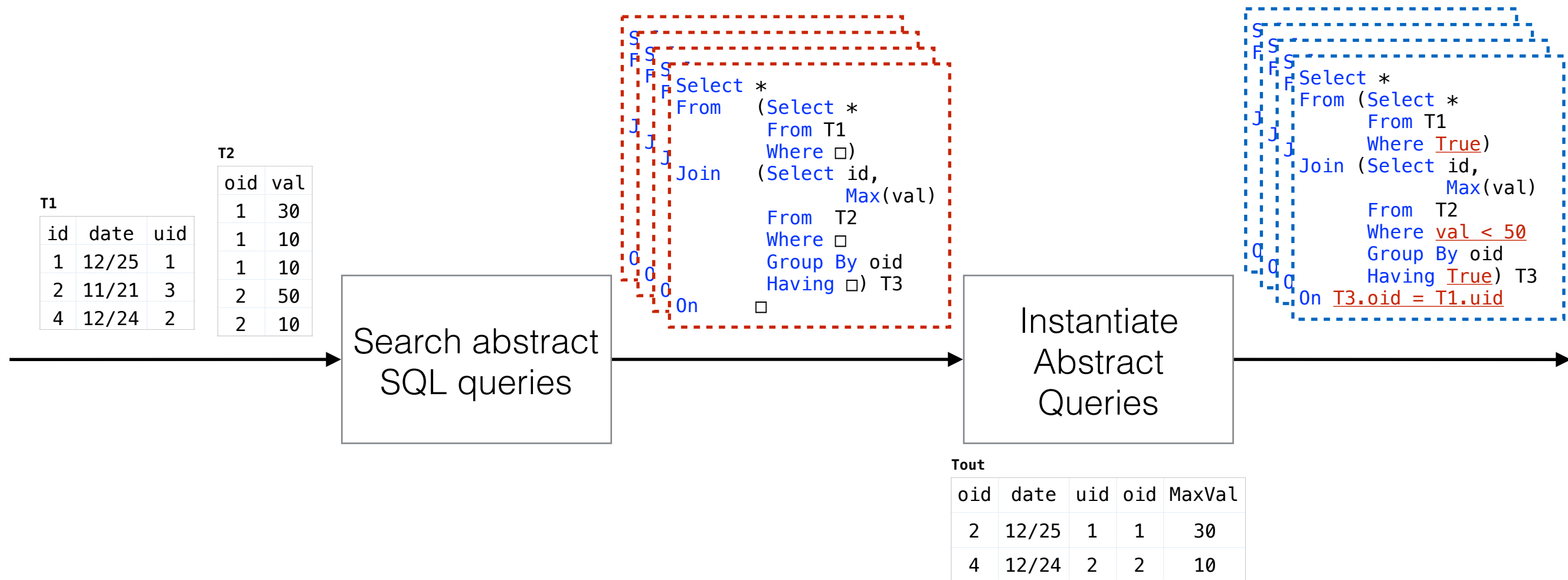
Some Related Work

- **Enumerative search**
 - Value-based Memoization [Udupa et al. PLDI'13]
- **Search optimization with approximation**
 - Synthesizing regex from examples [Lee et al. GPCE'16]
 - Monotonicity [Hu et al. PLDI'17]
- **Synthesizing table manipulation programs**
 - Pruning search space using partial programs [Feng et al. PLDI'17]

	Pruning Approach	Pruning Overhead	Pruning Power
Scythe	Over-approximation	Higher	Higher
Feng et al.	Constraint encoded properties	Lower	Lower

Benefit from value-based search space compression.

Algorithm: Decompose Search Process With Abstract Queries



Try demo on <http://scythe.cs.washington.edu!>