

# Topic 1: Data models and query languages

## Unit 1: SQL

### Lecture 1

Wolfgang Gatterbauer

CS7240 Principles of scalable data management (sp24)

<https://northeastern-datalab.github.io/cs7240/sp24/>

1/9/2024

## Topic 1: Data Models and Query Languages

- **Lecture 1 (Tue 1/19):** Course introduction / T1-U1 SQL / PostgreSQL setup / SQL Activities
- **Lecture 2 (Fri 1/12):** T1-U1 SQL
- **Lecture 3 (Tue 1/16):** T1-U1 SQL
- **Lecture 4 (Fri 1/19):** T1-U2 Logic & Relational Calculus
- **Lecture 5 (Tue 1/23):** T1-U2 Logic & Relational Calculus
- **Lecture 6 (Fri 1/26):** T1-U2 Logic & Relational Calculus
- **Lecture 7 (Tue 1/30):** T1-U3 Relational Algebra & Codd's Theorem
- **Lecture 8 (Fri 2/2):** T1-U3 Relational Algebra & Codd's Theorem
- **Lecture 9 (Tue 2/6):** T1-U4 Datalog & Recursion & ASP
- **Lecture 10 (Tue 2/9):** T1-U4 Datalog & Recursion & ASP
- **Lecture 11 (Tue 2/13):** T1-U4 Datalog & Recursion & ASP
- **Lecture 12 (Fri 2/16):** T1-U4 Datalog & Recursion & ASP

Pointers to relevant concepts & supplementary material:

- **Unit 1. SQL:** [SAMS'12], [CS 3200], [Cow'03] Ch3 & Ch5, [Complete'08] Ch6, [Silberschatz+'20] Ch3.8
- **Unit 2. Logic & Relational Calculus:** First-Order Logic (FOL), relational calculus (RC): [Barland+'08] 4.1.2 & 4.2.1 & 4.4, [Genesereth+] Ch6, [Halpern+'01], [Cow'03] Ch4.3 & 4.4, [Elmasri, Navathe'15] Ch8.6 & Ch8.7, [Silberschatz+'20] Ch27.1 & Ch27.2, [Alice'95] Ch3.1-3.3 & Ch4.2 & Ch4.4 & Ch5.3-5.4, [Barker-Plummer+'11] Ch11
- **Unit 3. Relational Algebra & Codd's Theorem:** Relational Algebra (RA), Codd's theorem: [Cow'03] Ch4.2, [Complete'08] Ch2.4 & Ch5.1-5.2, [Elmasri, Navathe'15] Ch8, [Silberschatz+'20] Ch2.6, [Alice'95] Ch4.4 & Ch5.4
- **Unit 4. Datalog & Recursion & Answer Set Programming:** Datalog, recursion, Stratified Datalog with negation, Datalog evaluation strategies, Stable Model semantics, Answer Set Programming (ASP): [Complete'08] Ch5.3, [Cow'03] Ch 24, [x755'19], [Soufflé], [DLV], [G., Suciu'10], [Eiter+'09]
- **(Unit 5. Alternative Data Models: NoSQL:** [Hellerstein, Stonebraker'05], [Sadalage, Fowler'12], [Harrison'16])

# Outline: T1-U1: SQL

- SQL

- Schema, keys, referential integrity
- Joins
- Aggregates and grouping
- Nested queries (Subqueries)
- Union and Theta Joins
- Nulls & Outer joins
- Top-k
- [Recursion: moved to T1-U4: Datalog]

# Top Programming Languages 2022 > Python's still No. 1, but employers love to see SQL skills

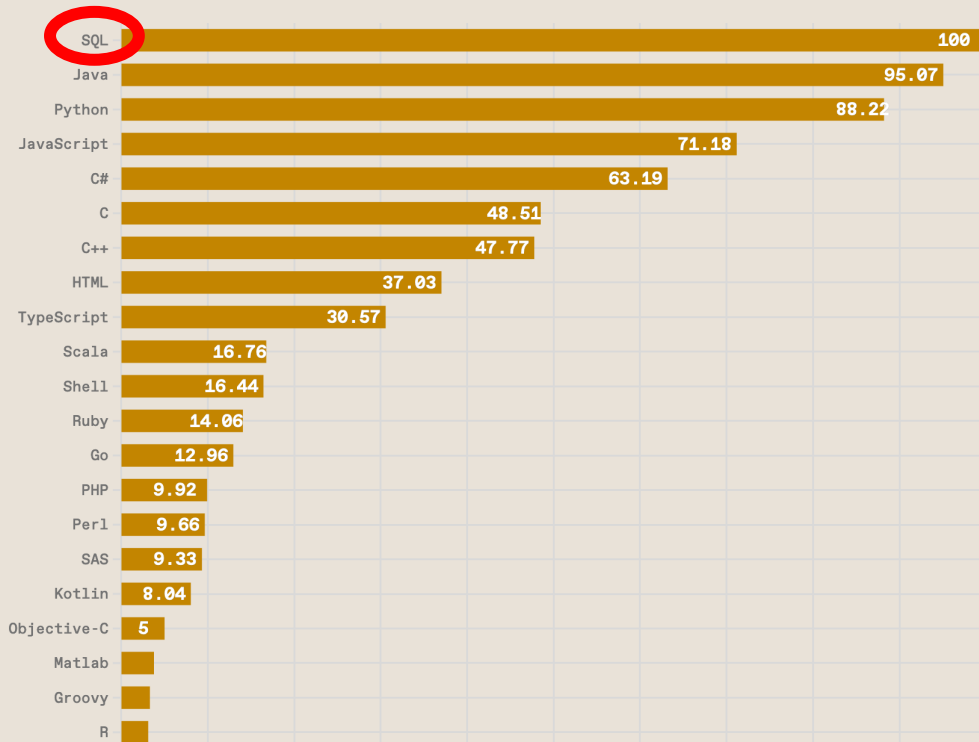
BY STEPHEN CASS | 23 AUG 2022 | 4 MIN READ |

## IEEE Spectrum's Top Programming Languages 2022

### Top Programming Languages 2022

Click a button to see a differently weighted ranking

Spectrum **Jobs** Trending



But among these stalwarts is the rising popularity of SQL. In fact, it's at No. 1 in our Jobs ranking, which looks solely at metrics from the IEEE Job Site and CareerBuilder. Having looked through literally hundreds and hundreds of job listings in the course of compiling these rankings for you, dear reader, I can say that the strength of the SQL signal is not because there are a lot of employers looking for *just* SQL coders, in the way that they advertise for Java experts or C++ developers. They want a given language *plus* SQL. And lots of them want that “plus SQL.”



**It may not be the most glamorous language...but some experience with SQL is a valuable arrow to have in your quiver.**

Source: <https://spectrum.ieee.org/top-programming-languages-2022>

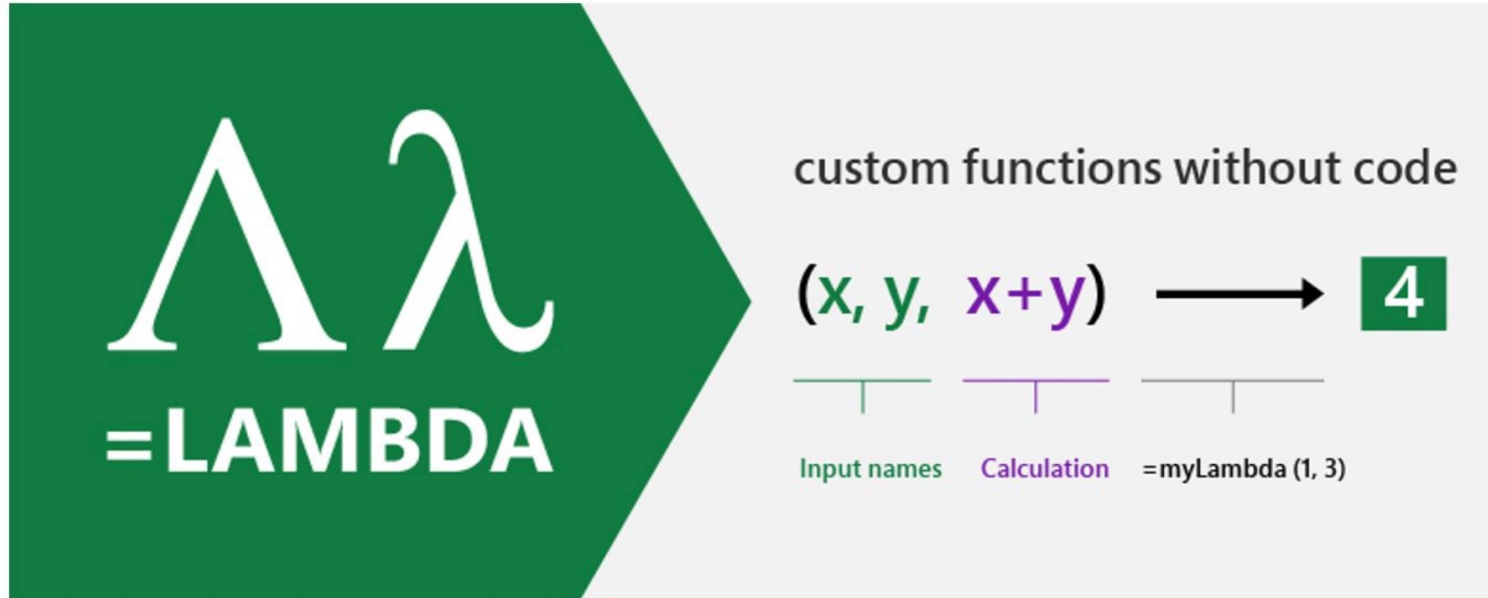
Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>



Fun question: What is the most popular PL?



# Fun question: What is the most popular PL?



Possibly interesting  
class scribe: Why is Excel  
Turing-complete?

Ever since it was released in the 1980s, Microsoft Excel has changed how people organize, analyze, and visualize their data, providing a basis for decision-making for the millions of people who use it each day. It's also the world's most widely used programming language. Excel formulas are written by an order of magnitude more users than all the C, C++, C#, Java, and Python programmers in the world combined. Despite its success, considered as a *programming language* Excel has fundamental weaknesses. Over the years, two particular shortcomings have stood out: (1) the Excel formula language really only supported scalar values—numbers, strings, and Booleans—and (2) it didn't let users define new functions.

Until now.

Source: <https://www.microsoft.com/en-us/research/blog/lambda-the-ultimaet-excel-worksheet-function/>

Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>

# Structured Query Language: SQL

- Influenced by relational calculus (= First Order Logic)
- SQL is a **declarative** query language
  - We say what we want to get
  - We don't say how we should get it ("**separation of concerns**")

# SQL: Declarative Programming

SQL

```
select (e.salary / (e.age - 18)) as comp
from employee as e
where e.name = "Jones"
```

Declarative Language: you say what you want without having to say how to do it.

Procedural Language: you have to specify exact steps to get the result.

# SQL: was not the only Attempt

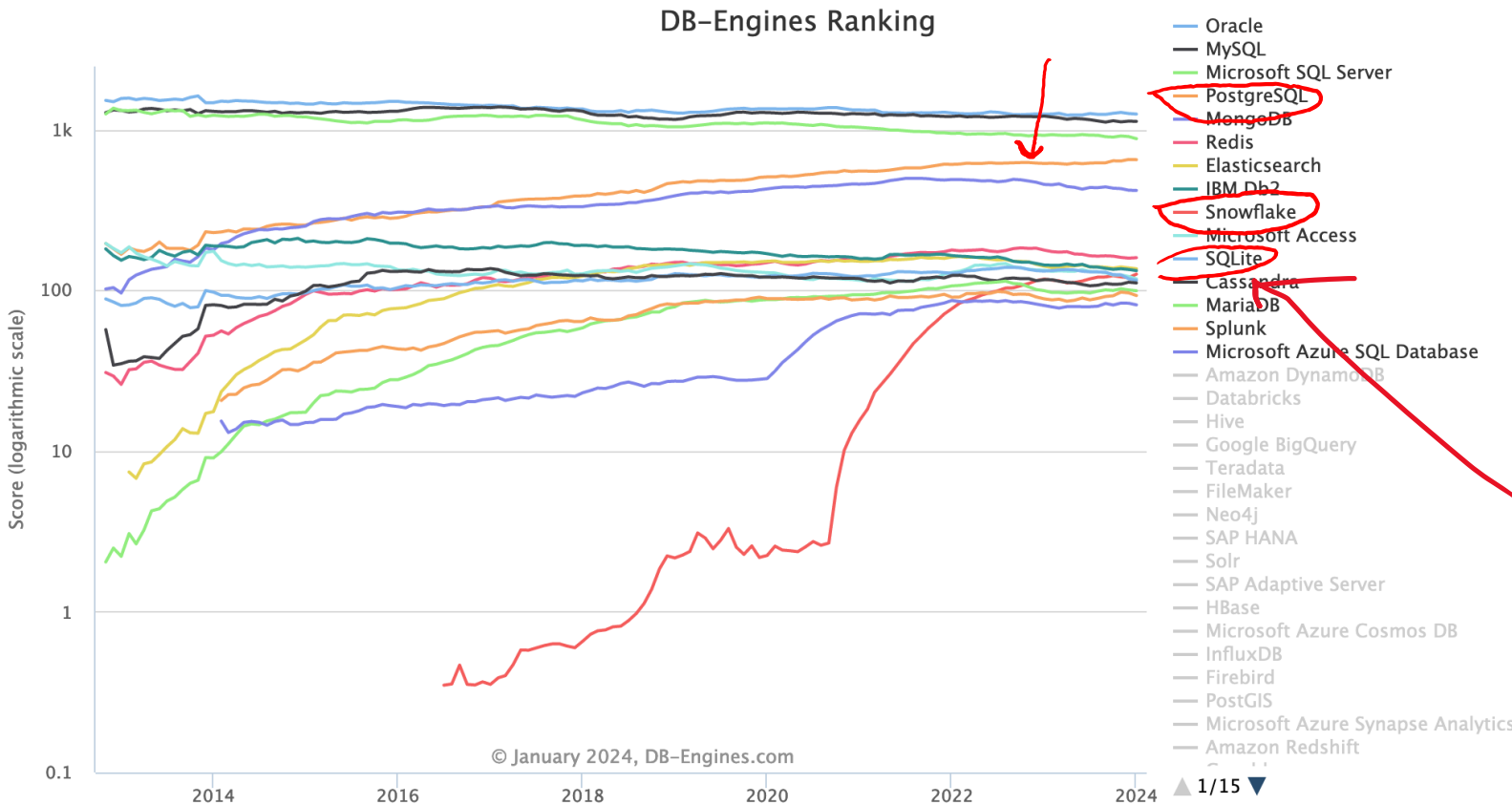
*reading order:*

```
SQL 3 select (e.salary / (e.age - 18)) as comp  
1 from employee as e  
2 where e.name = "Jones"
```

```
QUEL 1 range of e is employee  
3 retrieve (comp = e.salary / (e.age - 18))  
2 where e.name = "Jones"
```

*Commercially not used anymore since ~1980*

# Why PostgreSQL instead of MariaDB (or MySQL)



## Method of calculating the scores of the DB-Engines Ranking

The DB-Engines Ranking is a list of database management systems ranked by their current popularity. We measure the popularity of a system by using the following parameters:

- **Number of mentions of the system on websites**, measured as number of results in search engines queries. At the moment, we use [Google](#) and [Bing](#) for this measurement. In order to count only relevant results, we are searching for <system name> together with the term database, e.g. "Oracle" and "database".
- **General interest in the system**. For this measurement, we use the frequency of searches in [Google Trends](#).
- **Frequency of technical discussions about the system**. We use the number of related questions and the number of interested users on the well-known IT-related Q&A sites [Stack Overflow](#) and [DBA Stack Exchange](#).
- **Number of job offers, in which the system is mentioned**. We use the number of offers on the leading job search engines [Indeed](#) and [Simply Hired](#).
- **Number of profiles in professional networks, in which the system is mentioned**. We use the internationally most popular professional network [LinkedIn](#).
- **Relevance in social networks**. We count the number of [Twitter \(X\)](#) tweets, in which the system is mentioned.

The DB-Engines Ranking does not measure the number of installations of the systems, or their use within IT systems. It can be expected, that an increase of the popularity of a system as measured by the DB-Engines Ranking (e.g. in discussions or job offers) precedes a corresponding broad use of the system by a certain time factor. Because of this, the DB-Engines Ranking can act as an early indicator.

SQLite likely has the most number of installations: its is an embedded serverless database (not a server-client databas)



# Why PostgreSQL instead of MariaDB (or MySQL)



Although PostgreSQL has been around for a while, the relative **decline of MySQL** has made it a serious contender for the title of most used open source database. Since it works very similarly to MySQL, developers who prefer open source software are converting in droves.

## Advantages

- By far, PostgreSQL's most mentioned advantage is the efficiency of its central algorithm, which means it outperforms many databases that are advertised as more advanced. This is especially useful if you are working with large datasets, for which I/O processes can otherwise become a bottleneck.
- It is also one of the most flexible open source databases around; you can write functions in a wide range of server-side languages: Python, Perl, Java, Ruby, C, and R.
- As one of the most commonly used open source databases, PostgreSQL's community support is some of the best around.

*I also prefer PostgreSQL over MySQL because it has a more principled interpretation of SQL (and a powerful EXPLAIN command)*

Source: [https://db-engines.com/en/ranking\\_trend](https://db-engines.com/en/ranking_trend)

Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>

The screenshot shows a web browser displaying a Guardian article. The URL is <https://www.theguardian.com/info/2018/nov/30/bye-bye-mongo-hello-postgres>. The page has a dark blue header with the Guardian logo and navigation links for News, Opinion, Sport, Culture, Lifestyle, and More. The article title is 'Bye bye Mongo, Hello Postgres' and the sub-headline is 'In April the Guardian switched off the Mongo DB cluster used to store our content after completing a migration to PostgreSQL on Amazon RDS. This post covers why and how'. The author is listed as Philip McMahon, Maria-Livia Chiorean, Susie Coleman and Akash Askoolum. The article is dated Fri 30 Nov 2018 05:36 EST and has 438 likes. A photograph of an elephant holding a large bundle of greenery in its trunk is featured below the text.

▲ An elephant picking up some greenery. Photograph: Michael Schwarz/AP

Source: <https://www.theguardian.com/info/2018/nov/30/bye-bye-mongo-hello-postgres>

# Simple SQL Query

*Our friend here shows that you can follow along in Postgres. Just install the database from the text file "302 - ..." available in our sql folder from our course web page*

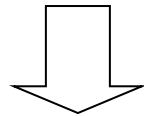


302

## Product

PName	Price	Category	Manufacturer
Gizmo	\$19.99	Gadgets	GizmoWorks
Powergizmo	\$29.99	Gadgets	GizmoWorks
SingleTouch	\$149.99	Photography	Canon
MultiTouch	\$203.99	Household	Hitachi

```
SELECT pName, price
FROM Product
WHERE price > 100
```



# Simple SQL Query

Our friend here shows that you can follow along in Postgres. Just install the database from the text file "302 - ..." available in our sql folder from our course web page

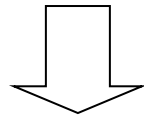


302

## Product

PName	Price	Category	Manufacturer
Gizmo	\$19.99	Gadgets	GizmoWorks
Powergizmo	\$29.99	Gadgets	GizmoWorks
SingleTouch	\$149.99	Photography	Canon
MultiTouch	\$203.99	Household	Hitachi

```
3 SELECT pName, price  
1 FROM Product  
2 WHERE price > 100
```



PName	Price
SingleTouch	\$149.99
MultiTouch	\$203.99

Selection  
& Projection

# How to install PostgreSQL?

As always: if you find something that does not work, PLEASE let me know to fix it!

## Topic 1: Data Models and Query Languages

- **Lecture 1 (Tue 1/19):** Course introduction / T1-U1 SQL / PostgreSQL setup / SQL Activities

The grid contains 16 numbered screenshots from a tutorial document. The first row shows the document's overview and outline. The second row shows the PostgreSQL app being downloaded and installed. The third row shows the app icon on the desktop and in the dock. The fourth row shows the app icon in Spotlight search. The fifth row shows the app icon in Spotlight search. The sixth row shows the app icon in Spotlight search. The seventh row shows the app icon in Spotlight search. The eighth row shows the app icon in Spotlight search. The ninth row shows the app icon in Spotlight search. The tenth row shows the app icon in Spotlight search. The eleventh row shows the app icon in Spotlight search. The twelfth row shows the app icon in Spotlight search. The thirteenth row shows the app icon in Spotlight search. The fourteenth row shows the app icon in Spotlight search. The fifteenth row shows the app icon in Spotlight search. The sixteenth row shows the app icon in Spotlight search.

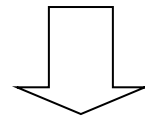
The screenshot shows a GitHub repository page for 'northeastern-dat... / cs3200-activiti...'. The page has a navigation bar with 'Code', 'Issues', 'Pull requests', 'Actions', 'Projects', 'Wiki', 'Security', 'Insights', and 'Settings'. Below the navigation bar is a 'Files' section with a search bar and a file browser. The file browser shows a directory structure with folders 'clingo', 'dlv', 'jupyter', and 'souffle', and a subdirectory 'sql'. The 'sql' directory is expanded, showing a list of files with their names, last commit messages, and last commit dates. The files listed are: '300-SmallIMDB.txt', '302-Simpleproducts.txt', '304-Worker.txt', '305-Conceptualevaluationst...', '306-NestedLoopJoin.py', '308-Purchase.txt', '312-Personandaddress.txt', '315-OtherProducts.txt', '316-OtherProducts.txt', and '317-OtherProducts.txt'. The repository is owned by 'wolfandthegang' and has 27124b7 as the last commit.

# Selection vs. Projection

## Product

PName	Price	Category	Manufacturer
Gizmo	\$19.99	Gadgets	GizmoWorks
Powergizmo	\$29.99	Gadgets	GizmoWorks
SingleTouch	\$149.99	Photography	Canon
MultiTouch	\$203.99	Household	Hitachi

```
SELECT pName, price
FROM Product
WHERE price > 100
```



PName	Price
SingleTouch	\$149.99
MultiTouch	\$203.99

*Selection  
& Projection*

*Where does the  
selection happen?*

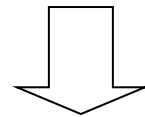


# Selection vs. Projection

**Product**

PName	Price	Category	Manufacturer
Gizmo	\$19.99	Gadgets	GizmoWorks
Powergizmo	\$29.99	Gadgets	GizmoWorks
SingleTouch	\$149.99	Photography	Canon
MultiTouch	\$203.99	Household	Hitachi

```
SELECT pName, price
FROM Product
WHERE price > 100
```



PName	Price
SingleTouch	\$149.99
MultiTouch	\$203.99

One **selects** certain  
entires=tuples (rows)  
-> happens in the  
**WHERE** clause  
-> acts like a **filter**



# Selection vs. Projection

Product			
PName	Price	Category	Manufacturer
Gizmo	\$19.99	Gadgets	GizmoWorks
Powergizmo	\$29.99	Gadgets	GizmoWorks
SingleTouch	\$149.99	Photography	Canon
MultiTouch	\$203.99	Household	Hitachi

One **projects** onto some attributes (columns)  
-> happens in the **SELECT** clause

```
SELECT pName, price
FROM Product
WHERE price > 100
```

One **selects** certain entires=tuples (rows)  
-> happens in the **WHERE** clause  
-> acts like a **filter**

PName	Price
SingleTouch	\$149.99
MultiTouch	\$203.99

# Eliminating Duplicates

## Product

PName	Price	Category	Manufacturer
Gizmo	\$19.99	Gadgets	GizmoWorks
Powergizmo	\$29.99	Gadgets	GizmoWorks
SingleTouch	\$149.99	Photography	Canon
MultiTouch	\$203.99	Household	Hitachi

```
SELECT category  
FROM Product
```



# Eliminating Duplicates

## Product

PName	Price	Category	Manufacturer
Gizmo	\$19.99	Gadgets	GizmoWorks
Powergizmo	\$29.99	Gadgets	GizmoWorks
SingleTouch	\$149.99	Photography	Canon
MultiTouch	\$203.99	Household	Hitachi

```
SELECT category  
FROM Product
```



Category
Gadgets
Gadgets
Photography
Household



Category
Gadgets
Photography
Household

*Set vs. Bag semantics*

*Think of a dictionary:  
keys mapping to  
# of occurrences*

*Gadgets : 2  
Photography : 1  
Household : 1*

*underlying set also  
called the "support"  
of the bag*

# Eliminating Duplicates

## Product

PName	Price	Category	Manufacturer
Gizmo	\$19.99	Gadgets	GizmoWorks
Powergizmo	\$29.99	Gadgets	GizmoWorks
SingleTouch	\$149.99	Photography	Canon
MultiTouch	\$203.99	Household	Hitachi

```
SELECT category  
FROM Product
```



Category
Gadgets
Gadgets
Photography
Household

```
SELECT DISTINCT category  
FROM Product
```



Category
Gadgets
Photography
Household

*Set vs. Bag  
semantics*

*Think of a  
dictionary:  
keys mapping to  
# of occurrences*

*Gadgets : 2  
Photography :  
Household : 1*

*underlying set also  
called the "support"  
of the bag*

# Outline: T1-U1: SQL

- SQL

- Schema, keys, referential integrity
- Joins
- Aggregates and grouping
- Nested queries (Subqueries)
- Union and Theta Joins
- Nulls & Outer joins
- Top-k
- [Recursion: moved to T1-U4: Datalog]

## Product

<u>PName</u>	Price	Category	Manufacturer
Gizmo	\$19.99	Gadgets	GizmoWorks
Powergizmo	\$29.99	Gadgets	GizmoWorks
SingleTouch	\$149.99	Photography	Canon
MultiTouch	\$203.99	Household	Hitachi

## Company

<u>CName</u>	StockPrice	Country
GizmoWorks	25	USA
Canon	65	Japan
Hitachi	15	Japan

What is here  
a key vs.  
a foreign key?





# Keys and Foreign Keys

## Product

<u>PName</u>	Price	Category	Manufacturer
Gizmo	\$19.99	Gadgets	GizmoWorks
Powergizmo	\$29.99	Gadgets	GizmoWorks
SingleTouch	\$149.99	Photography	Canon
MultiTouch	\$203.99	Household	Hitachi

Foreign key

## Company

<u>CName</u>	StockPrice	Country
GizmoWorks	25	USA
Canon	65	Japan
Hitachi	15	Japan

Keys

Keys and foreign keys  
are special cases of  
more general  
constraints. Which? ?

# Keys and Foreign Keys



In the following,  $R(U)$  denotes the schema of a relation with name  $R$  and set of attributes  $U$ .

## Functional Dependencies

A *functional dependency (FD)* on relations of schema  $R(U)$  is an expression of the form

$$R : X \rightarrow Y, \quad (1)$$

where  $X \subseteq U$  and  $Y \subseteq U$  are subsets of  $R$ 's attributes. Instance  $r$  of schema  $R(U)$  is said to *satisfy* FD  $fd$ , denoted  $r \models fd$ , if whenever tuples  $t_1 \in r$  and  $t_2 \in r$  agree on all attributes in  $X$ , they also agree on all attributes in  $Y$ :

$$r \models fd \iff \text{for every } t_1, t_2 \in r \text{ if } \pi_X(t_1) = \pi_X(t_2) \text{ then } \pi_Y(t_1) = \pi_Y(t_2)$$

Here,  $\pi_X(t)$  denotes the projection of tuple  $t$  on the attributes in  $X$ .

## Key Dependencies

In the particular case when  $Y=U$ , a functional dependency of form (1) is called a *key dependency*, and the set of attributes  $X$  is called a *key* for  $R$ .

## Product

<u>PName</u>	Price	Category	Manufacturer
Gizmo	\$19.99	Gadgets	GizmoWorks
Powergizmo	\$29.99	Gadgets	GizmoWorks
SingleTouch	\$149.99	Photography	Canon
MultiTouch	\$203.99	Household	Hitachi

## Company

<u>CName</u>	StockPrice	Country
GizmoWorks	25	USA
Canon	65	Japan
Hitachi	15	Japan

## Inclusion Dependencies

Functional and join dependencies and their special-case subclasses each pertain to single relations. The following class of dependencies can express connections between relations. An *inclusion dependency (IND)* on pairs of relations of schemas  $R(U)$  and  $S(V)$  (with  $R$  and  $S$  not necessarily distinct) is an expression of the form

$$R[X] \subseteq S[Z], \quad (4)$$

where  $X \subseteq U$  and  $Z \subseteq V$ . Inclusion dependencies are also known as *referential constraints*. Relations  $r$  and  $s$  of schemas  $R(U)$ , respectively  $S(V)$  satisfy inclusion dependency  $id$ , denoted  $r, s \models id$ , if the projection of  $r$  on  $X$  is included in the projection of  $s$  on  $Z$ :

$$r, s \models id \iff \Pi_X(r) \subseteq \Pi_Z(s).$$

When  $R$  and  $S$  refer to the same relation name, then  $r = s$  in the above definition of satisfaction.

## Foreign key

### Foreign Key Dependencies

In the particular case when  $Z$  is a key for relations of schema  $S$  ( $S: Z \rightarrow V$ ), INDs of form (4) are called *foreign key dependencies*. Intuitively, in this case the projection on  $X$  of every tuple  $t$  in  $r$  contains the key of a tuple from the "foreign" table  $s$ .

**Keys** →

# Keys and Foreign Keys

In the following,  $R(U)$  denotes the schema of a relation with name  $R$  and set of attributes  $U$ .

## Functional Dependencies

A *functional dependency (FD)* on relations of schema  $R(U)$  is an expression of the form

$$R : X \rightarrow Y, \quad (1)$$

where  $X \subseteq U$  and  $Y \subseteq U$  are subsets of  $R$ 's attributes. Instance  $r$  of schema  $R(U)$  is said to *satisfy* FD  $fd$ , denoted  $r \models fd$ , if whenever tuples  $t_1 \in r$  and  $t_2 \in r$  agree on all attributes in  $X$ , they also agree on all attributes in  $Y$ :

$$r \models fd \iff \text{for every } t_1, t_2 \in r \text{ if } \pi_X(t_1) = \pi_X(t_2) \text{ then } \pi_Y(t_1) = \pi_Y(t_2)$$

Here,  $\pi_X(t)$  denotes the projection of tuple  $t$  on the attributes in  $X$ .

## Key Dependencies

In the particular case when  $Y = U$ , a functional dependency of form (1) is called a *key dependency*, and the set of attributes  $X$  is called a *key* for  $R$ .

$R[X]$  functionally determines  $R[Y]$ :  
 $Y = f(X)$

**R**

...	X	Y	...
...	1	7	...
...	1	7	...
...	2	5	...
...	3	7	...

$R[X]$  is included in  $S[Z]$ :  
 $R[X] \subseteq S[Z]$

**S**

...	Z	...
...	1	...
...	2	...
...	2	...
...	3	...
...	4	...

## Inclusion Dependencies

Functional and join dependencies and their special-case subclasses each pertain to single relations. The following class of dependencies can express connections between relations. An *inclusion dependency (IND)* on pairs of relations of schemas  $R(U)$  and  $S(V)$  (with  $R$  and  $S$  not necessarily distinct) is an expression of the form

$$R[X] \subseteq S[Z], \quad (4)$$

where  $X \subseteq U$  and  $Z \subseteq V$ . Inclusion dependencies are also known as *referential constraints*. Relations  $r$  and  $s$  of schemas  $R(U)$ , respectively  $S(V)$  satisfy inclusion dependency  $id$ , denoted  $r, s \models id$ , if the projection of  $r$  on  $X$  is included in the projection of  $s$  on  $Z$ :

$$r, s \models id \iff \Pi_X(r) \subseteq \Pi_Z(s)$$

When  $R$  and  $S$  refer to the same relation name, then  $r = s$  in the above definition of satisfaction.

## Foreign Key Dependencies

In the particular case when  $Z$  is a key for relations of schema  $S$  ( $S: Z \rightarrow V$ ), INDs of form (4) are called *foreign key dependencies*. Intuitively, in this case the projection on  $X$  of every tuple  $t$  in  $r$  contains the key of a tuple from the "foreign" table  $s$ .

# Topic 1: Data models and query languages

## Unit 1: SQL (continued)

### Lecture 2

Wolfgang Gatterbauer

CS7240 Principles of scalable data management (sp24)

<https://northeastern-datalab.github.io/cs7240/sp24/>

1/12/2024

# Pre-class conversations

- Last class summary
- Class procedures based on past suggestions:
  - Your experience Canvas vs Piazza?
  - Example past scribe posted on Canvas; Secondary posting of class scribes to Piazza (optionally anonymous). I comment on both Canvas and Piazza
  - Already installed Postgres?
  - A downside of HWs with self-determined deadlines: you are in charge
  - Next week TUE online, no class FRI; see links for current research
- Today:
  - SQL continued

## *Topic 2: Complexity of Query Evaluation & Reverse Data Management*

- Unit 3: Provenance
  - [Simons'24]: 2023 Simons semester-long program on **Logic and Algorithms in Database Theory and AI**
  - [Dagstuhl'24]: 2024 Dagstuhl seminar on **Representation, Provenance, and Explanations in Database Theory and Logic**

# Outline: T1-U1: SQL

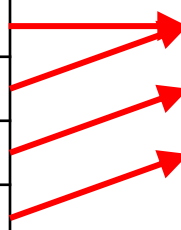
- SQL

- Schema, keys, referential integrity
- Joins
- Aggregates and grouping
- Nested queries (Subqueries)
- Union and Theta Joins
- Nulls & Outer joins
- Top-k
- [Recursion: moved to T1-U4: Datalog]



# Referential Integrity

Product				Company		
<u>PName</u>	Price	Category	Manufacturer	<u>CName</u>	StockPrice	Country
Gizmo	\$19.99	Gadgets	GizmoWorks	GizmoWorks	25	USA
Powergizmo	\$29.99	Gadgets	GizmoWorks	Canon	65	Japan
SingleTouch	\$149.99	Photography	Canon	Hitachi	15	Japan
MultiTouch	\$203.99	Household	Hitachi			

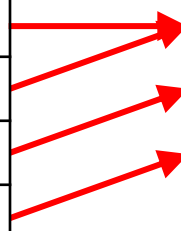


**Key constraint:** minimal subset of the attributes of a relation is a unique identifier for a tuple.

**Foreign key:** attribute in a relational table that matches a candidate key of another table

# Referential Integrity

Product				Company		
<u>PName</u>	Price	Category	Manufacturer	<u>CName</u>	StockPrice	Country
Gizmo	\$19.99	Gadgets	GizmoWorks	GizmoWorks	25	USA
Powergizmo	\$29.99	Gadgets	GizmoWorks	Canon	65	Japan
SingleTouch	\$149.99	Photography	Canon	Hitachi	15	Japan
MultiTouch	\$203.99	Household	Hitachi			



**Key constraint:** minimal subset of the attributes of a relation is a unique identifier for a tuple.

**Insert** into Product values ('Gizmo', 14.99, 'Gadgets', 'Hitachi');

Gizmo	\$14.99	Gadgets	Hitachi
-------	---------	---------	---------



**Foreign key:** attribute in a relational table that matches a candidate key of another table

# Referential Integrity

Product				Company		
<u>PName</u>	Price	Category	Manufacturer	<u>CName</u>	StockPrice	Country
Gizmo	\$19.99	Gadgets	GizmoWorks	GizmoWorks	25	USA
Powergizmo	\$29.99	Gadgets	GizmoWorks	Canon	65	Japan
SingleTouch	\$149.99	Photography	Canon	Hitachi	15	Japan
MultiTouch	\$203.99	Household	Hitachi			

**Key constraint:** minimal subset of the attributes of a relation is a unique identifier for a tuple.

`Insert into Product values ('Gizmo', 14.99, 'Gadgets', 'Hitachi');`

**tuple violates key constraint**

Gizmo	\$14.99	Gadgets	Hitachi
-------	---------	---------	---------

**Foreign key:** attribute in a relational table that matches a candidate key of another table

# Referential Integrity

Product				Company		
<u>PName</u>	Price	Category	Manufacturer	<u>CName</u>	StockPrice	Country
Gizmo	\$19.99	Gadgets	GizmoWorks	GizmoWorks	25	USA
Powergizmo	\$29.99	Gadgets	GizmoWorks	Canon	65	Japan
SingleTouch	\$149.99	Photography	Canon	Hitachi	15	Japan
MultiTouch	\$203.99	Household	Hitachi			

**Key constraint:** minimal subset of the attributes of a relation is a unique identifier for a tuple.

`Insert into Product values ('Gizmo', 14.99, 'Gadgets', 'Hitachi');`

Gizmo	\$14.99	Gadgets	Hitachi
-------	---------	---------	---------

tuple violates key constraint

**Foreign key:** attribute in a relational table that matches a candidate key of another table

`Insert into Product values ('SuperTouch', 249.99, 'Computer', 'NewCom');`

SuperTouch	\$249.99	Computer	NewCom
------------	----------	----------	--------



# Referential Integrity

Product				Company		
<u>PName</u>	Price	Category	Manufacturer	<u>CName</u>	StockPrice	Country
Gizmo	\$19.99	Gadgets	GizmoWorks	GizmoWorks	25	USA
Powergizmo	\$29.99	Gadgets	GizmoWorks	Canon	65	Japan
SingleTouch	\$149.99	Photography	Canon	Hitachi	15	Japan
MultiTouch	\$203.99	Household	Hitachi			

**Key constraint:** minimal subset of the attributes of a relation is a unique identifier for a tuple.

`Insert into Product values ('Gizmo', 14.99, 'Gadgets', 'Hitachi');`

Gizmo	\$14.99	Gadgets	Hitachi
-------	---------	---------	---------

tuple violates key constraint

**Foreign key:** attribute in a relational table that matches a candidate key of another table

`Insert into Product values ('SuperTouch', 249.99, 'Computer', 'NewCom');`

SuperTouch	\$249.99	Computer	NewCom
------------	----------	----------	--------

tuple violates foreign key constraint

# Referential Integrity

Product				Company		
<u>PName</u>	Price	Category	Manufacturer	<u>CName</u>	StockPrice	Country
Gizmo	\$19.99	Gadgets	GizmoWorks	GizmoWorks	25	USA
Powergizmo	\$29.99	Gadgets	GizmoWorks	Canon	65	Japan
SingleTouch	\$149.99	Photography	Canon	Hitachi	15	Japan
MultiTouch	\$203.99	Household	Hitachi			

**Key constraint:** minimal subset of the attributes of a relation is a unique identifier for a tuple.

```
Insert into Product values ('Gizmo', 14.99, 'Gadgets', 'Hitachi');
```

Gizmo	\$14.99	Gadgets	Hitachi
-------	---------	---------	---------

**Foreign key:** attribute in a relational table that matches a candidate key of another table

```
Insert into Product values ('SuperTouch', 249.99, 'Computer', 'NewCom');
```

SuperTouch	\$249.99	Computer	NewCom
------------	----------	----------	--------

tuple violates key constraint

tuple violates foreign key constraint

```
Delete from Company where CName = 'Canon';
```



# Referential Integrity

Product				Company		
<u>PName</u>	Price	Category	Manufacturer	<u>CName</u>	StockPrice	Country
Gizmo	\$19.99	Gadgets	GizmoWorks	GizmoWorks	25	USA
Powergizmo	\$29.99	Gadgets	GizmoWorks	Canon	65	Japan
SingleTouch	\$149.99	Photography	Canon	Hitachi	15	Japan
MultiTouch	\$203.99	Household	Hitachi			

**Key constraint:** minimal subset of the attributes of a relation is a unique identifier for a tuple.

`Insert into Product values ('Gizmo', 14.99, 'Gadgets', 'Hitachi');`

Gizmo	\$14.99	Gadgets	Hitachi
-------	---------	---------	---------

tuple violates key constraint

**Foreign key:** attribute in a relational table that matches a candidate key of another table

`Insert into Product values ('SuperTouch', 249.99, 'Computer', 'NewCom');`

SuperTouch	\$249.99	Computer	NewCom
------------	----------	----------	--------

tuple violates foreign key constraint

`Delete from Company where CName = 'Canon';`

# Schema specification in SQL



northeastern-datalab / cs3200-activities Public

Code Issues Pull requests 1 Actions Projects ...

master

cs3200-activities / sql /

wolfandthegang on May 23

- 300-SmallIMDB.txt 4 months ago
- 302-Simpleproducts.txt 4 months ago
- 304-Worker.txt 4 months ago
- 305-Conceptualevaluationstrategy.txt 4 months ago
- 306-NestedLoopJoin.py 10 months ago
- 308-Purchase.txt 4 months ago

```
-----  
-- Create the tables  
-----  
  
create table Company (  
    CName char(20) PRIMARY KEY,  
    StockPrice int,  
    Country char(20) );  
  
create table Product (  
    PName char(20),  
    Price decimal(9, 2),  
    Category char(20),  
    Manufacturer char(20),  
    PRIMARY KEY (PName),  
    FOREIGN KEY (Manufacturer) REFERENCES Company(CName) );  
  
-----  
-- Populate the tables  
-----  
  
insert into Company values ('GizmoWorks', 25, 'USA');  
insert into Company values ('Canon', 65, 'Japan');  
insert into Company values ('Hitachi', 15, 'Japan');  
  
insert into Product values ('Gizmo', 19.99, 'Gadgets', 'GizmoWorks');  
insert into Product values ('PowerGizmo', 29.99, 'Gadgets', 'GizmoWorks');
```



# Outline: T1-U1: SQL

- SQL

- Schema, keys, referential integrity
- **Joins**
- Aggregates and grouping
- Nested queries (Subqueries)
- Union and Theta Joins
- Nulls & Outer joins
- Top-k
- [Recursion: moved to T1-U4: Datalog]

## Product

PName	Price	Category	Manufacturer
Gizmo	\$19.99	Gadgets	GizmoWorks
Powergizmo	\$29.99	Gadgets	GizmoWorks
SingleTouch	\$149.99	Photography	Canon
MultiTouch	\$203.99	Household	Hitachi

## Company

CName	StockPrice	Country
GizmoWorks	25	USA
Canon	65	Japan
Hitachi	15	Japan

*Q: Find all products under \$200 manufactured in Japan;  
return their names and prices!*

?

Product				Company		
PName	Price	Category	Manufacturer	CName	StockPrice	Country
Gizmo	\$19.99	Gadgets	GizmoWorks	GizmoWorks	25	USA
Powergizmo	\$29.99	Gadgets	GizmoWorks	Canon	65	Japan
SingleTouch	\$149.99	Photography	Canon	Hitachi	15	Japan
MultiTouch	\$203.99	Household	Hitachi			

Q: Find all products under \$200 manufactured in Japan; return their names and prices!

```
SELECT pName, price
FROM Product, Company
WHERE manufacturer=cName
and country='Japan'
and price <= 200
```

Join b/w Product  
and Company

PName	Price
SingleTouch	\$149.99

Product (pName, price, category, manufacturer)  
Company (cName, stockPrice, country)

**Product**

PName	Price	Category	Manufacturer
Gizmo	\$19.99	Gadgets	GizmoWorks
Powergizmo	\$29.99	Gadgets	GizmoWorks
SingleTouch	\$149.99	Photography	Canon
MultiTouch	\$203.99	Household	Hitachi

**Company**

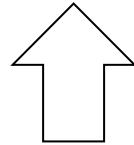
CName	StockPrice	Country
GizmoWorks	25	USA
Canon	65	Japan
Hitachi	15	Japan

```
SELECT *  
FROM Product, Company  
WHERE manufacturer=cName
```



Product (pName, price, category, manufacturer)  
Company (cName, stockPrice, country)

PName	Price	Category	Manufacturer	CName	StockPrice	Country
Gizmo	\$19.99	Gadgets	GizmoWorks	GizmoWorks	25	USA
Powergizmo	\$29.99	Gadgets	GizmoWorks	GizmoWorks	25	USA
SingleTouch	\$149.99	Photography	Canon	Canon	65	Japan
MultiTouch	\$203.99	Household	Hitachi	Hitachi	15	Japan



```
SELECT *  
FROM Product, Company  
WHERE manufacturer=cName
```

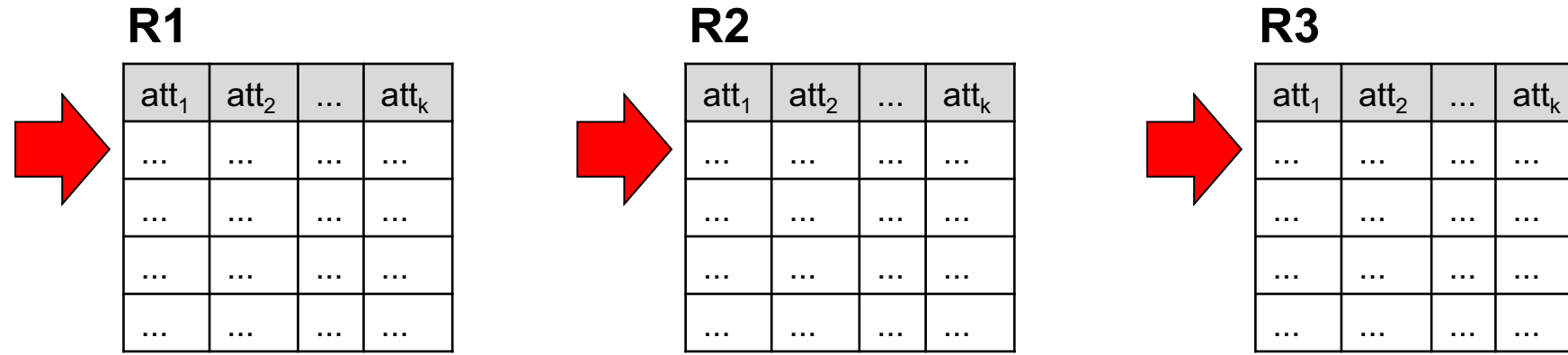
# Meaning (Semantics) of SELECT-FROM-WHERE queries

```
3 SELECT a1, a2, ..., ak  
1 FROM R1 as x1, R2 as x2, ..., Rn as xn  
2 WHERE Conditions
```

**Conceptual evaluation strategy (nested for loops):**

```
Answer = {}  
for x1 in R1 do  
  for x2 in R2 do  
    .....  
    for xn in Rn do  
      if Conditions  
        then Answer = Answer ∪ {(a1, ..., ak)}  
return Answer
```

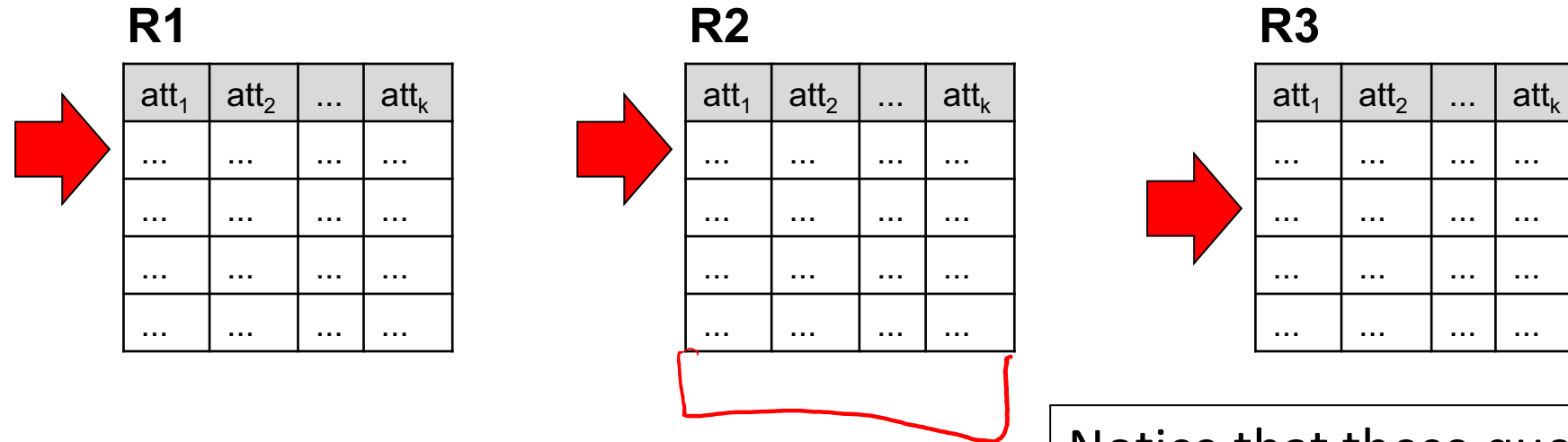
# Meaning (Semantics) of SELECT-FROM-WHERE queries



**Conceptual evaluation strategy (nested for loops):**

```
Answer = {}  
for  $x_1$  in  $R_1$  do  
  for  $x_2$  in  $R_2$  do  
    .....  
    for  $x_n$  in  $R_n$  do  
      if Conditions  
        then Answer = Answer  $\cup$   $\{(a_1, \dots, a_k)\}$   
return Answer
```

# Meaning (Semantics) of SELECT-FROM-WHERE queries

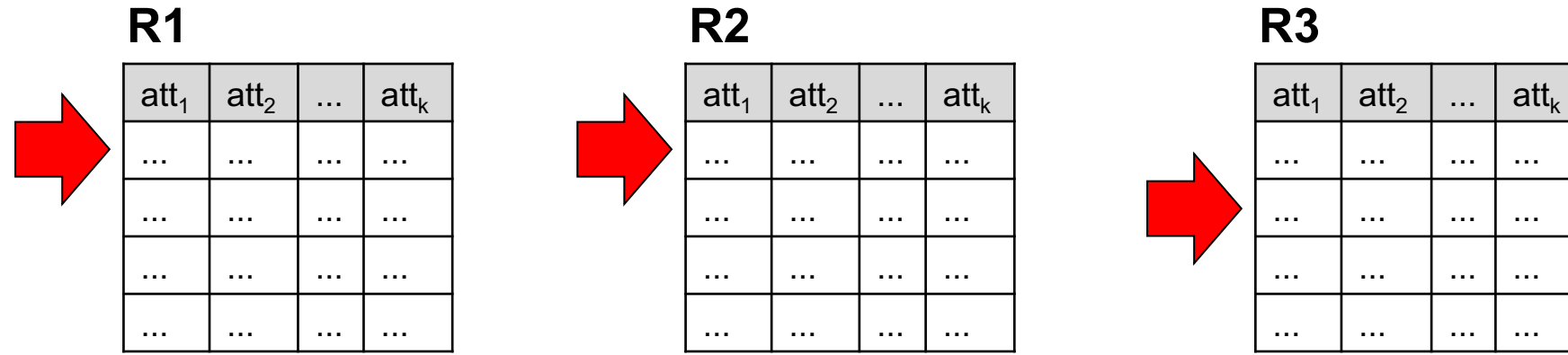


```
Answer = {}  
for x1 in R1 do  
  for x2 in R2 do  
    .....  
    for xn in Rn do  
      if Conditions  
        then Answer = Answer ∪ {(a1, ..., ak)}  
return Answer
```

Notice that these queries are "**monotone**":  
whenever we add tuples to the input,  
the output can never decrease:  
**if**  $R_1 \subseteq R'_1, R_2 \subseteq R'_2, R_3 \subseteq R'_3$   
**then**  $Q(R_1, R_2, R_3) \subseteq Q(R'_1, R'_2, R'_3)$



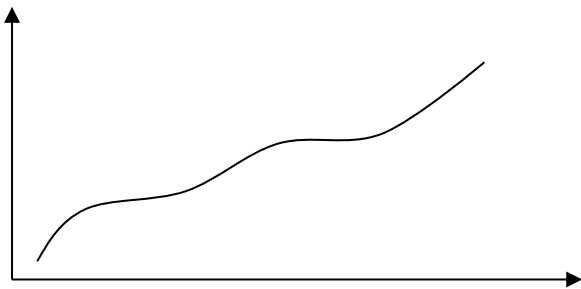
# Meaning (Semantics) of conjunctive SQL Queries



DEFINITION: A function  $f(x)$  is "**monotone**" (or better "**monotonically increasing**") if:

**if**  $x \leq y$

**then**  $f(x) \leq f(y)$



Notice that these queries are "**monotone**": whenever we add tuples to the input, the output can never decrease:

**if**  $R_1 \subseteq R'_1, R_2 \subseteq R'_2, R_3 \subseteq R'_3$

**then**  $Q(R_1, R_2, R_3) \subseteq Q(R'_1, R'_2, R'_3)$

# Conceptual Evaluation Strategy



- Semantics of an SQL query defined in terms of the following conceptual evaluation strategy:
  - **FROM**: Compute the cross-product of the relations. This is a new set of larger tuples.
  - **WHERE**: Only keep the tuples that pass the qualifications ("selection", filter)
  - **SELECT**: Delete attributes that are not in listed attributes
  - If **DISTINCT** is specified, eliminate duplicate rows.
- This strategy is probably the least efficient way to compute a query! An optimizer will find (**algebraically equivalent** but) more efficient strategies to compute the same answers.
- We say “semantics” not “execution order”. Why?



# Conceptual Evaluation Strategy

- Semantics of an SQL query defined in terms of the following conceptual evaluation strategy:
  - **FROM**: Compute the cross-product of the relations. This is a new set of larger tuples.
  - **WHERE**: Only keep the tuples that pass the qualifications ("selection", filter)
  - **SELECT**: Delete attributes that are not in listed attributes
  - If **DISTINCT** is specified, eliminate duplicate rows.
- This strategy is probably the least efficient way to compute a query! An optimizer will find (**algebraically equivalent** but) more efficient strategies to compute the same answers.
- We say “semantics” not “execution order”. Why?
  - The preceding slides show **what** a join means (**semantics = meaning**): "the logic"
  - Not actually **how** the DBMS actually executes it (separation of concerns): **algebra**

# Table Alias (Tuple Variables)



```
Person (pName, address, works_for)
University (uName, address)
```

```
SELECT DISTINCT pName, address
FROM Person, University
WHERE works_for = uName
```

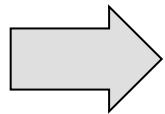
What will this query return ?

# Table Alias (Tuple Variables)

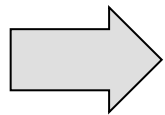
```
Person (pName, address, works_for)
University (uName, address)
```

```
SELECT DISTINCT pName, address
FROM Person, University
WHERE works_for = uName
```

*which address?  
Error!*



```
SELECT DISTINCT pName, University.address
FROM Person, University
WHERE Person.works_for = University.uName
```



```
SELECT DISTINCT X.pName, Y.address
FROM Person as X, University Y
WHERE X.works_for = Y.uName
```

*Notice that the use of "as" is not necessary, it is optional !!*

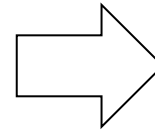
# Using the Formal Semantics

R(a), S(a), T(a)

*What do these queries compute?*

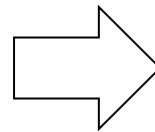
R	S	T
a	a	a
1	1	2
2		

```
SELECT R.a
FROM   R, S
WHERE  R.a=S.a
```



?

```
SELECT R.a
FROM   R, S, T
WHERE  R.a=S.a
      or R.a=T.a
```



?

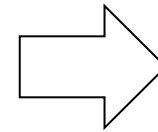
# Using the Formal Semantics

R(a), S(a), T(a)

*What do these queries compute?*

R	S	T
a	a	a
1	1	2
2		

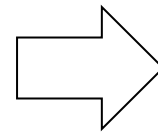
```
SELECT R.a
FROM R, S
WHERE R.a=S.a
```



a
1

Returns  $R \cap S$   
(intersection)

```
SELECT R.a
FROM R, S, T
WHERE R.a=S.a
      or R.a=T.a
```



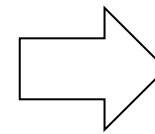
# Using the Formal Semantics

R(a), S(a), T(a)

*What do these queries compute?*

R	S	T
a	a	a
1	1	2
2		

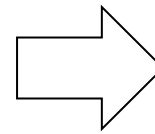
```
SELECT R.a
FROM R, S
WHERE R.a=S.a
```



a
1

Returns  $R \cap S$   
(intersection)

```
SELECT R.a
FROM R, S, T
WHERE R.a=S.a
      or R.a=T.a
```



a
1
2

Returns  $R \cap (S \cup T)$   
if  $S \neq \emptyset$  and  $T \neq \emptyset$



# Using the Formal Semantics

R(a), S(a), T2(a)

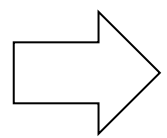


What do these queries compute?

R	S	T2
a	a	a
1	1	<del>2</del>
2		

Our colorful hands represent "team exercises" If we are online, please make a screenshot!

```
SELECT R.a
FROM R, S
WHERE R.a=S.a
```



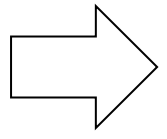
a
1

Returns  $R \cap S$  (intersection)



Next, we are removing the input tuple "(2)"

```
SELECT R.a
FROM R, S, T2 as T
WHERE R.a=S.a
or R.a=T.a
```



a
1
2

Returns  $R \cap (S \cup T)$  if  $S \neq \emptyset$  and  $T \neq \emptyset$



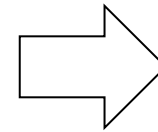
# Using the Formal Semantics

R(a), S(a), T2(a)

*What do these queries compute?*

R	S	T2
a	a	a
1	1	<del>1</del>
2		

```
SELECT R.a
FROM R, S
WHERE R.a=S.a
```

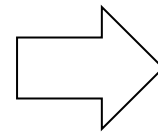


a
1

Returns  $R \cap S$   
(intersection)

Next, we are removing the input tuple "(2)"

```
SELECT R.a
FROM R, S, T2 as T
WHERE R.a=S.a
or R.a=T.a
```



a
1
2

Returns  $R \cap (S \cup T)$   
if  $S \neq \emptyset$  and  $T \neq \emptyset$



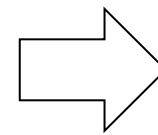
# Using the Formal Semantics

What do these queries compute?

R(a), S(a), T2(a)

R	S	T2
a	a	a
1	1	<del>2</del>
2	3	

```
SELECT R.a
FROM R, S
WHERE R.a=S.a
```

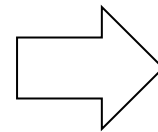


a
1

Returns  $R \cap S$   
(intersection)

Next, we are removing the input tuple "(2)"

```
SELECT R.a
FROM R, S, T2 as T
WHERE R.a=S.a
or R.a=T.a
```



a
---

Returns  $\emptyset$   
if  $S = \emptyset$  or  $T = \emptyset$

Can seem counterintuitive! But remember conceptual evaluation strategy: Nested loops. If one table is empty -> no looping

# Illustration with Python



```
1 '''
2 Created on 3/23/2015
3 Illustrates nested Loop Join in SQL
4 __author__ = 'gatt'
5 '''
6
7 print "--- 1st nested loop ---"
8 for i in xrange(2):
9     for j in xrange(3):
10        for k in xrange(2):
11            print "i=%d, j=%d, k=%d: " % (i, j, k),
12            if i == j or i == k:
13                print "TRUE",
14            print
15
16 print "\n--- 2nd nested loop ---"
17 for i in xrange(2):
18     for j in xrange(3):
19        for k in xrange(1):
20            print "i=%d, j=%d, k=%d: " % (i, j, k),
21            if i == j or i == k:
22                print "TRUE",
23            print
24
25 print "\n--- 3rd nested loop ---"
26 for i in xrange(2):
27     for j in xrange(3):
28        for k in xrange(0):
29            print "i=%d, j=%d, k=%d: " % (i, j, k),
30            if i == j or i == k:
31                print "TRUE",
32            print
33
```

```
/Library/Frameworks/Python.framework/Versio
--- 1st nested loop ---
i=0, j=0, k=0: TRUE
i=0, j=0, k=1: TRUE
i=0, j=1, k=0: TRUE
i=0, j=1, k=1:
i=0, j=2, k=0: TRUE
i=0, j=2, k=1:
i=1, j=0, k=0:
i=1, j=0, k=1: TRUE
i=1, j=1, k=0: TRUE
i=1, j=1, k=1: TRUE
i=1, j=2, k=0:
i=1, j=2, k=1: TRUE

--- 2nd nested loop ---
i=0, j=0, k=0: TRUE
i=0, j=1, k=0: TRUE
i=0, j=2, k=0: TRUE
i=1, j=0, k=0:
i=1, j=1, k=0: TRUE
i=1, j=2, k=0:

--- 3rd nested loop ---
Process finished with exit code 0
```

"Premature optimization is the root of all evil."  
Donald Knuth (1974)

"When you are diagnosing problems, don't think about how you will solve them—just diagnose them. Blurring the steps leads to suboptimal outcomes because it interferes with uncovering the true problems."  
Ray Dalio (Principles, 2017)

The comparison gets never evaluated

Our colorful hands represent "team exercises"  
If we are online, please make a screenshot!



Product (<sup>PK</sup> pName, price, category, manufacturer)  
Company (<sup>FK</sup> cName, stockPrice, country)

Q: Find all US companies that manufacture both a product below \$20 and a product above \$25.

```
SELECT DISTINCT cName
FROM
WHERE
```

# Quiz: Answer 1

Our colorful hands represent "team exercises"  
If we are online, please make a screenshot!



302

Product (pName, price, category, manufacturer)  
Company (cName, stockPrice, country)

*Q: Find all US companies that manufacture both a product below \$20 and a product above \$25.*

```
SELECT DISTINCT cName
FROM Product as P, Company
WHERE country = 'USA'
      and P.price < 20
      and P.price > 25
      and P.manufacturer = cName
```

What about this query?



# Quiz: Answer 1



302

Product (pName, price, category, manufacturer)  
Company (cName, stockPrice, country)

*Q: Find all US companies that manufacture both a product below \$20 and a product above \$25.*

```
SELECT DISTINCT cName  
FROM Product as P, Company  
WHERE country = 'USA'  
and P.price < 20  
and P.price > 25  
and P.manufacturer = cName
```

*Wrong! Gives empty result: There is no product with price <20 and >25*

# Quiz: Answer 2



302

Product (pName, price, category, manufacturer)  
Company (cName, stockPrice, country)

*Q: Find all US companies that manufacture both a product below \$20 and a product above \$25.*

```
SELECT DISTINCT cName
FROM Product as P, Company
WHERE country = 'USA'
      and (P.price < 20
      or P.price > 25)
      and P.manufacturer = cName
```

What about this query?





# Quiz: Answer 2



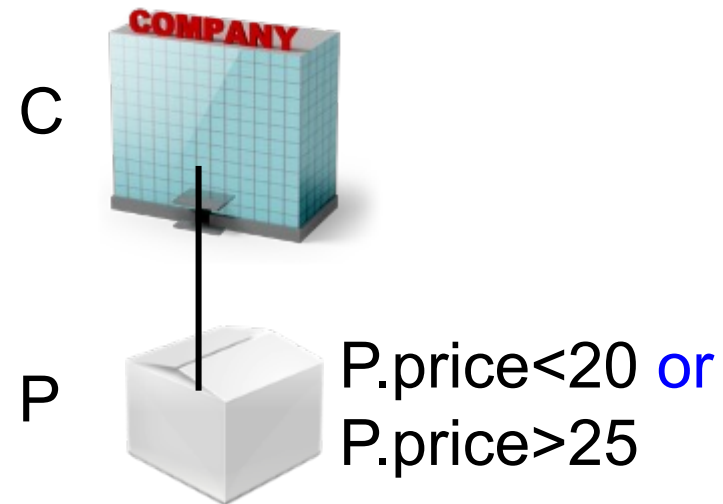
302

Product (pName, price, category, manufacturer)  
Company (cName, stockPrice, country)

*Q: Find all US companies that manufacture both a product below \$20 and a product above \$25.*

```
SELECT DISTINCT cName  
FROM Product as P, Company  
WHERE country = 'USA'  
and (P.price < 20  
or P.price > 25)  
and P.manufacturer = cName
```

Returns companies  
with single product  
w/price (<20 or >25)



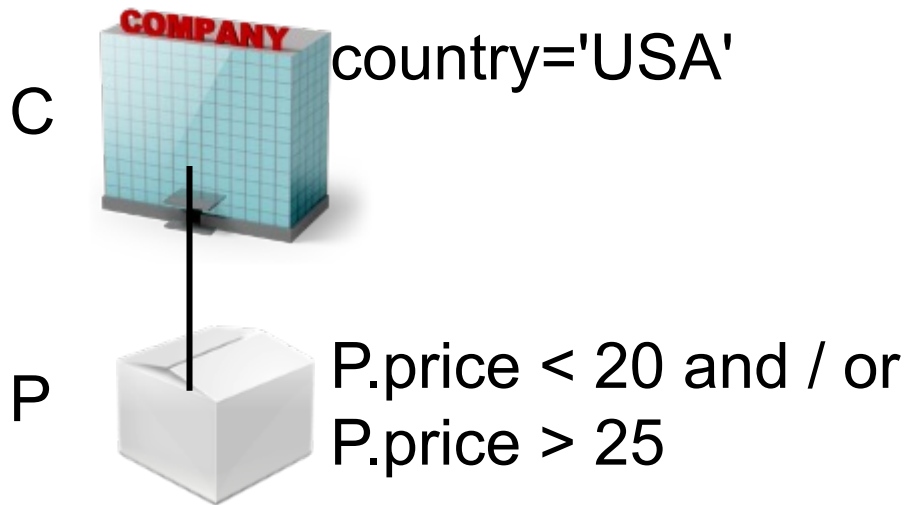
# Quiz: Answer 1



302

Product (pName, price, category, manufacturer)  
Company (cName, stockPrice, country)

*Q: Find all US companies that manufacture both a product below \$20 and a product above \$25.*



What do we actually want?



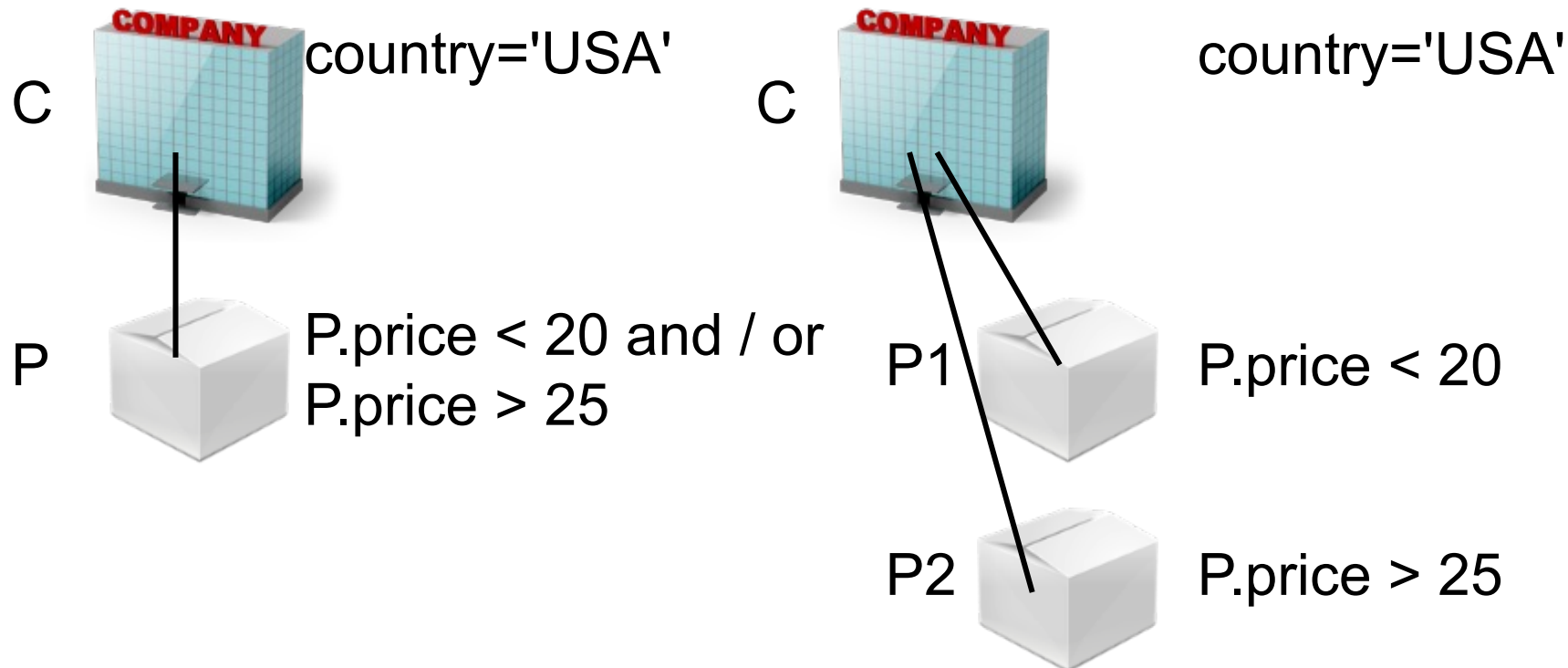
# Quiz: Answer 1 vs. what we actually want



302

Product (pName, price, category, manufacturer)  
Company (cName, stockPrice, country)

*Q: Find all US companies that manufacture both a product below \$20 and a product above \$25.*



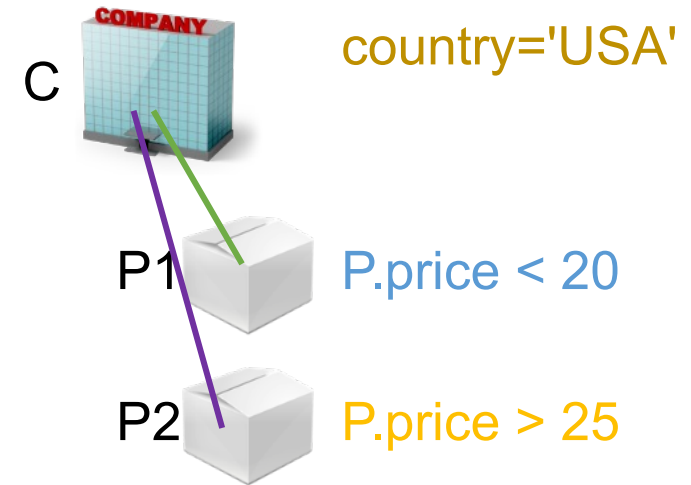
# Quiz: correct answer: we need "self-joins"!



Product (pName, price, category, manufacturer)  
Company (cName, stockPrice, country)

*Q: Find all US companies that manufacture both a product below \$20 and a product above \$25.*

```
SELECT DISTINCT cName
FROM Product as P1, Product as P2, Company
WHERE country = 'USA'
      and P1.price < 20
      and P2.price > 25
      and P1.manufacturer = cName
      and P2.manufacturer = cName
```



# Quiz response: we need "self-joins"!

P1




PName	Price	Category	Manufacturer
Gizmo	\$19.99	Gadgets	GizmoWorks
Powergizmo	\$29.99	Gadgets	GizmoWorks
SingleTouch	\$149.99	Photography	Canon
MultiTouch	\$203.99	Household	Hitachi

P2



PName	Price	Category	Manufacturer
Gizmo	\$19.99	Gadgets	GizmoWorks
Powergizmo	\$29.99	Gadgets	GizmoWorks
SingleTouch	\$149.99	Photography	Canon
MultiTouch	\$203.99	Household	Hitachi

Company

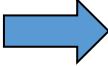


CName	StockPrice	Country
GizmoWorks	25	USA
Canon	65	Japan
Hitachi	15	Japan

```
SELECT DISTINCT cName
FROM Product as P1, Product as P2, Company
WHERE country = 'USA'
      and P1.price < 20
      and P2.price > 25
      and P1.manufacturer = cName
      and P2.manufacturer = cName
```


# Quiz response: we need "self-joins"!

P1




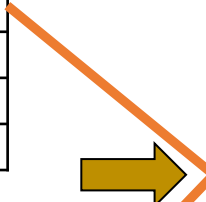
PName	Price	Category	Manufacturer
Gizmo	\$19.99	Gadgets	GizmoWorks
Powergizmo	\$29.99	Gadgets	GizmoWorks
SingleTouch	\$149.99	Photography	Canon
MultiTouch	\$203.99	Household	Hitachi

P2



PName	Price	Category	Manufacturer
Gizmo	\$19.99	Gadgets	GizmoWorks
Powergizmo	\$29.99	Gadgets	GizmoWorks
SingleTouch	\$149.99	Photography	Canon
MultiTouch	\$203.99	Household	Hitachi

Company



CName	StockPrice	Country
GizmoWorks	25	USA
Canon	65	Japan
Hitachi	15	Japan

```
SELECT DISTINCT cName
FROM Product as P1, Product as P2, Company
WHERE country = 'USA'
    and P1.price < 20
    and P2.price > 25
    and P1.manufacturer = cName
    and P2.manufacturer = cName
```



CName
GizmoWorks

# Outline: T1-U1: SQL

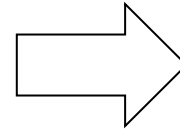
- SQL

- Schema, keys, referential integrity
- Joins
- **Aggregates and grouping**
- Nested queries (Subqueries)
- Union and Theta Joins
- Nulls & Outer joins
- Top-k
- [Recursion: moved to T1-U4: Datalog]

# Grouping and Aggregation

## Purchase

Product	Price	Quantity
Bagel	3	20
Bagel	2	20
Banana	1	50
Banana	2	10
Banana	4	10



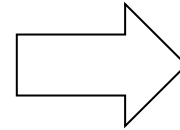
*Q: For each product, find Total Quantities (TQ = sum of quantities) purchased, for all products with price >1.*



# Grouping and Aggregation

## Purchase

Product	Price	Quantity
Bagel	3	20
Bagel	2	20
Banana	1	50
Banana	2	10
Banana	4	10



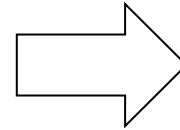
Product	TQ
Bagel	?
Banana	?

*Q: For each product, find Total Quantities (TQ = sum of quantities) purchased, for all products with price >1.*

# Grouping and Aggregation

## Purchase

Product	Price	Quantity
Bagel	3	20
Bagel	2	20
<del>Banana</del>	<del>1</del>	<del>50</del>
Banana	2	10
Banana	4	10



Product	TQ
Bagel	40
Banana	20

*Q: For each product, find Total Quantities (TQ = sum of quantities) purchased, for all products with price >1.*

# From → Where → Group By → Select

## Purchase

Product	Price	Quantity
Bagel	3	20
Bagel	2	20
<del>Banana</del>	<del>1</del>	<del>50</del>
Banana	2	10
Banana	4	10

→

Product	TQ
Bagel	40
Banana	20

Select contains

- grouped attributes
- and aggregates

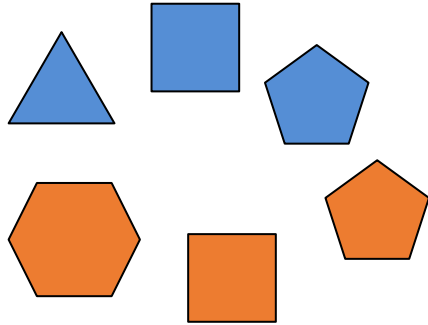
```
4 SELECT product, sum(quantity) as TQ
1 FROM Purchase
2 WHERE price > 1
3 GROUP BY product
```

Tuples grouped together need to share the same value for attribute "product"

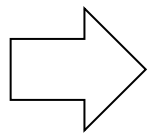


# Groupings illustrated with colored shapes

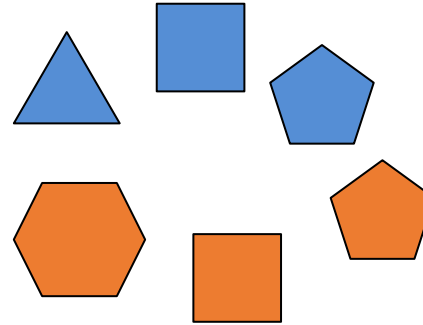
group by color



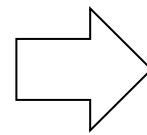
```
SELECT color,  
       avg(numc) anc  
FROM   Shapes  
GROUP BY color
```



group by numc (# of corners)

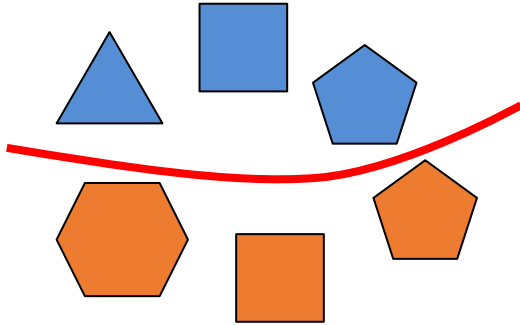


```
SELECT numc  
FROM   Shapes  
GROUP BY numc
```

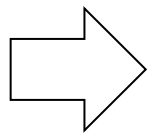


# Groupings illustrated with colored shapes

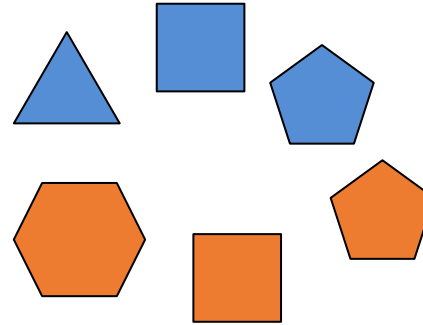
group by color



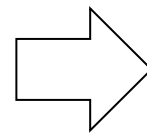
```
SELECT color,  
       avg(numc) anc  
FROM   Shapes  
GROUP BY color
```



group by numc (# of corners)



```
SELECT numc  
FROM   Shapes  
GROUP BY numc
```

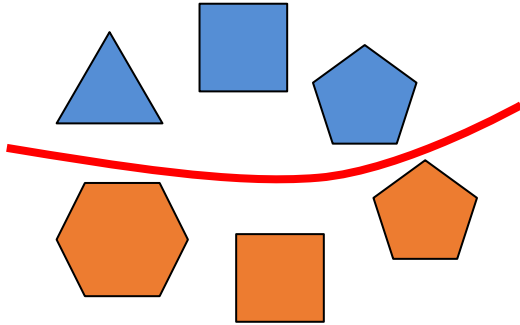


color	numc
blue	3
blue	4
blue	5
orange	4
orange	5
orange	6



# Groupings illustrated with colored shapes

group by color

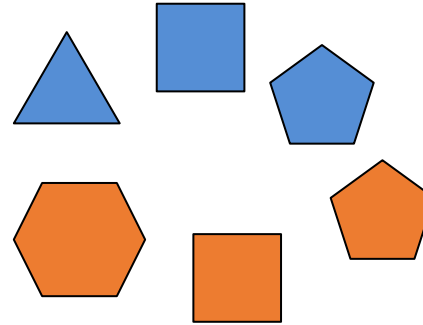


```
SELECT color,  
       avg(numc) anc  
FROM   Shapes  
GROUP BY color
```

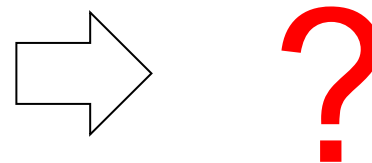
⇒

color	anc
blue	4
orange	5

group by numc (# of corners)



```
SELECT numc  
FROM   Shapes  
GROUP BY numc
```



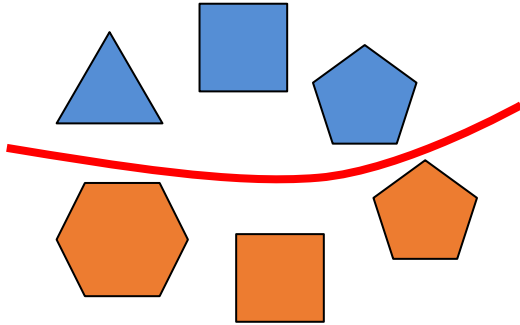
color	numc
blue	3
blue	4
blue	5
orange	4
orange	5
orange	6



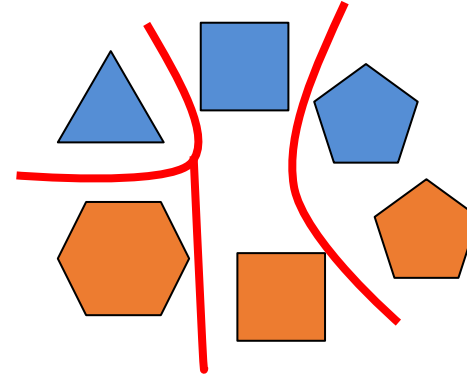
# Groupings illustrated with colored shapes



group by color



group by numc (# of corners)



```
SELECT color,  
       avg(numc) anc  
FROM   Shapes  
GROUP BY color
```

⇒

color	anc
blue	4
orange	5

```
SELECT numc  
FROM   Shapes  
GROUP BY numc
```

⇒

numc
3
4
5
6

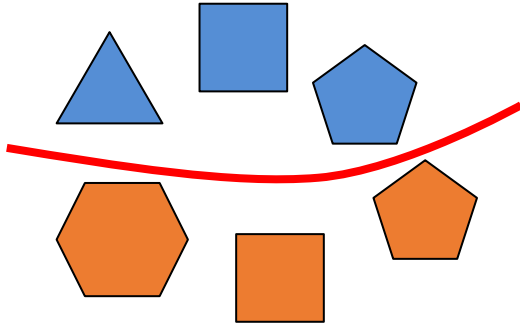
Without group by ?

color	numc
blue	3
blue	4
blue	5
orange	4
orange	5
orange	6

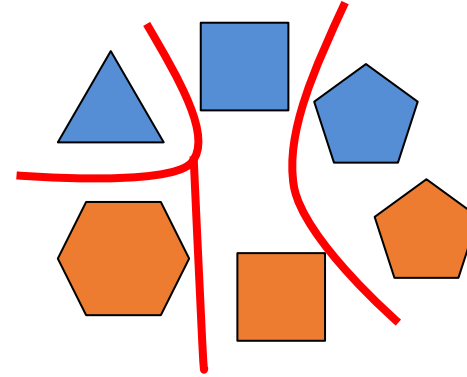
# Groupings illustrated with colored shapes



group by color



group by numc (# of corners)



color	numc
blue	3
blue	4
blue	5
orange	4
orange	5
orange	6

```
SELECT color,  
       avg(numc) anc  
FROM   Shapes  
GROUP BY color
```

➔

color	anc
blue	4
orange	5

```
SELECT numc  
FROM   Shapes  
GROUP BY numc
```

➔

numc
3
4
5
6

Same as:

```
SELECT DISTINCT numc  
FROM   Shapes
```

Without group by!



# Outline: T1-U1: SQL

- SQL

- Schema, keys, referential integrity
- Joins
- Aggregates and grouping
- **Nested queries (Subqueries)**
- Union and Theta Joins
- Nulls & Outer joins
- Top-k
- [Recursion: moved to T1-U4: Datalog]

# Subqueries = Nested queries

Outer block

```
SELECT ...  
FROM ...  
WHERE ...  
HAVING ...
```

Inner block

```
(SELECT ...  
FROM ...  
WHERE ... )
```

We focus mainly on nestings in the WHERE clause, which are the most expressive type of nesting.

- We can nest queries because SQL is **compositional**:
  - **Input & Output** are represented as **relations (multisets)**
  - Subqueries also return relations; thus the output of one query can thus be used as the input to another (**nesting**)
- This is extremely powerful (think in terms of input/output)
- A complication: subqueries can be **correlated** (not just in-/output)

# Subqueries in

SELECT clause

FROM clause

*(also called "derived tables")*

WHERE clause

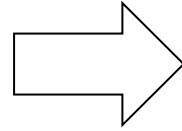
HAVING clause

# Subqueries in FROM clause = Derived tables

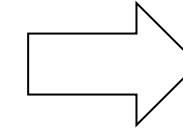


## Purchase

Product	Price	Quantity
Bagel	3	20
Bagel	2	20
Banana	1	50
Banana	2	10
Banana	4	10



Product	TQ
Bagel	40
Banana	70



MTQ
70

Q1: For each product, find total quantities (sum of quantities) purchased.

Q2: Find the maximal total quantities purchased across all products.

```
SELECT product, SUM(quantity) as TQ
FROM Purchase
GROUP BY product
```

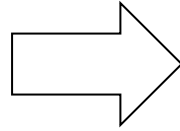


# Subqueries in FROM clause = Derived tables



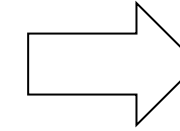
## Purchase

Product	Price	Quantity
Bagel	3	20
Bagel	2	20
Banana	1	50
Banana	2	10
Banana	4	10



X

Product	TQ
Bagel	40
Banana	70



MTQ
70

Q1: For each product, find total quantities (sum of quantities) purchased.

Q2: Find the maximal total quantities purchased across all products.

```
SELECT product, SUM(quantity) as TQ
FROM Purchase
GROUP BY product
```

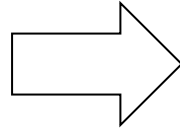


# Subqueries in FROM clause = Derived tables

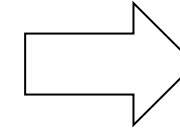


## Purchase

Product	Price	Quantity
Bagel	3	20
Bagel	2	20
Banana	1	50
Banana	2	10
Banana	4	10

**X**

Product	TQ
Bagel	40
Banana	70



MTQ
70

Q1: For each product, find total quantities (sum of quantities) purchased.

```
SELECT product, SUM(quantity) as TQ
FROM Purchase
GROUP BY product
```

Q2: Find the maximal total quantities purchased across all products.

```
SELECT MAX(TQ) as MTQ
FROM X
```

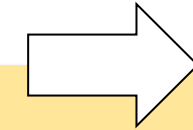
# Subqueries in FROM clause = Derived tables



## Purchase

Product	Price	Quantity
Bagel	3	20
Bagel	2	20
Banana	1	50
Banana	2	10
Banana	4	10

```
SELECT MAX(TQ) as MTQ
FROM (SELECT product, SUM(quantity) as TQ
      FROM Purchase
      GROUP BY product) X
```



MTQ
70

Q1: For each product, find total quantities (sum of quantities) purchased.

```
SELECT product, SUM(quantity) as TQ
FROM Purchase
GROUP BY product
```

Q2: Find the maximal total quantities purchased across all products.

```
SELECT MAX(TQ) as MTQ
FROM X
```

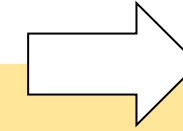
# Common Table Expressions (CTE): WITH clause



## Purchase

Product	Price	Quantity
Bagel	3	20
Bagel	2	20
Banana	1	50
Banana	2	10
Banana	4	10

SELECT MAX(TQ) as MTQ  
FROM (SELECT product, SUM(quantity) as TQ  
FROM Purchase  
GROUP BY product) X



MTQ
70

CTE (Common  
Table Expression)

Query using CTE

WITH X as  
(SELECT product, SUM(quantity) as TQ  
FROM Purchase  
GROUP BY product)  
SELECT MAX(TQ) as MTQ  
FROM X

The WITH clause defines a temporary relation that is available only to the query in which it occurs. Sometimes easier to read. Very useful for queries that need to access the same intermediate result multiple times



# Subqueries in

SELECT clause

FROM clause

WHERE clause

*(including IN, ANY, ALL)*

HAVING clause

# Subqueries in WHERE clause



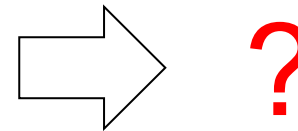
305

What do these queries return?

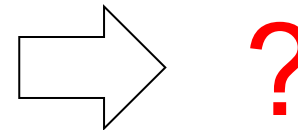
R
a
1
2

W	
a	b
2	0
3	0
4	0

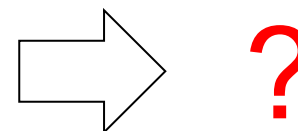
```
SELECT a
FROM R
WHERE a IN
      (SELECT a FROM W)
```



```
SELECT a
FROM R
WHERE a < ANY
      (SELECT a FROM W)
```



```
SELECT a
FROM R
WHERE a < ALL
      (SELECT a FROM W)
```



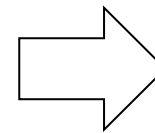
# Subqueries in WHERE clause

What do these queries return?

R
a
1
2

W	
a	b
2	0
3	0
4	0

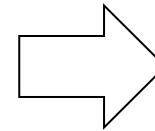
```
SELECT a
FROM R
WHERE a IN
      (SELECT a FROM W)
```



a
2

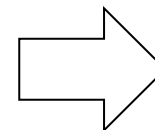
Since 2 is in the set (bag)  
(2, 3, 4)

```
SELECT a
FROM R
WHERE a < ANY
      (SELECT a FROM W)
```



?

```
SELECT a
FROM R
WHERE a < ALL
      (SELECT a FROM W)
```



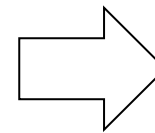
?

# Subqueries in WHERE clause

What do these queries return?

R	W	
a	a	b
1	2	0
2	3	0
	4	0

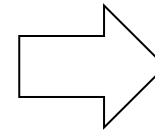
```
SELECT a
FROM R
WHERE a IN
      (SELECT a FROM W)
```



a
2

Since 2 is in the set (bag)  
(2, 3, 4)

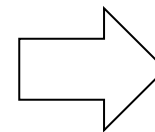
```
SELECT a
FROM R
WHERE a < ANY
      (SELECT a FROM W)
```



a
1
2

Since 1 and 2 are <  
than at least one  
("any") of 2, 3 or 4

```
SELECT a
FROM R
WHERE a < ALL
      (SELECT a FROM W)
```



?

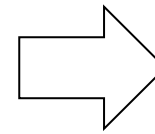
# Subqueries in WHERE clause

What do these queries return?

SQLite does not support "ANY" or "ALL" ☹️

R	W	
a	a	b
1	2	0
2	3	0
	4	0

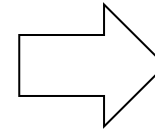
```
SELECT a
FROM R
WHERE a IN
      (SELECT a FROM W)
```



a
2

Since 2 is in the set (bag)  
(2, 3, 4)

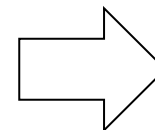
```
SELECT a
FROM R
WHERE a < ANY
      (SELECT a FROM W)
```



a
1
2

Since 1 and 2 are <  
than at least one  
("any") of 2, 3 or 4

```
SELECT a
FROM R
WHERE a < ALL
      (SELECT a FROM W)
```



a
1

Since 1 is < than  
each ("all") of 2, 3,  
and 4

# Correlated subqueries

- In all previous cases, the nested subquery in the inner select block could be entirely evaluated before processing the outer select block.
  - Recall the "**compositional**" nature of relational queries
  - This is no longer the case for **correlated nested queries**.
- Whenever a condition in the WHERE clause of a nested query references some column of a table declared in the outer query, the two queries are said to be correlated.
  - The nested query is then evaluated once for each tuple (or combination of tuples) in the outer query (that's the **conceptual evaluation strategy**)

# Correlated subquery (existential $\exists$ )



**Product**

PName	Price	Category	cid
Gizmo	\$19.99	Gadgets	1
Powergizmo	\$29.99	Gadgets	1
SingleTouch	\$14.99	Photography	2
MultiTouch	\$203.99	Household	3

**Company**

cid	CName	StockPrice	Country
1	GizmoWorks	25	USA
2	Canon	65	Japan
3	Hitachi	15	Japan

slightly  
different  
product  
database!

Q<sub>1</sub>: Find all companies that make some product(s) with price < 25

Using **IN**: Set / Bag membership

```
SELECT DISTINCT C.cname
FROM Company C
WHERE C.cid IN ( SELECT P.cid
                 FROM Product P
                 WHERE P.price < 25)
```

Is this a correlated  
nested query



# Correlated subquery (existential $\exists$ )



**Product**

PName	Price	Category	cid
Gizmo	\$19.99	Gadgets	1
Powergizmo	\$29.99	Gadgets	1
SingleTouch	\$14.99	Photography	2
MultiTouch	\$203.99	Household	3

**Company**

cid	CName	StockPrice	Country
1	GizmoWorks	25	USA
2	Canon	65	Japan
3	Hitachi	15	Japan

slightly different product database!

Q<sub>1</sub>: Find all companies that make some product(s) with price < 25

Using **IN**: Set / Bag membership

```
SELECT DISTINCT C.cname
FROM Company C
WHERE C.cid IN ( SELECT P.cid
                  FROM Product P
                  WHERE P.price < 25)
```

Not a correlated nested query!

```
SELECT DISTINCT C.cname
FROM Company C
WHERE C.cid IN ( 1, 2 )
```

Inner query does not reference outer query! You could first evaluate the inner query by itself.



# Correlated subquery (existential $\exists$ )

**Product**

PName	Price	Category	cid
Gizmo	\$19.99	Gadgets	1
Powergizmo	\$29.99	Gadgets	1
SingleTouch	\$14.99	Photography	2
MultiTouch	\$203.99	Household	3

**Company**

cid	CName	StockPrice	Country
1	GizmoWorks	25	USA
2	Canon	65	Japan
3	Hitachi	15	Japan

Q<sub>1</sub>: Find all companies that make some product(s) with price < 25

Using **EXISTS**: TRUE if the subquery's result is **NOT empty**

```
SELECT DISTINCT C.cname
FROM Company C
WHERE EXISTS ( SELECT *
                FROM Product P
                WHERE P.cid = C.cid
                and P.price < 25)
```

Is this a correlated  
nested query



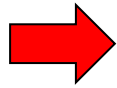
# Correlated subquery (existential $\exists$ )

**Product**



PName	Price	Category	cid
Gizmo	\$19.99	Gadgets	1
Powergizmo	\$29.99	Gadgets	1
SingleTouch	\$14.99	Photography	2
MultiTouch	\$203.99	Household	3

**Company**



cid	CName	StockPrice	Country
1	GizmoWorks	25	USA
2	Canon	65	Japan
3	Hitachi	15	Japan

Q<sub>1</sub>: Find all companies that make some product(s) with price < 25

Using **EXISTS**: TRUE if the subquery's result is **NOT empty**

```
SELECT DISTINCT C.cname
FROM Company C
WHERE EXISTS ( SELECT *
                FROM Product P
                WHERE P.cid = C.cid
                and P.price < 25)
```

*This is a correlated nested query!  
Notice the additional join condition  
referencing a relation from the  
outer query.*

*Recall our conceptual evaluation  
strategy!*

# Correlated subquery (existential $\exists$ )

**Product**

PName	Price	Category	cid
Gizmo	\$19.99	Gadgets	1
Powergizmo	\$29.99	Gadgets	1
SingleTouch	\$14.99	Photography	2
MultiTouch	\$203.99	Household	3

**Company**

cid	CName	StockPrice	Country
1	GizmoWorks	25	USA
2	Canon	65	Japan
3	Hitachi	15	Japan

Q<sub>1</sub>: Find all companies that make some product(s) with price < 25

Using **ANY** (also **SOME**): again **set / bag comparison**

```
SELECT DISTINCT C.cname
FROM Company C
WHERE 25 > ANY ( SELECT price
                  FROM Product P
                  WHERE P.cid = C.cid)
```

*But do we really need  
to write this query as  
nested query*



**SQLite does not support "ANY" ☹️**

# Correlated subquery (existential $\exists$ )

**Product**

PName	Price	Category	cid
Gizmo	\$19.99	Gadgets	1
Powergizmo	\$29.99	Gadgets	1
SingleTouch	\$14.99	Photography	2
MultiTouch	\$203.99	Household	3

**Company**

cid	CName	StockPrice	Country
1	GizmoWorks	25	USA
2	Canon	65	Japan
3	Hitachi	15	Japan

Q<sub>1</sub>: Find all companies that make some product(s) with price < 25

```
SELECT DISTINCT C.cname
FROM Company C, Product P
WHERE C.cid = P.cid
and P.price < 25
```

*We did not need to write nested queries;  
we can "unnest" it!*

*Existential quantifiers are easy 😊*

# Correlated subquery (universal $\forall$ )

**Product**

PName	Price	Category	cid
Gizmo	\$19.99	Gadgets	1
Powergizmo	\$29.99	Gadgets	1
SingleTouch	\$14.99	Photography	2
MultiTouch	\$203.99	Household	3

**Company**

cid	CName	StockPrice	Country
1	GizmoWorks	25	USA
2	Canon	65	Japan
3	Hitachi	15	Japan

~~Q<sub>1</sub>: Find all companies that make some product(s) with price < 25~~

Q<sub>2</sub>: Find all companies that make only products with price < 25

≡ Q<sub>2</sub>: Find all companies for which all products have price < 25

≡ Q<sub>2</sub>: Find all companies that do not have any product with price  $\geq$  25

*Universal quantifiers are more complicated! ☹*

*(Think about the companies that should not be returned)*

*All three formulations are equivalent: a company with no product will be returned!*

# Correlated subquery (universal $\forall$ = not exists $\nexists$ )



Q<sub>2</sub>: Find all companies that make only products with price < 25

Step 1: Q<sub>2</sub>': Find the other companies that make some product(s) with price ≥ 25

```
SELECT DISTINCT C.cname
FROM   Company C
WHERE  C.cid IN      ( SELECT P.cid
                      FROM   Product P
                      WHERE  P.price >= 25)
```

*First think about the companies that should not be returned!*

Step 2: Q<sub>2</sub>: Find all companies that make no products with price ≥ 25

```
SELECT DISTINCT C.cname
FROM   Company C
WHERE  C.cid NOT IN ( SELECT P.cid
                      FROM   Product P
                      WHERE  P.price >= 25)
```

# Correlated subquery (universal $\forall$ = not exists $\nexists$ )



Q<sub>2</sub>: Find all companies that make only products with price < 25

Step 1: Q<sub>2</sub>': Find the other companies that make some product(s) with price ≥ 25

```
SELECT DISTINCT C.cname
FROM   Company C
WHERE  EXISTS      ( SELECT *
                    FROM   Product P
                    WHERE  C.cid = P.cid
                    and    P.price >= 25)
```

*First think about the companies that should not be returned!*

Step 2: Q<sub>2</sub>: Find all companies that make no products with price ≥ 25

```
SELECT DISTINCT C.cname
FROM   Company C
WHERE  NOT EXISTS ( SELECT *
                   FROM   Product P
                   WHERE  C.cid = P.cid
                   and    P.price >= 25)
```

# Correlated subquery (universal $\forall$ = not exists $\nexists$ )



Q<sub>2</sub>: Find all companies that make only products with price < 25

Step 1: Q<sub>2</sub>': Find the other companies that make some product(s) with price ≥ 25

```
SELECT DISTINCT C.cname
FROM   Company C
WHERE  25 <= ANY ( SELECT P.price
                  FROM   Product P
                  WHERE  C.cid = P.cid)
```

First think about the companies that should not be returned!

Step 2: Q<sub>2</sub>: Find all companies that make no products with price ≥ 25

```
SELECT DISTINCT C.cname
FROM   Company C
WHERE  25 > ALL ( SELECT P.price
                 FROM   Product P
                 WHERE  C.cid = P.cid)
```



# Correlated subquery (universal $\forall$ = not exists $\nexists$ )



Q<sub>2</sub>: Find all companies that make only products with price < 25

Step 1: Q<sub>2</sub>': Find the other companies that make some product(s) with price ≥ 25

```
SELECT DISTINCT C.cname
FROM   Company C
WHERE  25 <= ANY ( SELECT P.price
                  FROM   Product P
                  WHERE  C.cid = P.cid)
```

First think about the companies that should not be returned!

Step 2: Q<sub>2</sub>: Find all companies that make no products with price ≥ 25

```
SELECT DISTINCT C.cname
FROM   Company C
WHERE  25 > ALL ( SELECT P.price
                 FROM   Product P
                 WHERE  C.cid = P.cid)
```




# A natural question

Q<sub>2</sub>: Find all companies that make only products with price < 25

- How can we unnest (no GROUP BY) the universal quantifier query ?

```
SELECT ...  
FROM ...  
WHERE ...
```



# Topic 1: Data models and query languages

## Unit 1: SQL (continued)

### Lecture 3

Wolfgang Gatterbauer

CS7240 Principles of scalable data management (sp24)

<https://northeastern-datalab.github.io/cs7240/sp24/>

1/16/2024




# A natural question

Q<sub>2</sub>: Find all companies that make only products with price < 25

- How can we unnest (no GROUP BY) the universal quantifier query ?

```
SELECT ...  
FROM ...  
WHERE ...
```





# Pre-class conversations

- Last class summary
- New class members: next time quick introduction
- Please point out any errors on slides, web page, Canvas, etc.
- Any questions on class policies
  
- Today:
  - SQL continued
  - Recall: no class this FRI, we resume TUE in person

# CS 7240: Topics and approximate agenda (Spring'24)

This schedule will be updated regularly as the class progresses. Check back frequently. I will usually post lecture slides by the end of the day following a lecture (thus the next day), or latest two days after class. **Notice that I post one single slide deck for each unit (e.g. Topic 1 - Unit 1- SQL), and I keep those slide decks updated as we progress with the unit across lectures.** I post them here on this website (or in Canvas if I think they are not yet ready to be released in public). Please also check our [DATA lab seminar](#) for talks of interest.

## Topic 1: Data Models and Query Languages

- **Lecture 1 (Tue 1/9):** Course introduction / T1-U1 SQL / PostgreSQL setup / SQL Activities
- **Lecture 2 (Fri 1/12):** T1-U1 SQL T1-U1 SQL
- **Lecture 3 (Tue 1/16) via Zoom:** T1-U1 SQL
- **Lecture 4 (Fri 1/19):** no class
- **Lecture 5 (Tue 1/23):** T1-U2 Logic & Relational Calculus
- **Lecture 6 (Fri 1/26):** T1-U2 Logic & Relational Calculus
- **Lecture 7 (Tue 1/30):** T1-U3 Relational Algebra & Codd's Theorem
- **Lecture 8 (Fri 2/2):** T1-U3 Relational Algebra & Codd's Theorem
- **Lecture 9 (Tue 2/6):** T1-U4 Datalog & Recursion & ASP
- **Lecture 10 (Tue 2/9):** T1-U4 Datalog & Recursion & ASP
- **Lecture 11 (Tue 2/13):** T1-U4 Datalog & Recursion & ASP
- **Lecture 12 (Fri 2/16):** T1-U4 Datalog & Recursion & ASP

## Topic 2: Complexity of Query Evaluation & Reverse Data Management

- **Lecture 11 (Tue 2/14):** T2-U1 Conjunctive Queries
- **Lecture 12 (Fri 2/17):** T2-U1 Conjunctive Queries
- **Lecture 13 (Tue 2/21):** T2-U2 Beyond Conjunctive Queries
- **Lecture 14 (Fri 2/24):** T2-U3 Provenance
- **Lecture 15 (Tue 2/28):** T2-U3 Provenance
- **Lecture 16 (Fri 3/3):** T2-U4 Reverse Data Management

## Topic 3: Efficient Query Evaluation & Factorized Representations

- Spring break (Tue 3/7, Fri 3/10: [Northeast Database day 2023 @ Northeastern](#))
- **Lecture 17 (Tue 3/14):** T3-U1 Acyclic Queries
- **Lecture 18 (Fri 3/17):** T3-U1 Acyclic Queries
- **Lecture 19 (Tue 3/21):** T3-U2 Cyclic Queries
- **Lecture 20 (Fri 3/24):** T3-U2 Cyclic Queries
- **Lecture 21 (Tue 3/28):** T3-U2 Cyclic Queries
- **Lecture 22 (Fri 3/31):** T3-U2 Cyclic Queries
- **Lecture 23 (Tue 4/4):** T3-U3 Factorized Representations
- **Lecture 24 (Fri 4/7):** T3-U4 Optimization Problems & Top-k
- **Lecture 25 (Tue 4/11):** T3-U4 Optimization Problems & Top-k

## Topic 4: Normalization, Information Theory & Axioms for Uncertainty

- **Lecture:** Normal Forms & Information Theory
- **Lecture:** Axioms for Uncertainty

## Topic 5: Linear Algebra & Iterative Graph Algorithms

- **Lecture:** Graphs & Linear Algebra
- **Lecture:** Computation Graphs

## Project presentations

- **Lecture 26 (Fri 4/14):** P4 Project presentations
- **Lecture 27 (Tue 4/18):** P4 Project presentations

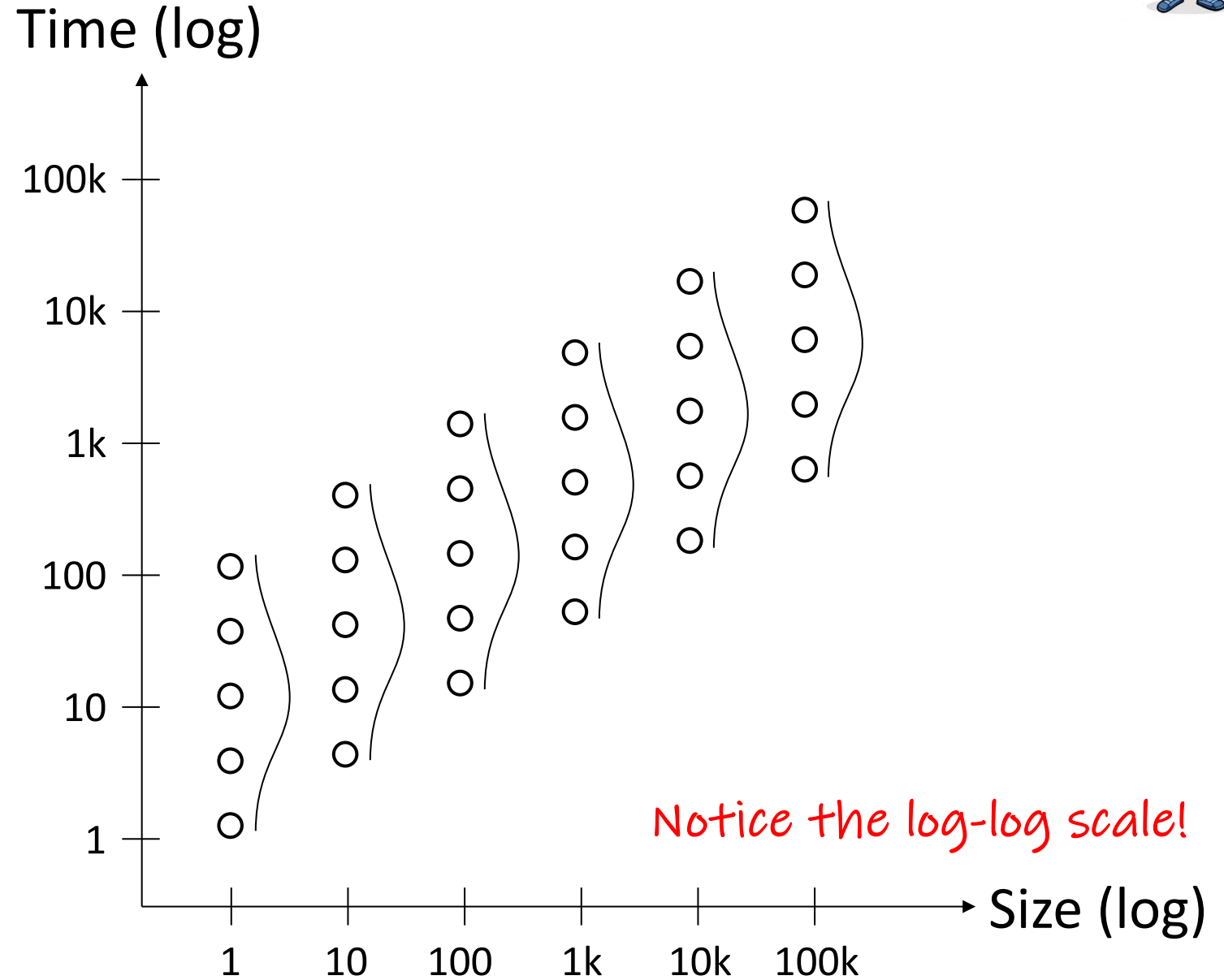
Revisiting our question  
from first class



# Question: How to deal with cut-offs when binning?



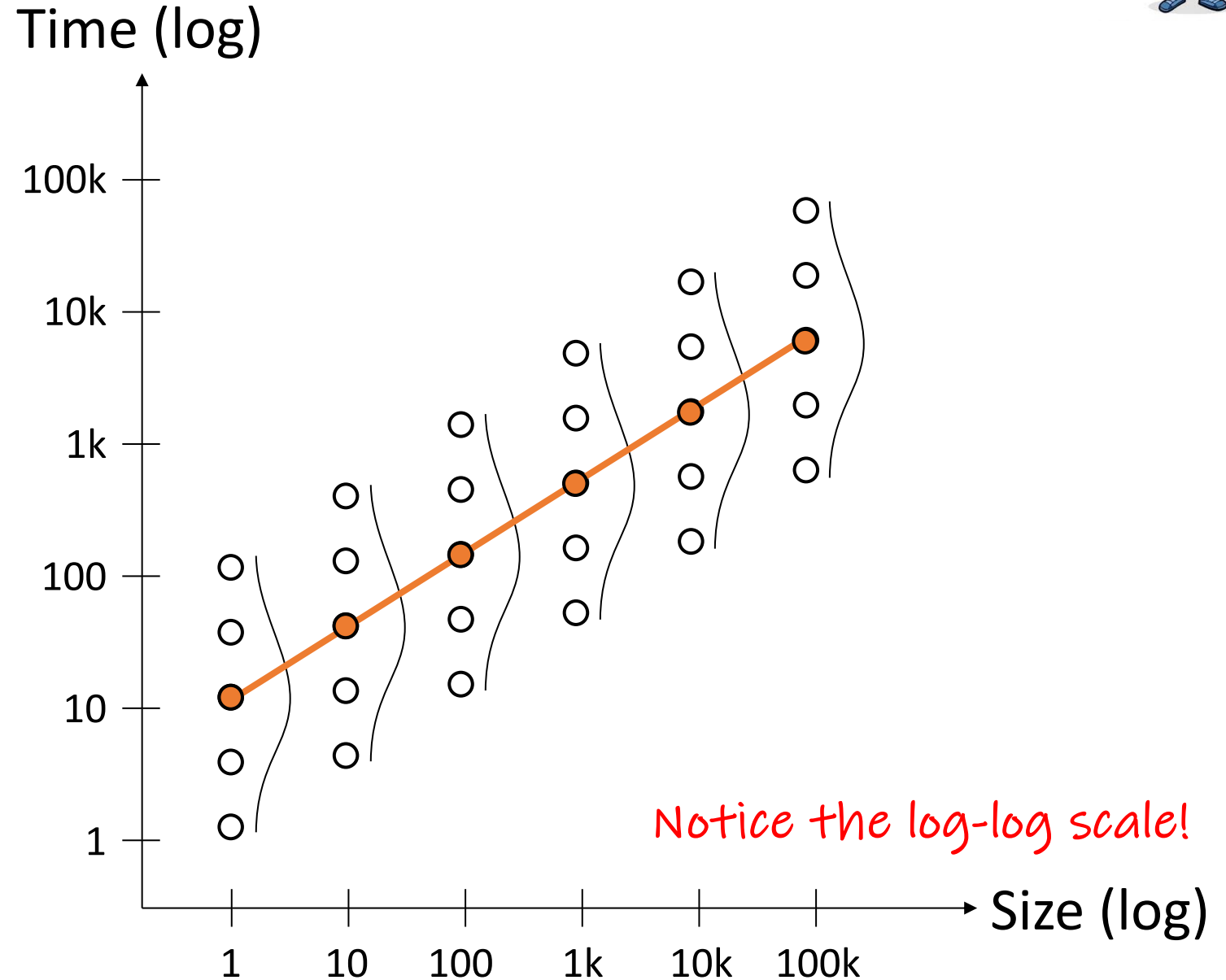
- These are the true points that you would get if you could run the experiments long enough.



# Question: How to deal with cut-offs when binning?



- These are the true points that you would get if you could run the experiments long enough.
- Here is what the aggregate would look like if we could get all points and then aggregated for each size

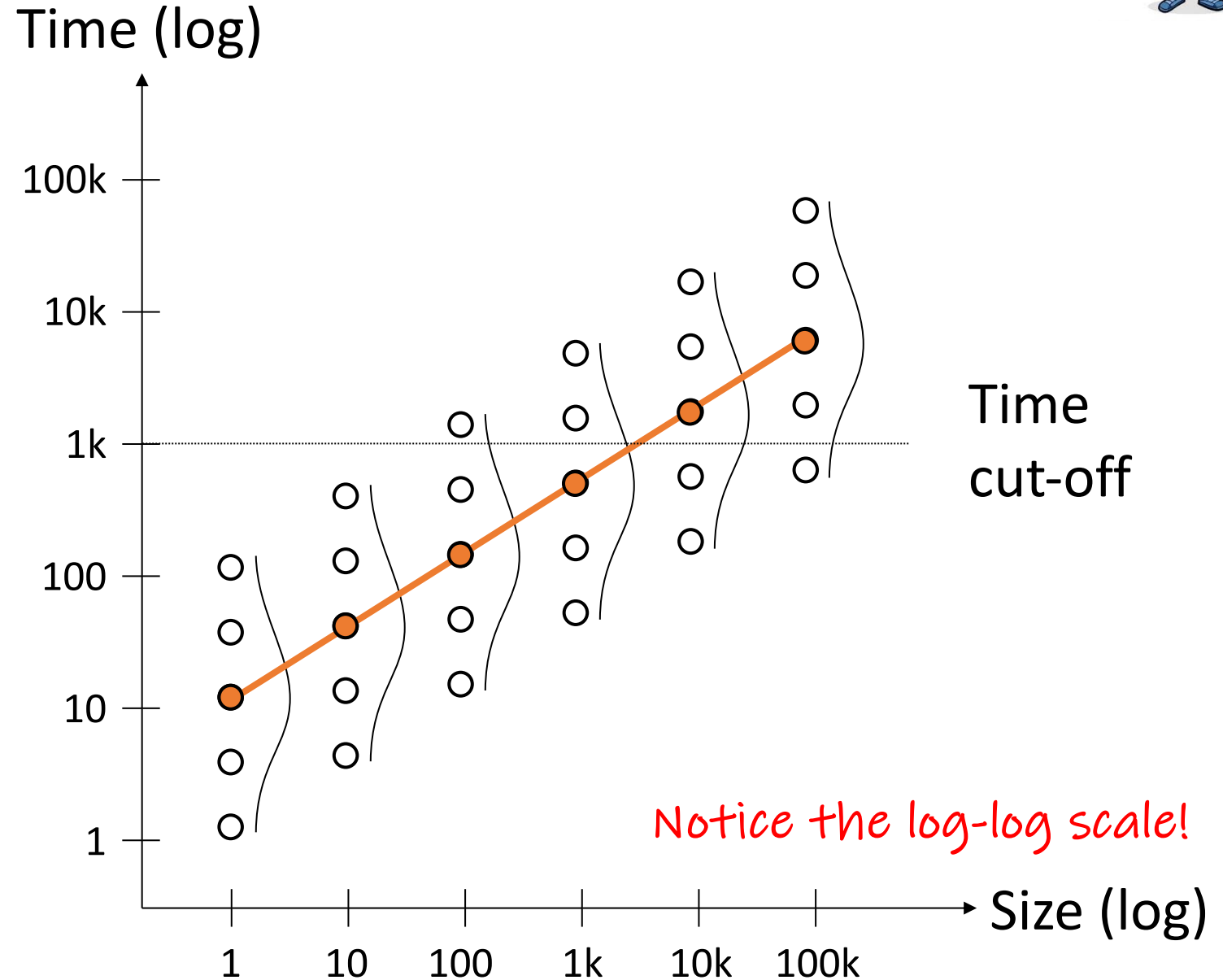


# Question: How to deal with cut-offs when binning?



- These are the true points that you would get if you could run the experiments long enough.
- Here is what the aggregate would look like if we could get all points and then aggregated for each size
- However, some experiments take too long and we thus have to cut them off after some time.

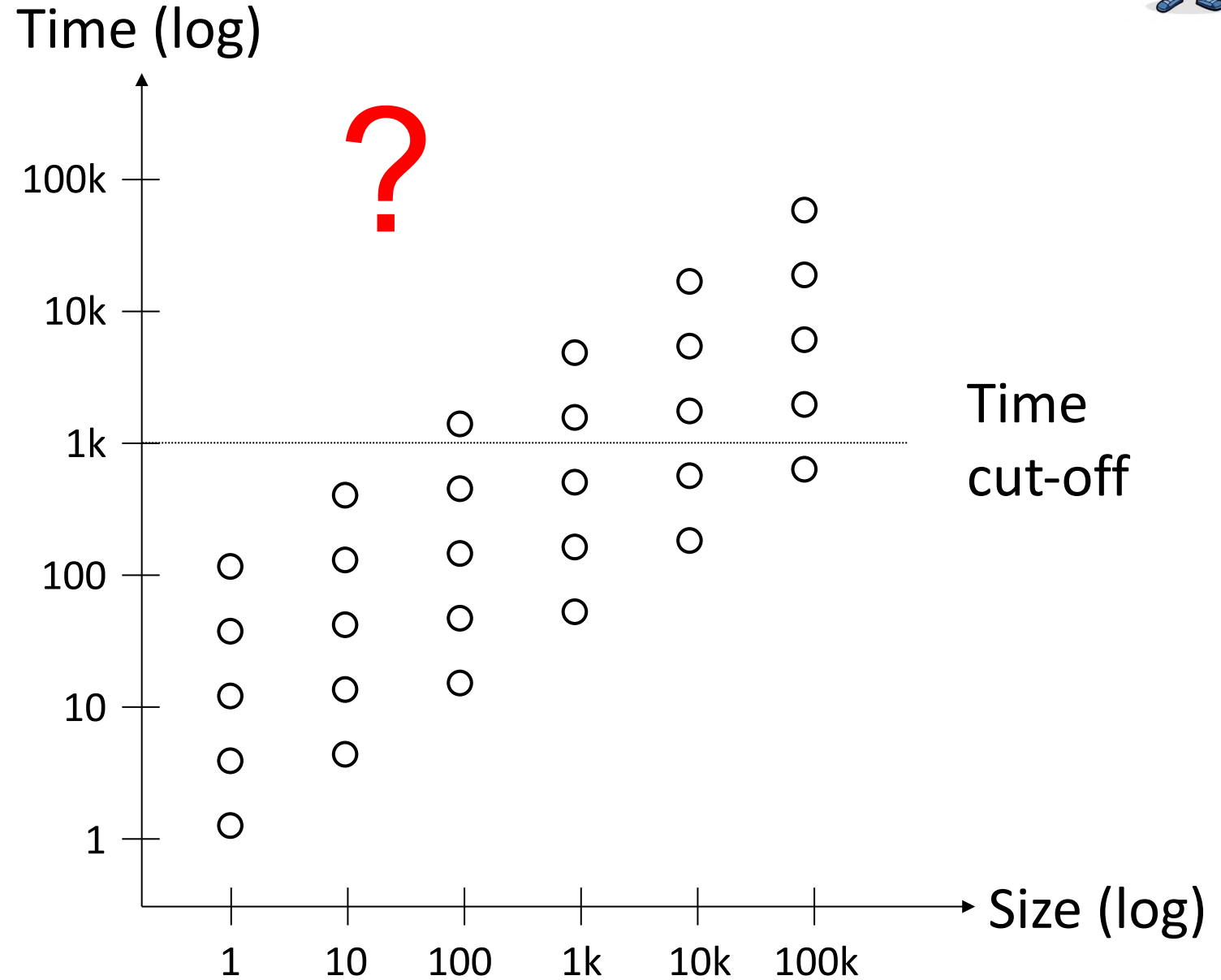
*Question: There is an overall trend, yet big variation for each experiment. We still like to capture the overall trend with some smart aggregations. What can we do?*



# Question: How to deal with cut-offs when binning?



- **Option 1:** What if we terminate experiments that take too long, and only average over the "seen points"?

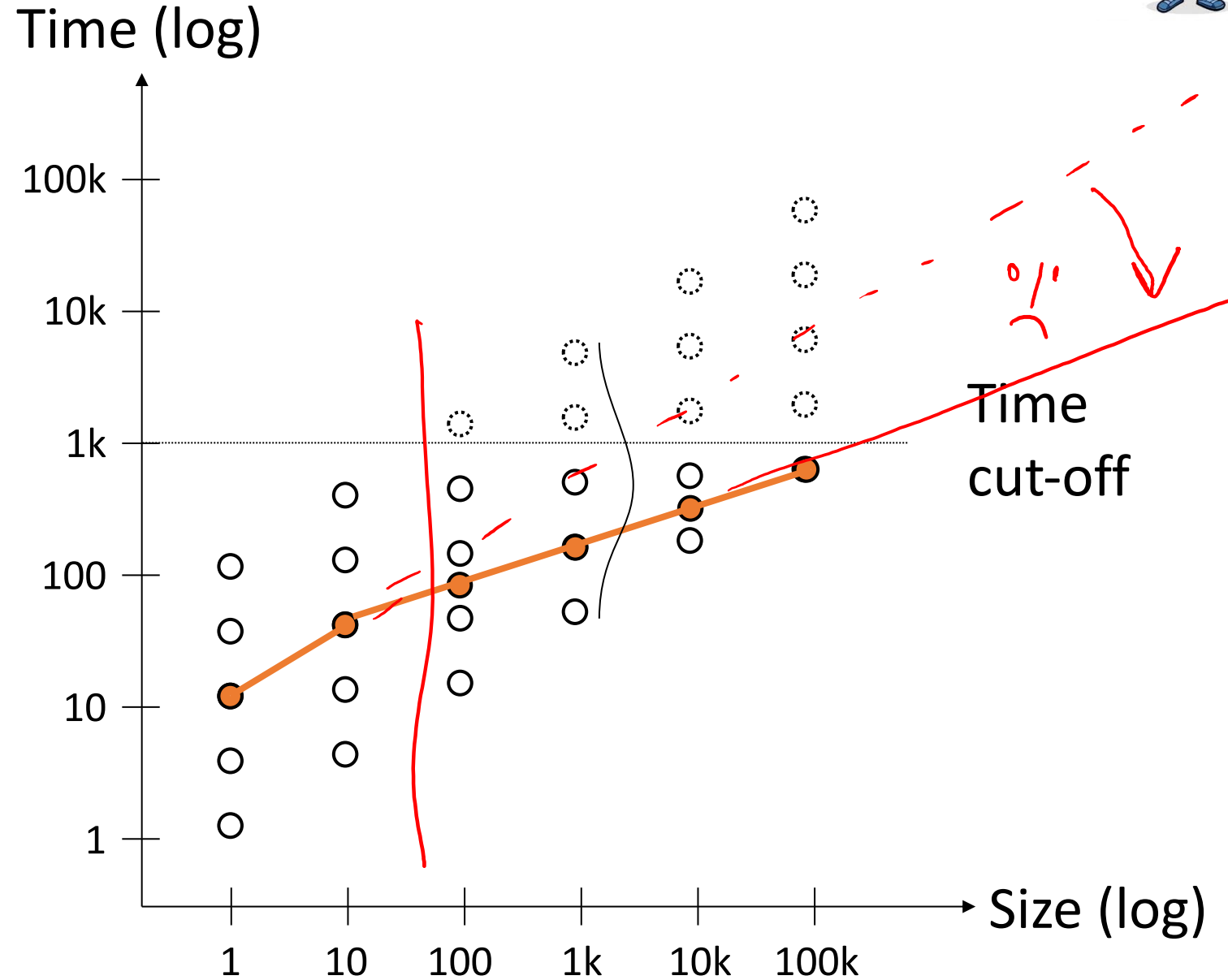


# Question: How to deal with cut-offs when binning



- **Option 1:** What if we terminate experiments that take too long, and only average over the "seen points"?

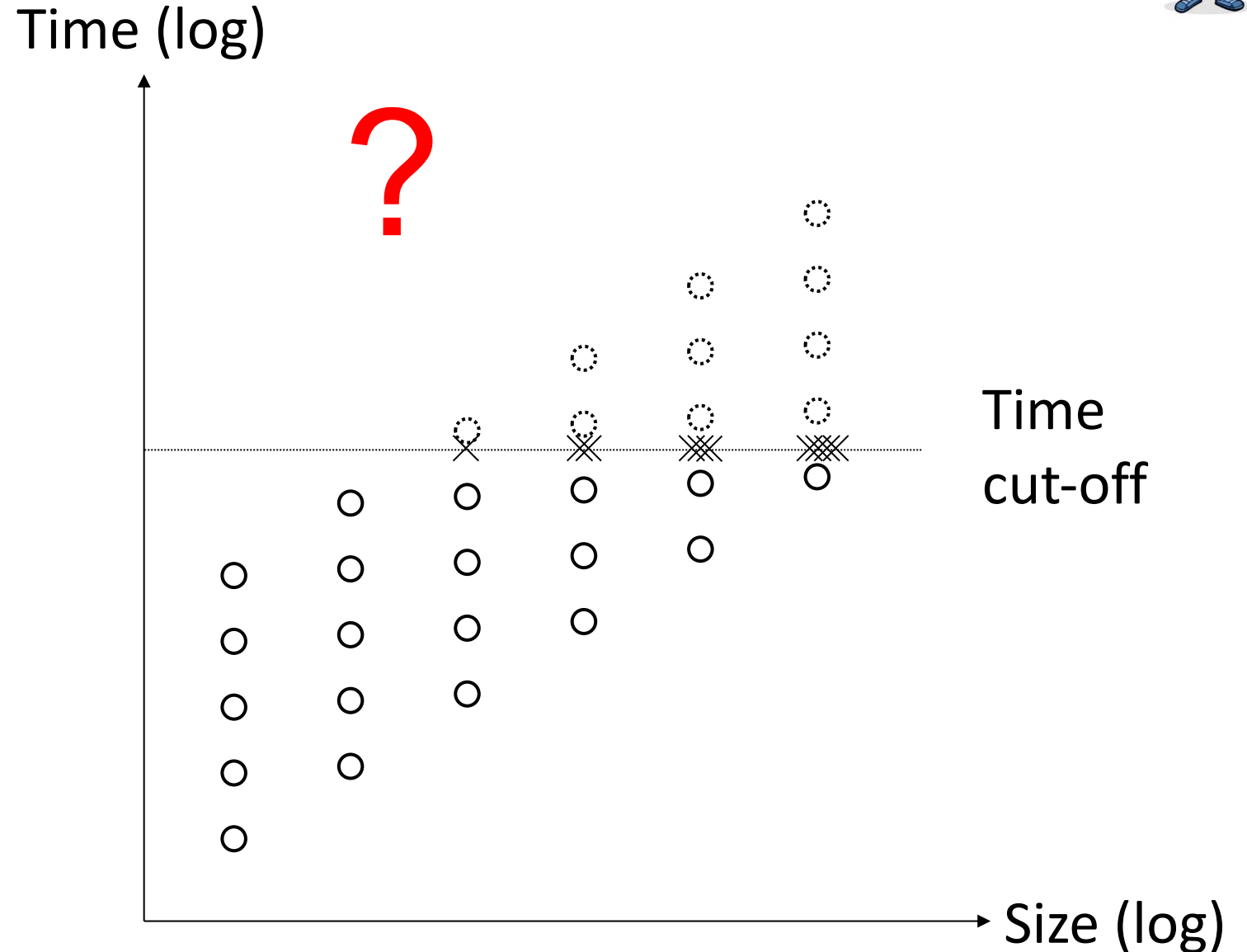
What would you do ?  
We will discuss next class



# How to deal with cut-offs when binning: Option 1



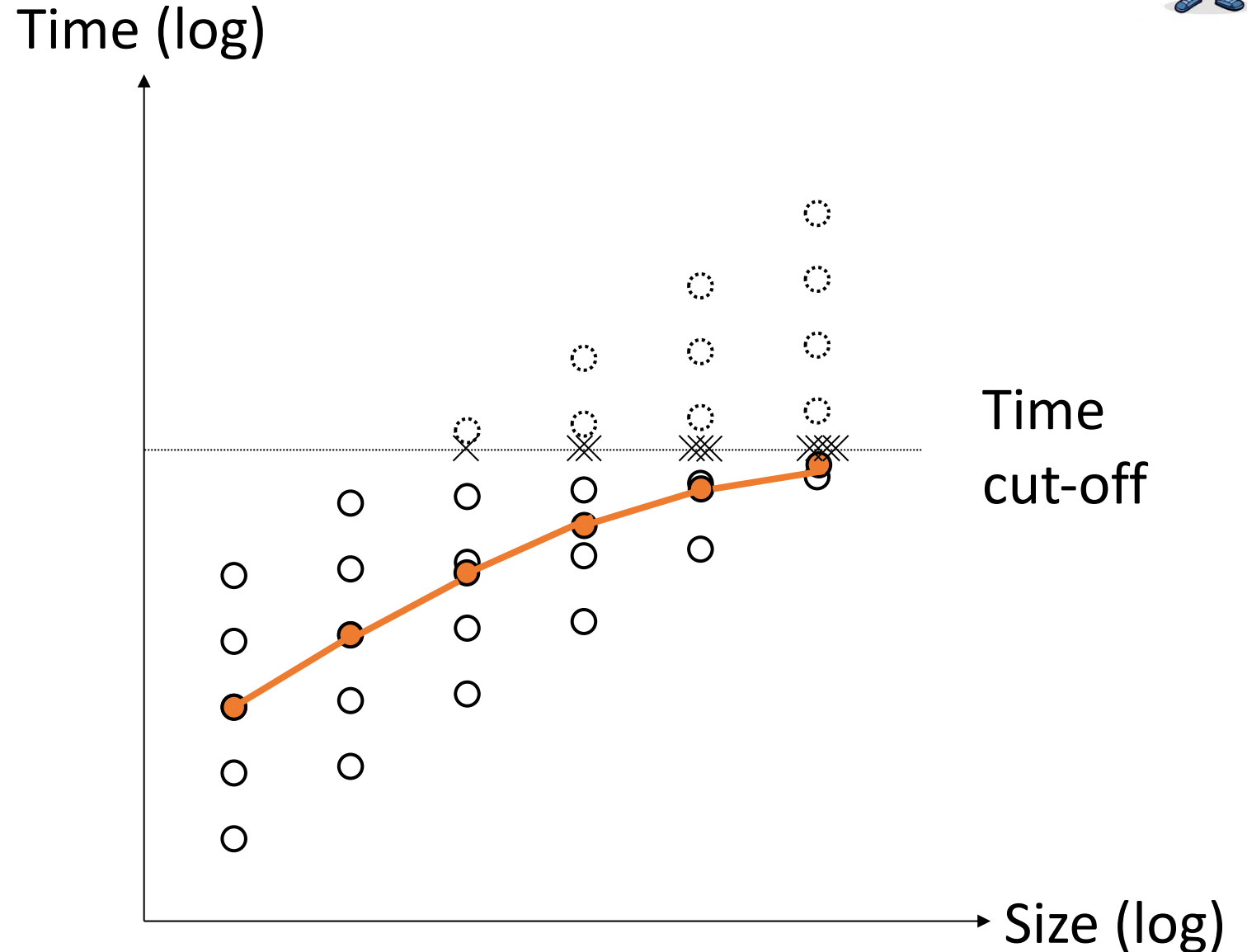
- **Option 1:** What if we terminate experiments that take too long, and only average over the "seen points"?
- **Option 2:** What if we cut the points off and still use the cut-off points, and then average?



# How to deal with cut-offs when binning: Option 1



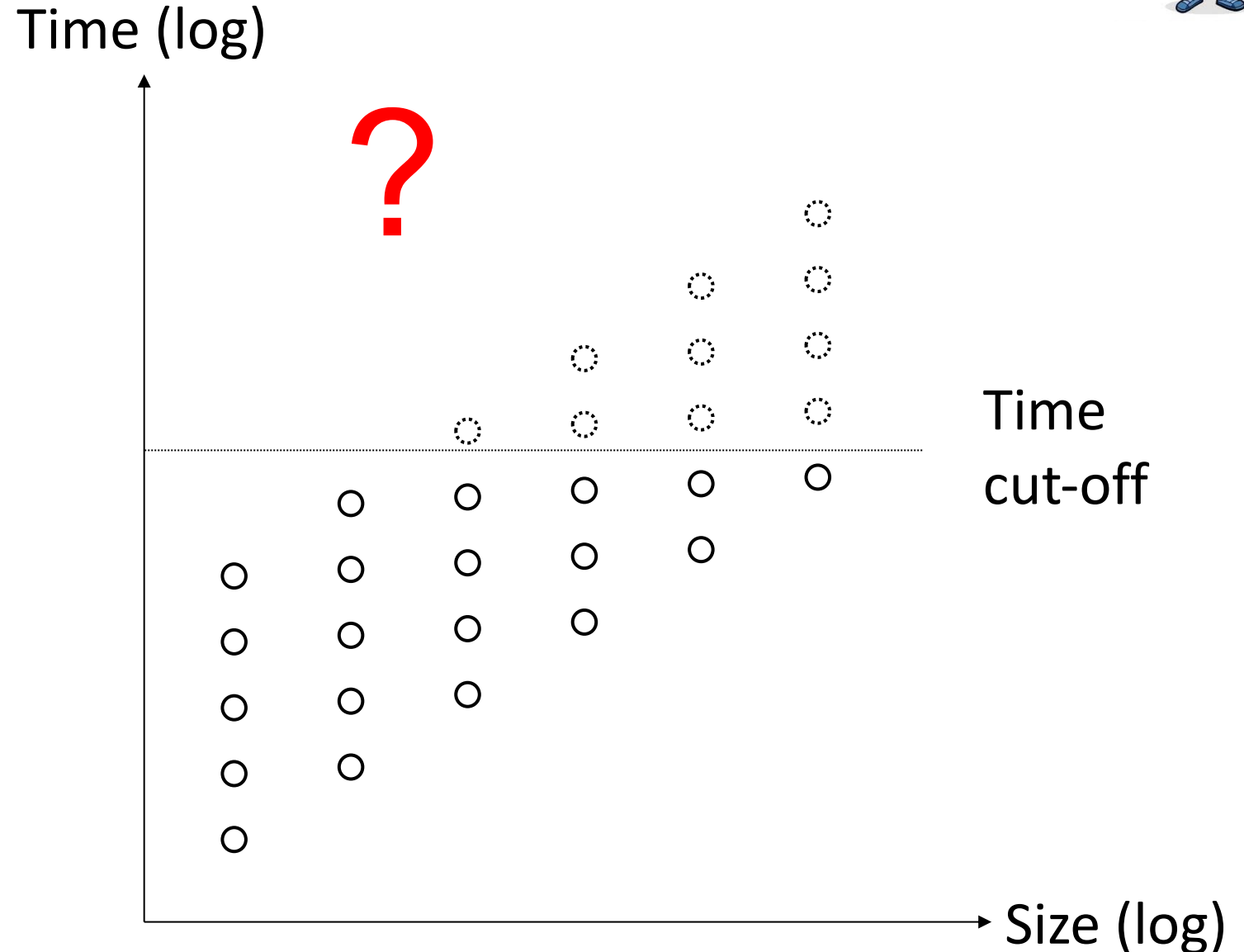
- **Option 1:** What if we terminate experiments that take too long, and only average over the "seen points"?
- **Option 2:** What if we cut the points off and still use the cut-off points, and then average?



# How to deal with cut-offs when binning: Option 2



- **Option 1:** What if we terminate experiments that take too long, and only average over the "seen points"?
- **Option 2:** What if we cut the points off and still use the cut-off points, and then average?
- **Option 3:** What if we \*only\* use those sizes (x-axis) for which all experiments finish in time?

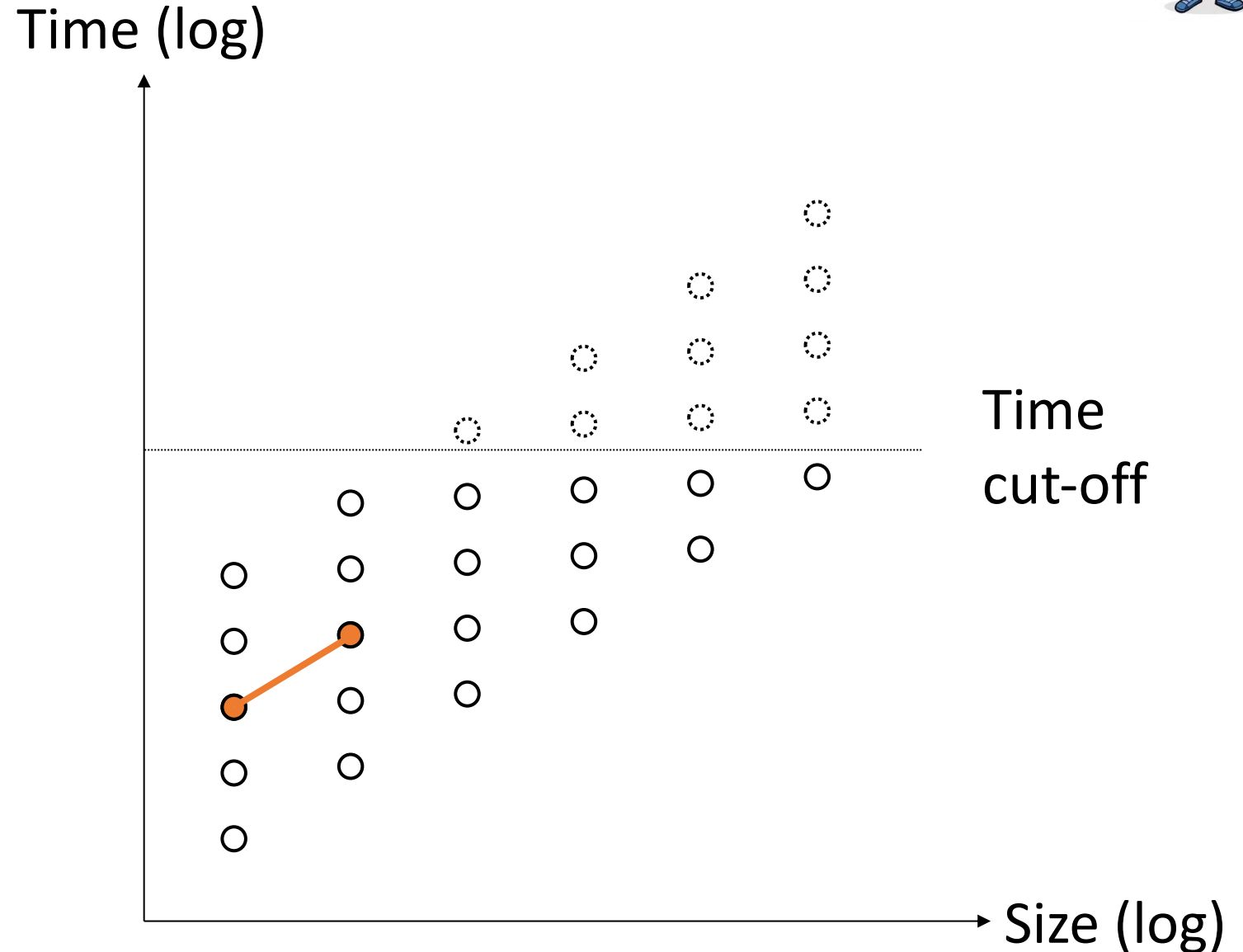




# How to deal with cut-offs when binning: Option 2



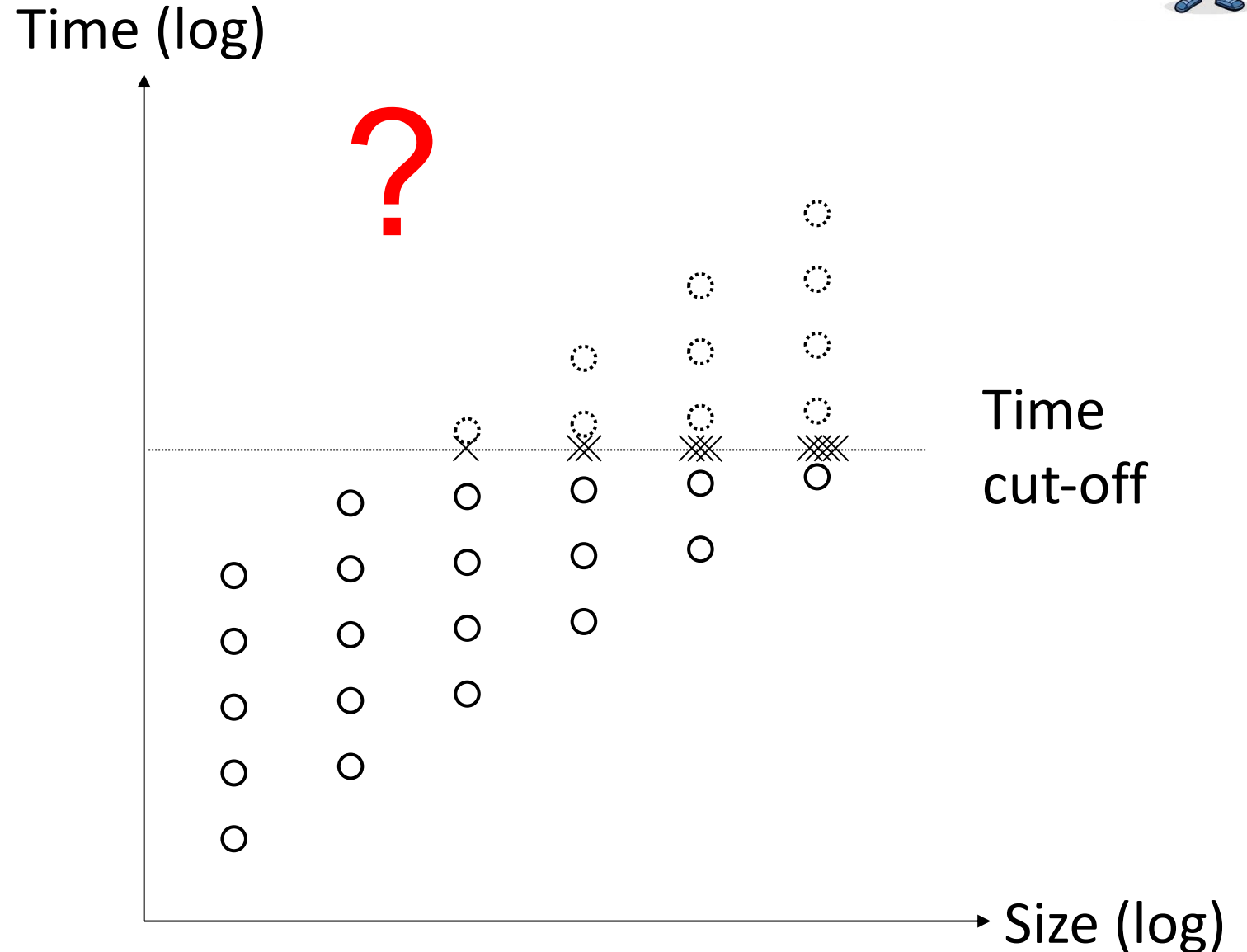
- **Option 1:** What if we terminate experiments that take too long, and only average over the "seen points"?
- **Option 2:** What if we cut the points off and still use the cut-off points, and then average?
- **Option 3:** What if we \*only\* use those sizes (x-axis) for which all experiments finish in time?



# How to deal with cut-offs when binning: Option 3



- **Option 1:** What if we terminate experiments that take too long, and only average over the "seen points"?
- **Option 2:** What if we cut the points off and still use the cut-off points, and then average?
- **Option 3:** What if we \*only\* use those sizes (x-axis) for which all experiments finish in time?
- **Option 4:** What if we take the median over all seen and cut-off points?



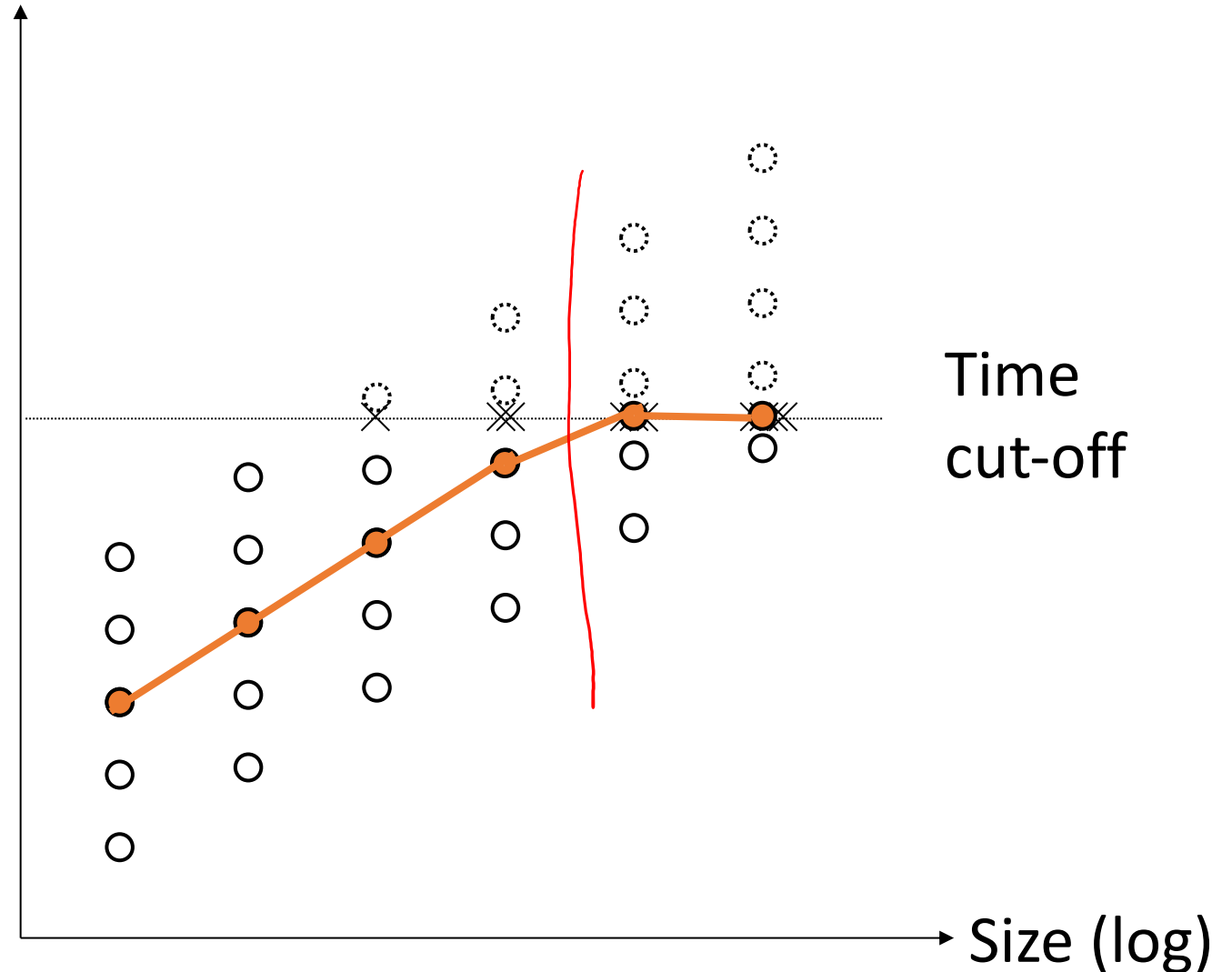
# How to deal with cut-offs when binning: Option 3



- **Option 1:** What if we terminate experiments that take too long, and only average over the "seen points"?
- **Option 2:** What if we cut the points off and still use the cut-off points, and then average?
- **Option 3:** What if we \*only\* use those sizes (x-axis) for which all experiments finish in time?
- **Option 4:** What if we take the median over all seen and cut-off points?

*Notice the informal "semantics" of median: If more points are "above you" then you are pulled by their number, not by their distance (in contrast to average where distance is kind of a weight)*

Time (log)

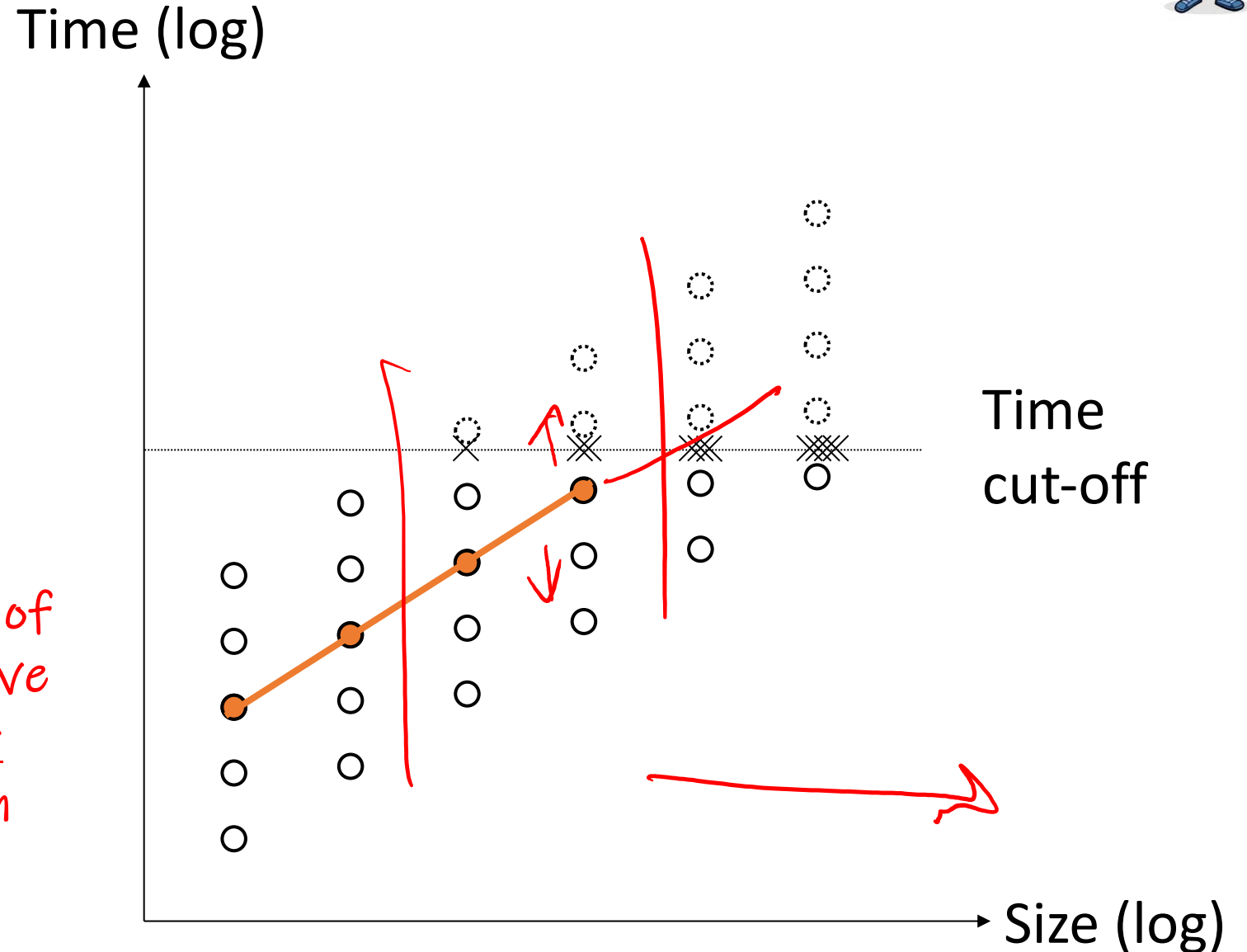


# How to deal with cut-offs when binning: Suggestion

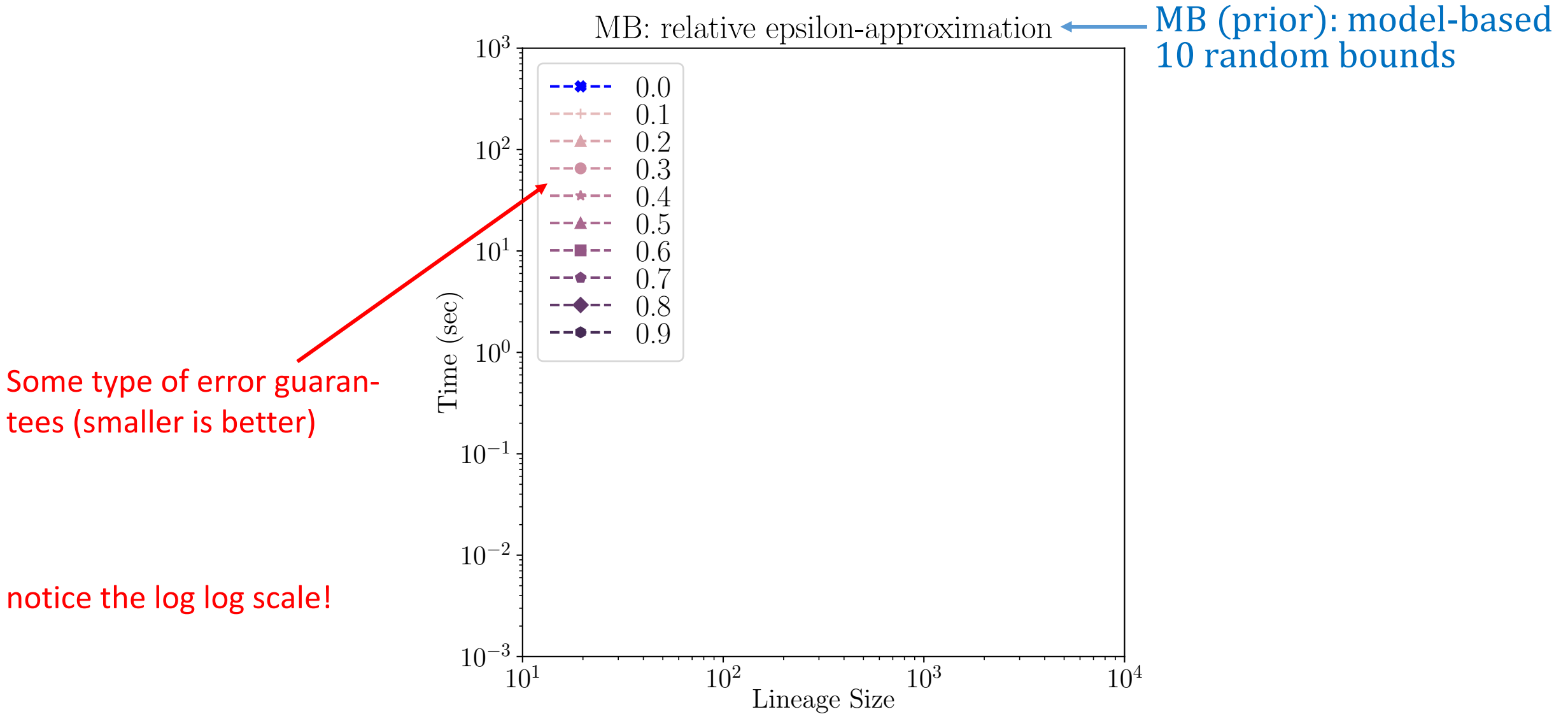


- **Suggestion:** Take the median over all seen and cut-off points, as long as there are <50% cut-off points!

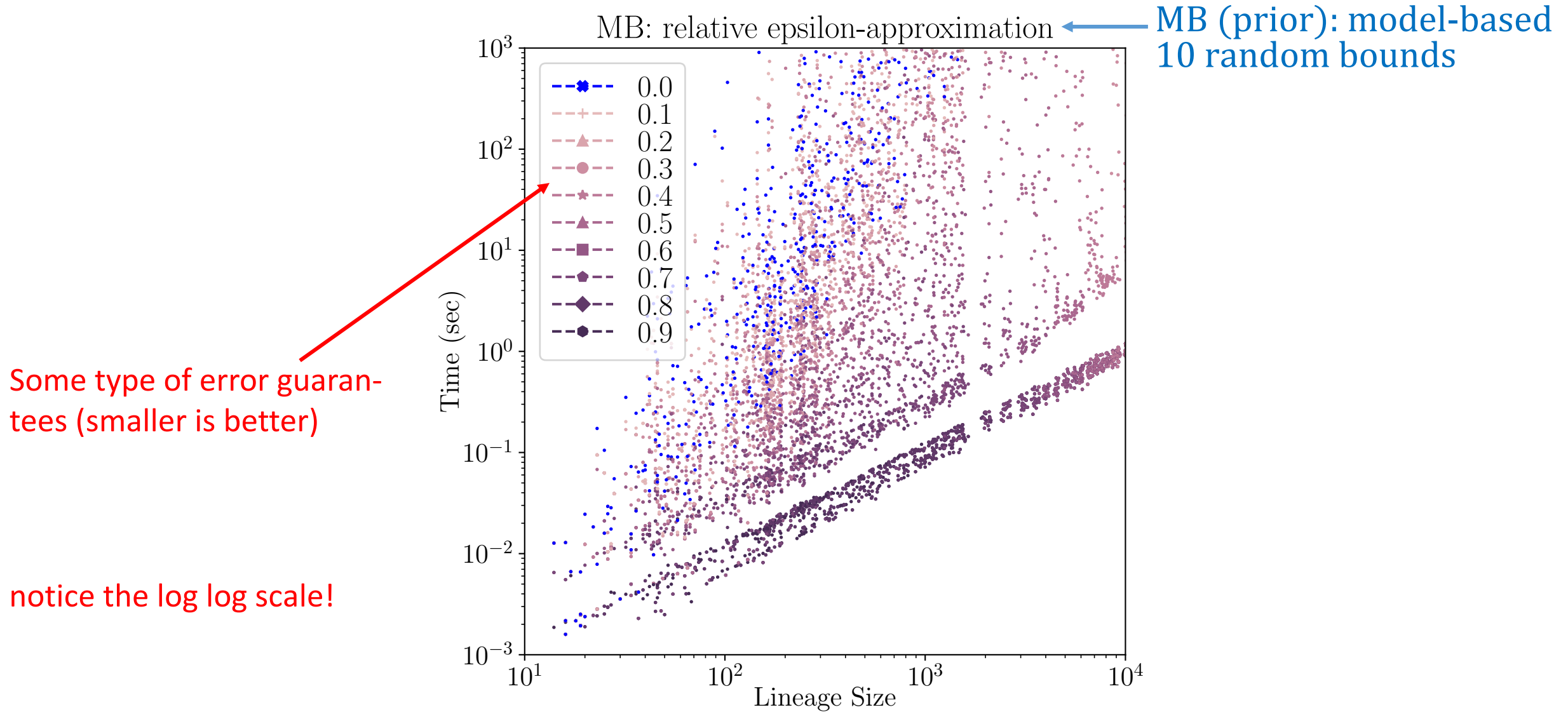
Notice the informal "semantics" of median: If more points are "above you" then you are pulled by their number, not by their distance (in contrast to average where distance is kind of a weight)



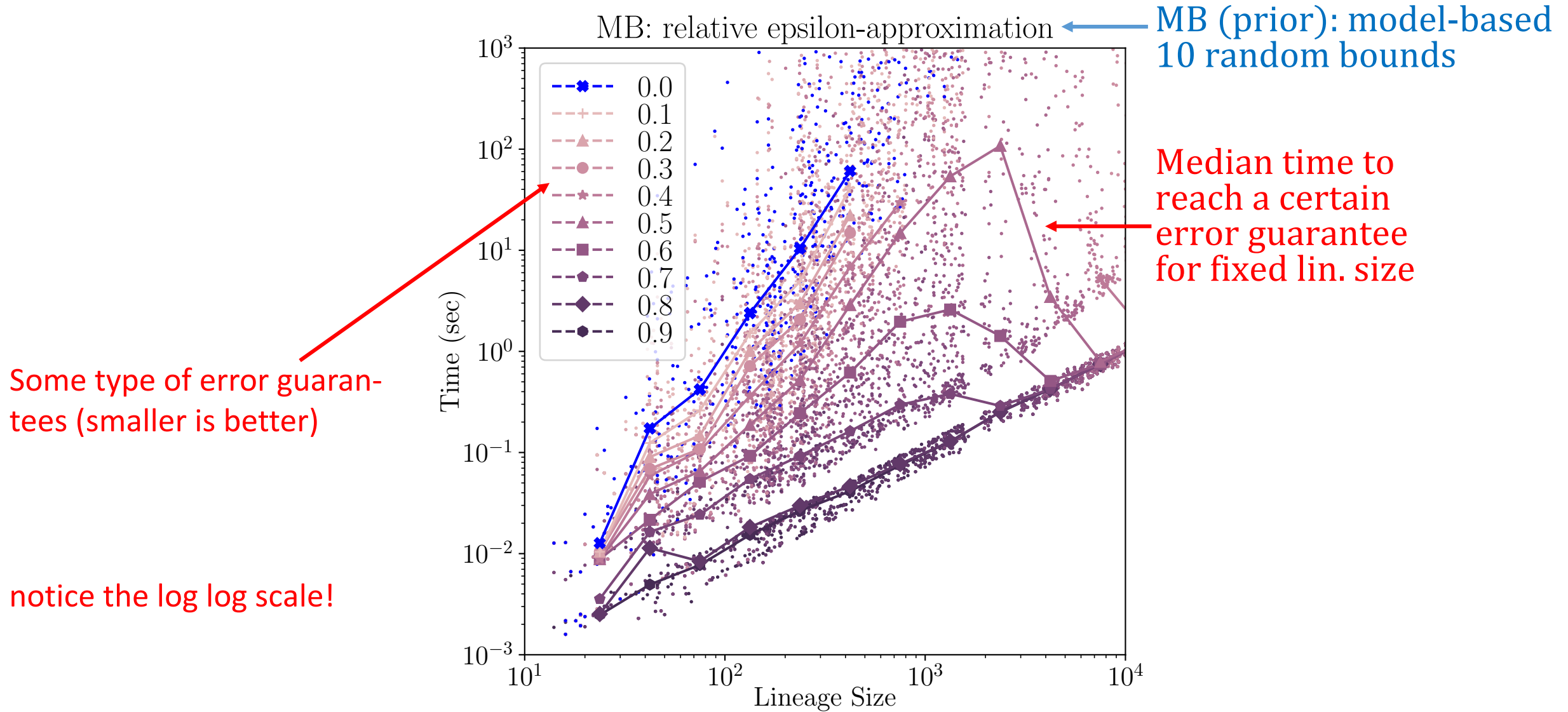
# Example: Experiments figures from [SIGMOD'19]



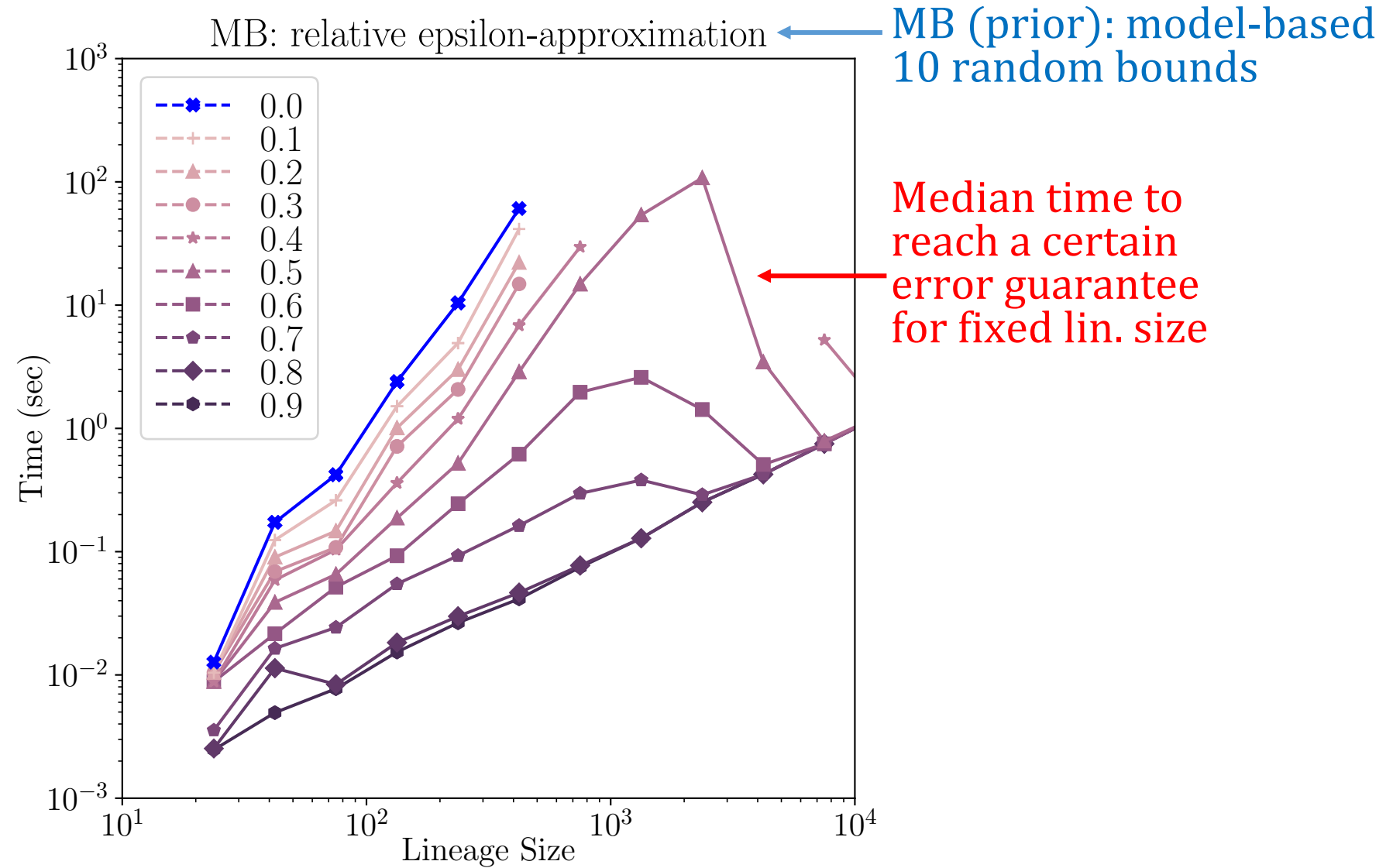
# Example: Experiments figures from [SIGMOD'19]



# Example: Experiments figures from [SIGMOD'19]

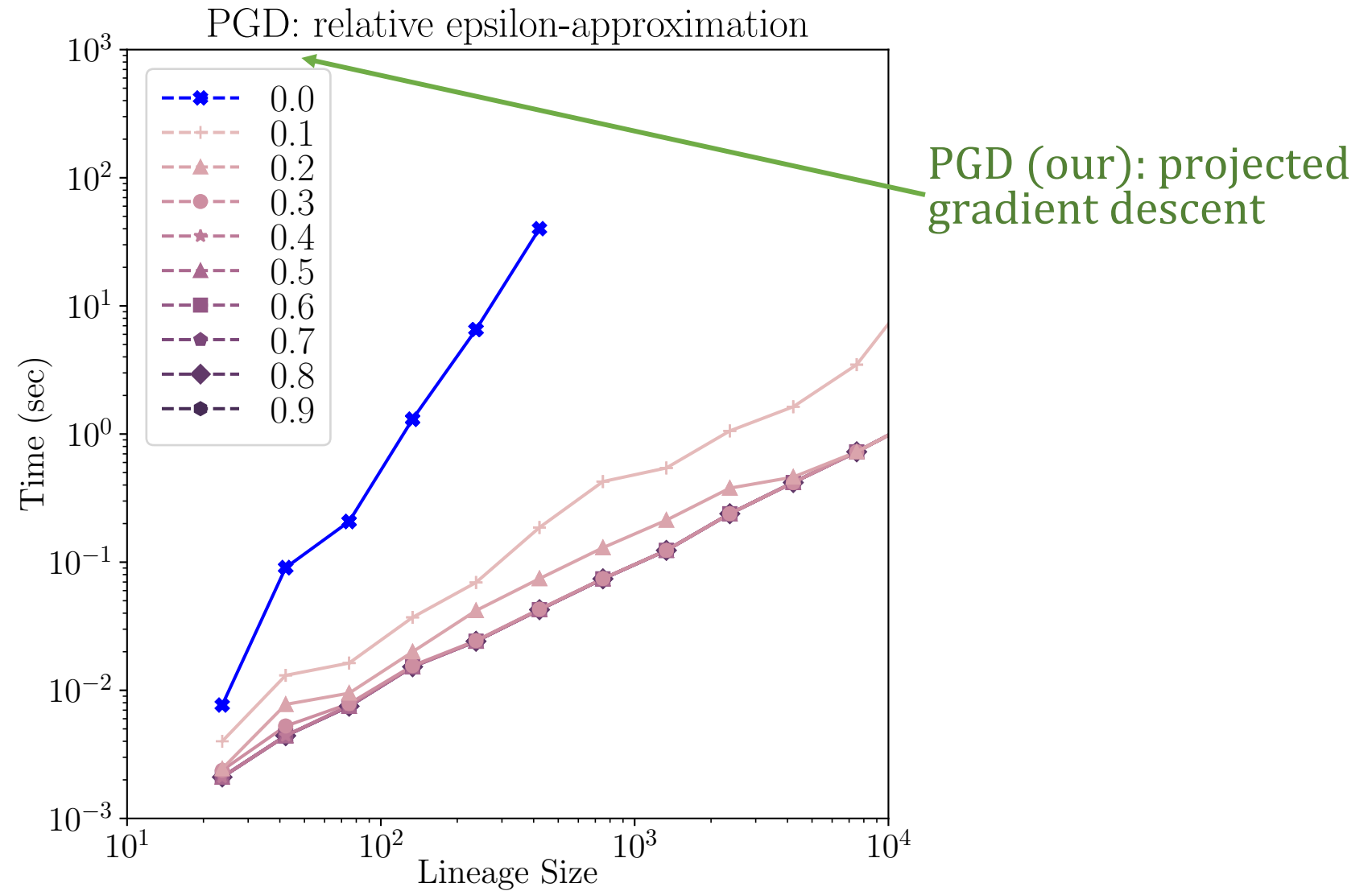


# Example: Experiments figures from [SIGMOD'19]

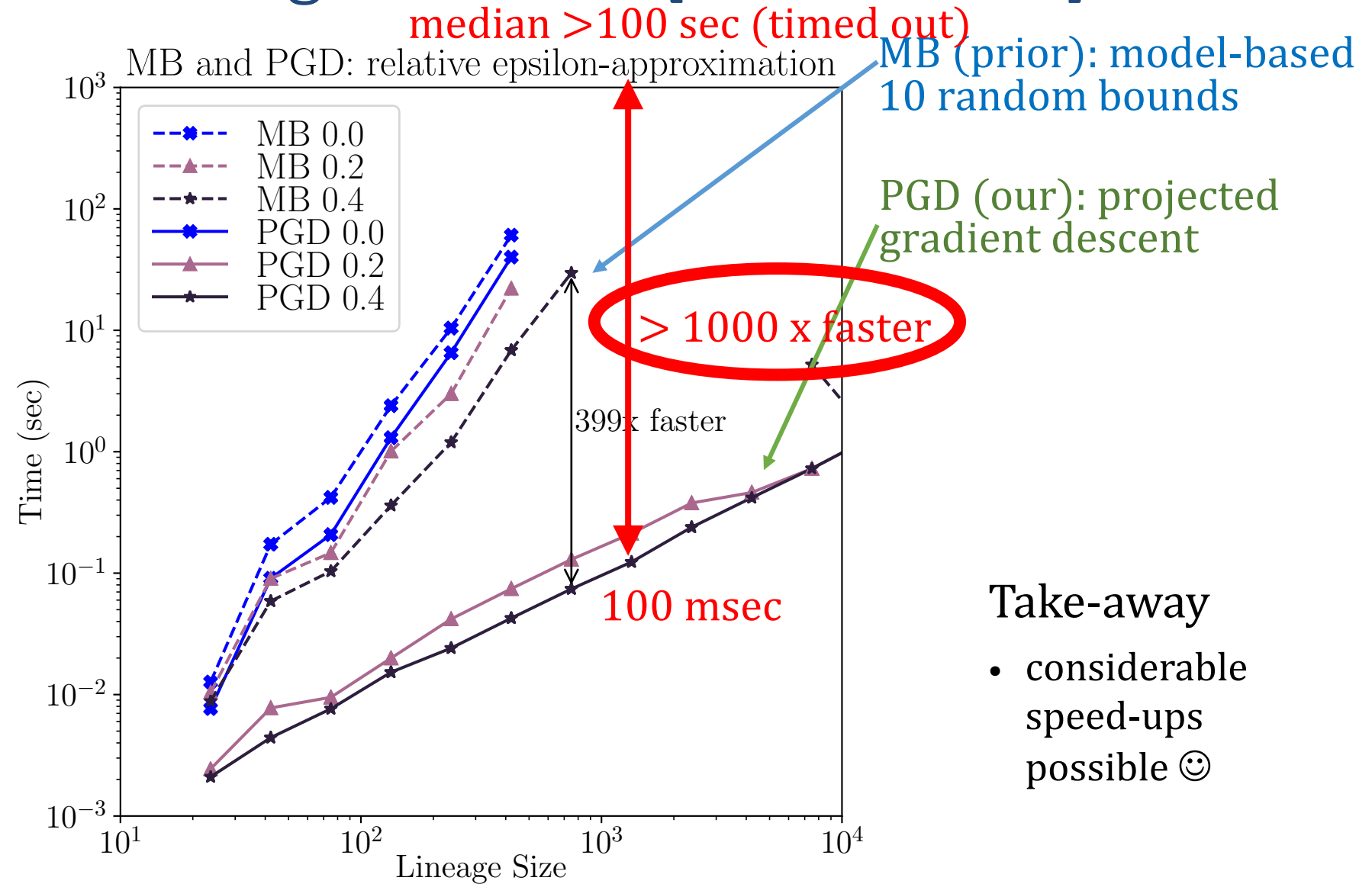




# Example: Experiments figures from [SIGMOD'19]



# Example: Experiments figures from [SIGMOD'19]



## Take-away

- considerable speed-ups possible 😊

# Viewpoint

## The End of Programming

*The end of classical computer science is coming, and most of us are dinosaurs waiting for the meteor to hit.*

Why do I think we should care about experimental setups, even in theory!

cations, most software, as we know it, will be replaced by AI systems that are *trained* rather than *programmed*. In situations

tension, I believe the computer scientists of the future will be so far removed from the classic definitions of “software” that they would be hard-pressed to reverse a linked list or implement Quicksort. (I am

mers.<sup>1</sup> I am talking about *replacing the entire concept of writing programs with training models*. In the future, CS stu-

will be about coming up with the right examples, the right training data, and the right ways to evaluate the training process. Suitably powerful models capable

lignant AI running amok. We currently have no way, apart from empirical study, to determine the limits of current AI systems. As for future AI models that are or-

tions; the building blocks of AI systems are much higher-level abstractions like attention layers, tokenizers, and datasets. A time traveler from even 20 years ago

decidability. AI-based computation has long since crossed the Rubicon of being amenable to static analysis and formal proof. We are rapidly moving toward a world where the fundamental building blocks of computation are temperamental, mysterious, adaptive agents.

*large AI models work*. People are publishing research papers<sup>3-5</sup> actually *discovering new behaviors* of existing large models, even though these systems have been “engineered” by humans. Large AI models are capable of doing things that they have not been explicitly trained to do,

# Outline: T1-U1: SQL

- SQL

- Schema, keys, referential integrity
- Joins
- Aggregates and grouping
- **Nested queries (Subqueries)**
- Union and Theta Joins
- Nulls & Outer joins
- Top-k
- [Recursion: moved to T1-U4: Datalog]

# Understanding nested queries with "Relational Diagrams"

# The sailors database

Sailor (sid, sname, rating, age)  
Reserves (sid, bid, day)  
Boat (bid, bname, color)



340

**Sailor**

<u>sid</u>	sname	rating	age
22	Dustin	7	45.0
29	Brutus	1	33.0
31	Lubber	8	55.5
32	Andy	8	25.5
58	Rusty	10	35.0
64	Horatio	7	35.0
71	Zorba	10	16.0
74	Horatio	9	35.0
85	Art	3	25.5
95	Bob	3	63.5

Figure 5.1 An Instance *S3* of Sailors

**Reserves**

<u>sid</u>	<u>bid</u>	<u>day</u>
22	101	10/10/98
22	102	10/10/98
22	103	10/8/98
22	104	10/7/98
31	102	11/10/98
31	103	11/6/98
31	104	11/12/98
64	101	9/5/98
64	102	9/8/98
74	103	9/8/98

Figure 5.2 An Instance *R2* of Reserves

**Boat**

<u>bid</u>	bname	color
101	Interlake	blue
102	Interlake	red
103	Clipper	green
104	Marine	red

Figure 5.3 An Instance *B1* of Boats

# Nested query 1



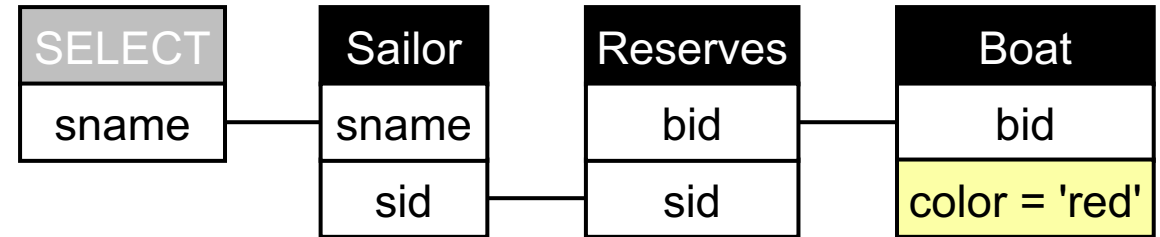
Q:

```
SELECT DISTINCT S.sname
FROM Sailor S
WHERE S.sid IN
  (SELECT R.sid
   FROM Reserves R
   WHERE R.bid IN
     (SELECT B.bid
      FROM Boat B
      WHERE B.color = 'red'))
```

Sailor (sid, sname, rating, age)  
Reserves (sid, bid, day)  
Boat (bid, bname, color)



340





# Nested query 1

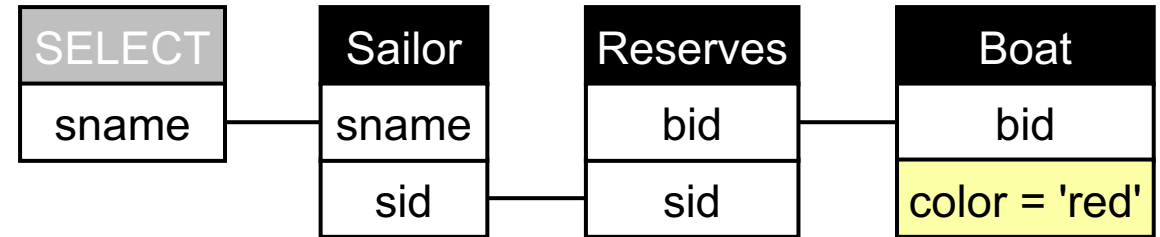
Sailor (sid, sname, rating, age)  
Reserves (sid, bid, day)  
Boat (bid, bname, color)



340

Q: Find the names of sailors who have reserved a red boat.

```
SELECT DISTINCT S.sname
FROM Sailor S
WHERE S.sid IN
  (SELECT R.sid
   FROM Reserves R
   WHERE R.bid IN
     (SELECT B.bid
      FROM Boat B
      WHERE B.color = 'red'))
```



$\{S.sname \mid \exists S \in \text{Sailor}.(\exists R \in \text{Reserves}.(R.sid=S.sid \wedge \exists B \in \text{Boat}.(B.bid=R.bid \wedge B.color='red'))))\}$



# Nested query 1

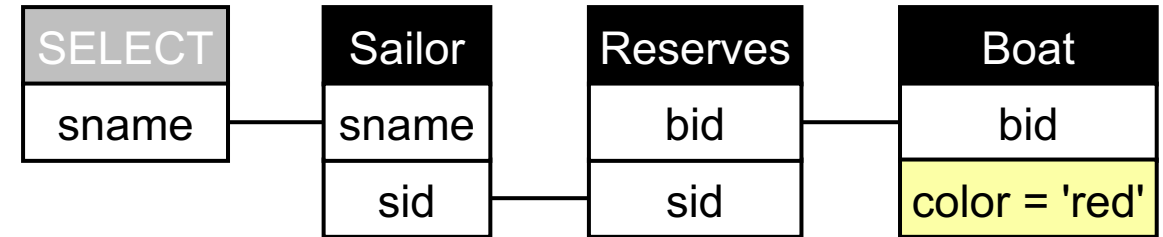
Sailor (sid, sname, rating, age)  
Reserves (sid, bid, day)  
Boat (bid, bname, color)



340

Q: Find sailors who have reserved a red boat.

```
SELECT DISTINCT S.sname
FROM Sailor S
WHERE EXISTS
  (SELECT R.sid
   FROM Reserves R
   WHERE R.sid=S.sid
   AND EXISTS
     (SELECT B.bid
      FROM Boat B
      WHERE B.bid = R.bid
      AND B.color = 'red'))
```



*This is an alternative way to write the previous query with EXISTS and correlated nested queries that matches the Relational Calculus below.*

$\{S.sname \mid \exists S \in \text{Sailor}.(\exists R \in \text{Reserves}.(R.sid=S.sid \wedge \exists B \in \text{Boat}.(B.bid=R.bid \wedge B.color='red')))\}$

# Nested query 2

Sailor (sid, sname, rating, age)  
Reserves (sid, bid, day)  
Boat (bid, bname, color)

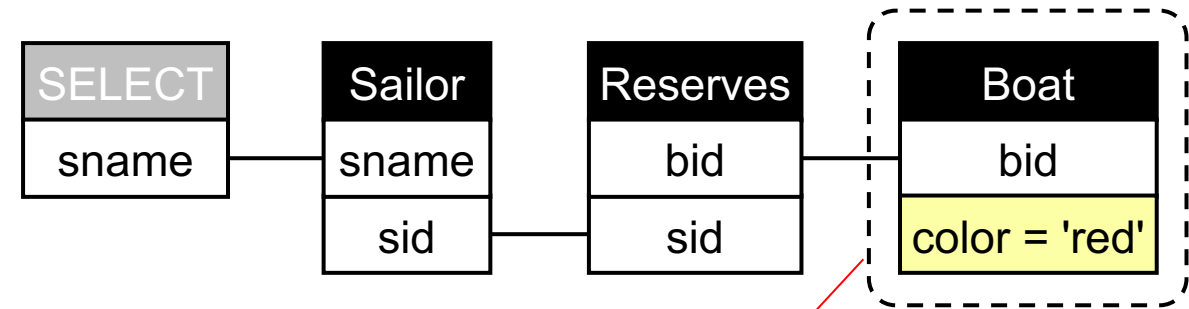


340



Q:

```
SELECT DISTINCT S.sname
FROM Sailor S
WHERE EXISTS
  (SELECT R.sid
   FROM Reserves R
   WHERE R.sid=S.sid
   AND NOT EXISTS
     (SELECT B.bid
      FROM Boat B
      WHERE B.bid = R.bid
      AND B.color = 'red'))
```



Dashed lines represent  
not exists ~~A~~

{S.sname |  $\exists S \in \text{Sailor}.(\exists R \in \text{Reserves}.(R.\text{sid}=S.\text{sid} \wedge \nexists B \in \text{Boat}.(B.\text{bid}=R.\text{bid} \wedge B.\text{color}='red')))$ }

# Nested query 2

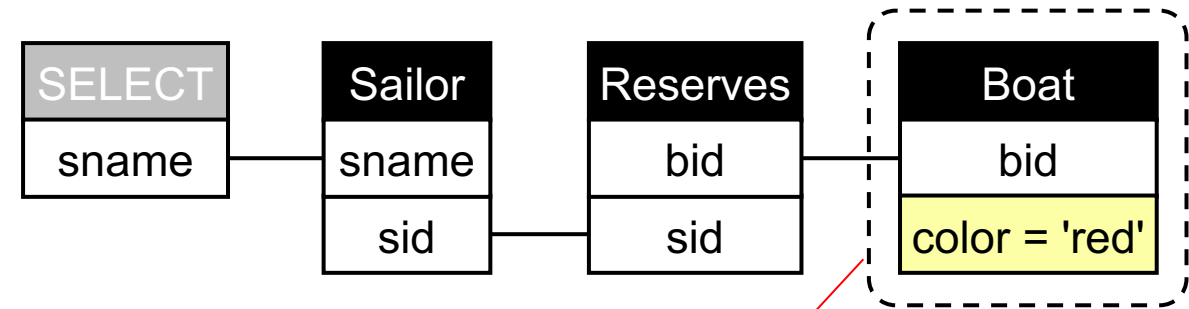
Sailor (sid, sname, rating, age)  
Reserves (sid, bid, day)  
Boat (bid, bname, color)



340

Q: Find sailors who have reserved a boat **that is not red**.

```
SELECT DISTINCT S.sname
FROM Sailor S
WHERE EXISTS
  (SELECT R.sid
   FROM Reserves R
   WHERE R.sid=S.sid
   AND NOT EXISTS
     (SELECT B.bid
      FROM Boat B
      WHERE B.bid = R.bid
      AND B.color = 'red'))
```



Dashed lines represent  
not exists ~~∃~~

They must have reserved at least one boat  
in another color. They can also have  
reserved a red boat in addition.

{S.sname |  $\exists S \in \text{Sailor}.(\exists R \in \text{Reserves}.(R.\text{sid}=S.\text{sid} \wedge \nexists B \in \text{Boat}.(B.\text{bid}=R.\text{bid} \wedge B.\text{color}='red'))$ )}

# Nested query 3

Sailor (sid, sname, rating, age)  
Reserves (sid, bid, day)  
Boat (bid, bname, color)

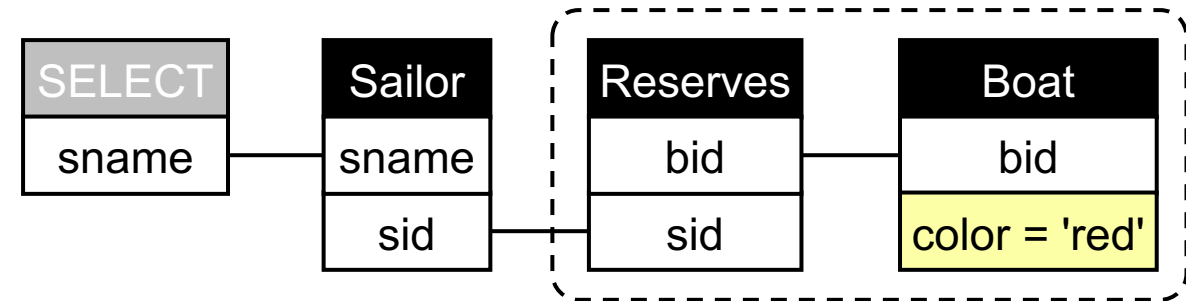


340



Q:

```
SELECT DISTINCT S.sname
FROM Sailor S
WHERE NOT EXISTS
  (SELECT R.sid
   FROM Reserves R
   WHERE R.sid=S.sid
   AND EXISTS
     (SELECT B.bid
      FROM Boat B
      WHERE B.bid = R.bid
      AND B.color = 'red'))
```



{S.sname |  $\exists S \in \text{Sailor}.(\nexists R \in \text{Reserves}.(R.\text{sid}=S.\text{sid} \wedge \exists B \in \text{Boat}.(B.\text{bid}=R.\text{bid} \wedge B.\text{color}='red')))$ }

# Nested query 3

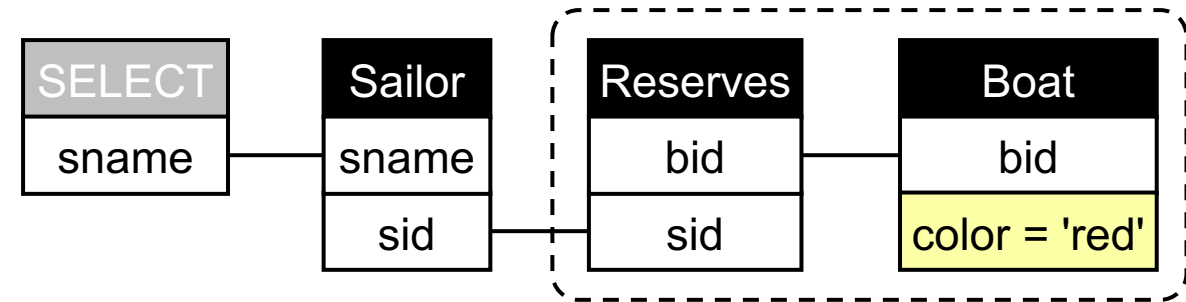
Sailor (sid, sname, rating, age)  
Reserves (sid, bid, day)  
Boat (bid, bname, color)



340

Q: Find sailors who have **not** reserved a red boat.

```
SELECT DISTINCT S.sname
FROM Sailor S
WHERE NOT EXISTS
  (SELECT R.sid
   FROM Reserves R
   WHERE R.sid=S.sid
   AND EXISTS
     (SELECT B.bid
      FROM Boat B
      WHERE B.bid = R.bid
      AND B.color = 'red'))
```



*They can have reserved 0 or more boats in another color, but must not have reserved any red boat.*

$\{S.sname \mid \exists S \in \text{Sailor}. (\nexists R \in \text{Reserves}. (R.sid=S.sid \wedge \exists B \in \text{Boat}. (B.bid=R.bid \wedge B.color='red'))))\}$

# Quiz: Dustin?



340

### Sailor

<u>sid</u>	sname	rating	age
22	Dustin	7	45.0
29	Brutus	1	33.0
31	Lubber	8	55.5
32	Andy	8	25.5
58	Rusty	10	35.0
64	Horatio	7	35.0
71	Zorba	10	16.0
74	Horatio	9	35.0
85	Art	3	25.5
95	Bob	3	63.5

Figure 5.1 An Instance *S3* of Sailors

### Reserves

<u>sid</u>	<u>bid</u>	<u>day</u>
22	101	10/10/98
22	102	10/10/98
22	103	10/8/98
22	104	10/7/98
31	102	11/10/98
31	103	11/6/98
31	104	11/12/98
64	101	9/5/98
64	102	9/8/98
74	103	9/8/98

Figure 5.2 An Instance *R2* of Reserves

### Boat

<u>bid</u>	bname	color
101	Interlake	blue
102	Interlake	red
103	Clipper	green
104	Marine	red

Figure 5.3 An Instance *B1* of Boats

Should Dustin be in the output of either of the two queries?

Q2: Find sailors who have reserved a boat that is not red.

Q3: Find sailors who have not reserved a red boat.



# Quiz: Dustin?



### Sailor

<u>sid</u>	sname	rating	age
22	Dustin	7	45.0
29	Brutus	1	33.0
31	Lubber	8	55.5
32	Andy	8	25.5
58	Rusty	10	35.0
64	Horatio	7	35.0
71	Zorba	10	16.0
74	Horatio	9	35.0
85	Art	3	25.5
95	Bob	3	63.5

Figure 5.1 An Instance *S3* of Sailors

### Reserves

<u>sid</u>	<u>bid</u>	<u>day</u>
22	101	10/10/98
22	102	10/10/98
22	103	10/8/98
22	104	10/7/98
31	102	11/10/98
31	103	11/6/98
31	104	11/12/98
64	101	9/5/98
64	102	9/8/98
74	103	9/8/98

Figure 5.2 An Instance *R2* of Reserves

### Boat

<u>bid</u>	bname	color
101	Interlake	blue
102	Interlake	red
103	Clipper	green
104	Marine	red

Figure 5.3 An Instance *B1* of Boats

Should Dustin be in the output of either of the two queries?

Q2: Find sailors who have reserved a boat that is not red.

Yes!

Q3: Find sailors who have not reserved a red boat.

No!

# Nested query 4

Sailor (sid, sname, rating, age)  
Reserves (sid, bid, day)  
Boat (bid, bname, color)

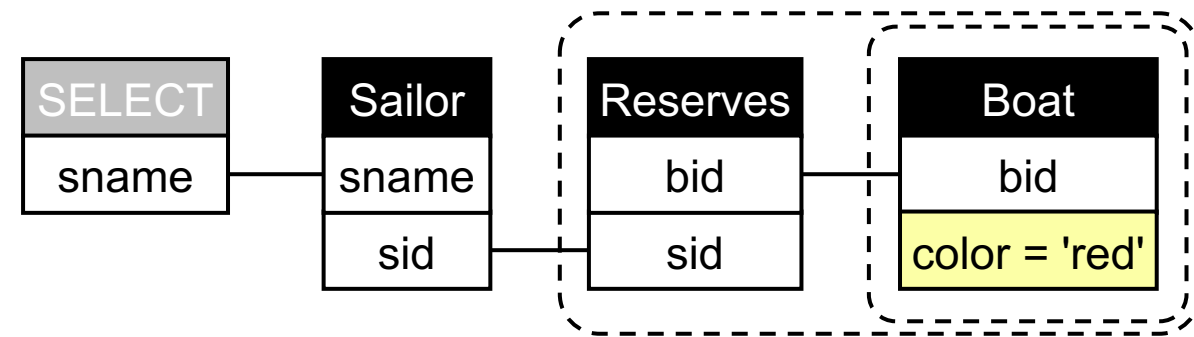


340



Q:

```
SELECT DISTINCT S.sname
FROM Sailor S
WHERE NOT EXISTS
  (SELECT R.sid
   FROM Reserves R
   WHERE R.sid=S.sid
   AND NOT EXISTS
     (SELECT B.bid
      FROM Boat B
      WHERE B.bid = R.bid
      AND B.color = 'red'))
```



$\{S.sname \mid \exists S \in \text{Sailor}.(\nexists R \in \text{Reserves}.(R.sid=S.sid \wedge \nexists B \in \text{Boat}.(B.bid=R.bid \wedge B.color='red')))\}$



# Nested query 4

Sailor (sid, sname, rating, age)  
Reserves (sid, bid, day)  
Boat (bid, bname, color)



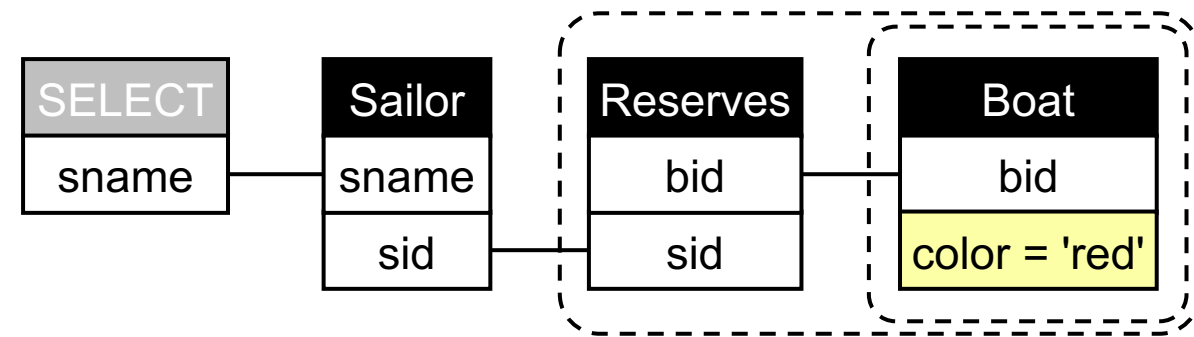
340

They can have reserved 0 or more boats in red, just no other color.

= Find sailors who have reserved **only** red boats

Q: Find sailors who have **not** reserved a boat **that is not red**.

```
SELECT DISTINCT S.sname
FROM Sailor S
WHERE NOT EXISTS
  (SELECT R.sid
   FROM Reserves R
   WHERE R.sid=S.sid
   AND NOT EXISTS
     (SELECT B.bid
      FROM Boat B
      WHERE B.bid = R.bid
      AND B.color = 'red'))
```



They can have reserved 0 or more boats in red, just no other color.

{S.sname |  $\exists S \in \text{Sailor} . (\nexists R \in \text{Reserves} . (R.\text{sid} = S.\text{sid} \wedge \nexists B \in \text{Boat} . (B.\text{bid} = R.\text{bid} \wedge B.\text{color} = \text{'red'})))$ }

# Nested query 4 (universal)

Sailor (sid, sname, rating, age)  
Reserves (sid, bid, day)  
Boat (bid, bname, color)



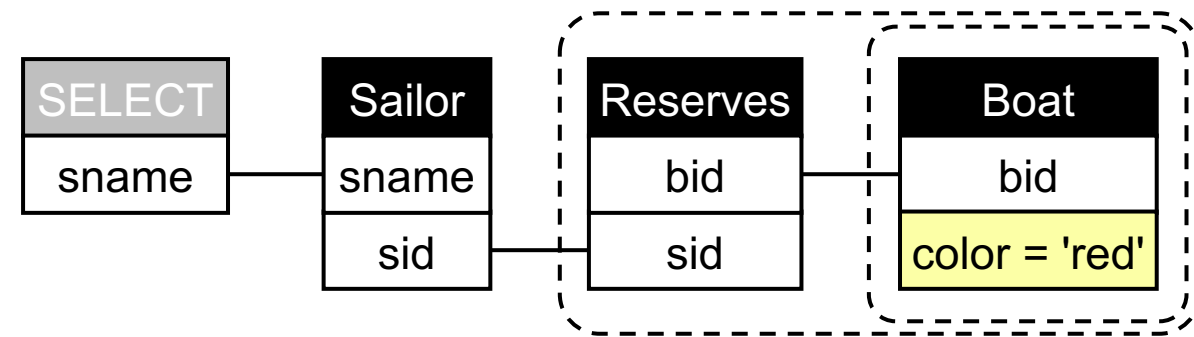
340

*They can have reserved 0 or more boats in red, just no other color.*

= Find sailors who have reserved **only** red boats

Q: Find sailors who have **not** reserved a boat **that is not red**.

```
SELECT DISTINCT S.sname
FROM Sailor S
WHERE NOT EXISTS
  (SELECT R.sid
   FROM Reserves R
   WHERE R.sid=S.sid
   AND NOT EXISTS
     (SELECT B.bid
      FROM Boat B
      WHERE B.bid = R.bid
      AND B.color = 'red'))
```



*They can have reserved 0 or more boats in red, just no other color.*

$\{S.sname \mid \exists S \in \text{Sailor}.(\forall R \in \text{Reserves}.(R.sid=S.sid \rightarrow \exists B \in \text{Boat}.(B.bid=R.bid \wedge B.color='red'))))\}$   
 $\{S.sname \mid \exists S \in \text{Sailor}.(\nexists R \in \text{Reserves}.(R.sid=S.sid \wedge \nexists B \in \text{Boat}.(B.bid=R.bid \wedge B.color='red'))))\}$

# Nested query 4 (another variant)

Sailor (sid, sname, rating, age)  
Reserves (sid, bid, day)  
Boat (bid, bname, color)

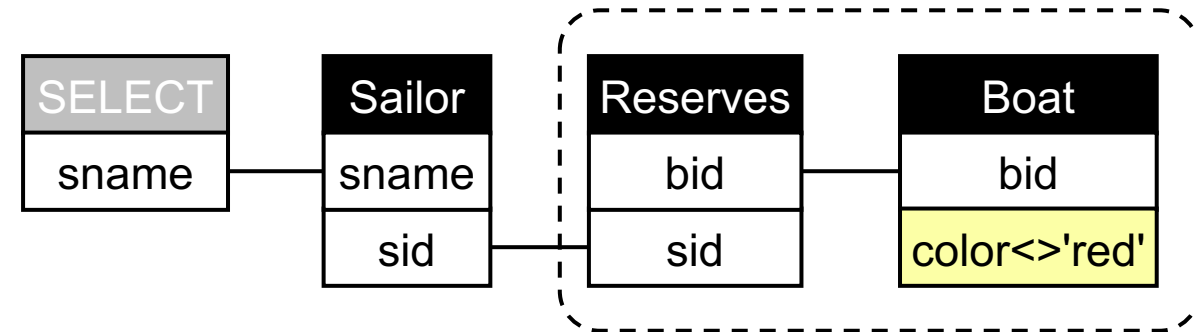


340

= Find sailors who have reserved **only** red boats

Q: Find sailors who have **not** reserved a boat **that is not red**.

```
SELECT DISTINCT S.sname
FROM Sailor S
WHERE NOT EXISTS
  (SELECT R.sid
   FROM Reserves R
   WHERE R.sid=S.sid
   AND EXISTS
     (SELECT B.bid
      FROM Boat B
      WHERE B.bid = R.bid
      AND B.color <> 'red'))
```



*They can have reserved 0 or more boats in red, just no other color.*

*Equivalence with previous variant only because of FK-PK constraint!*

$\{S.sname \mid \exists S \in \text{Sailor}. (\nexists R \in \text{Reserves}. (R.sid=S.sid \wedge \exists B \in \text{Boat}. (B.bid=R.bid \wedge B.color \neq 'red'))))\}$

# Nested query 5

Sailor (sid, sname, rating, age)  
Reserves (sid, bid, day)  
Boat (bid, bname, color)

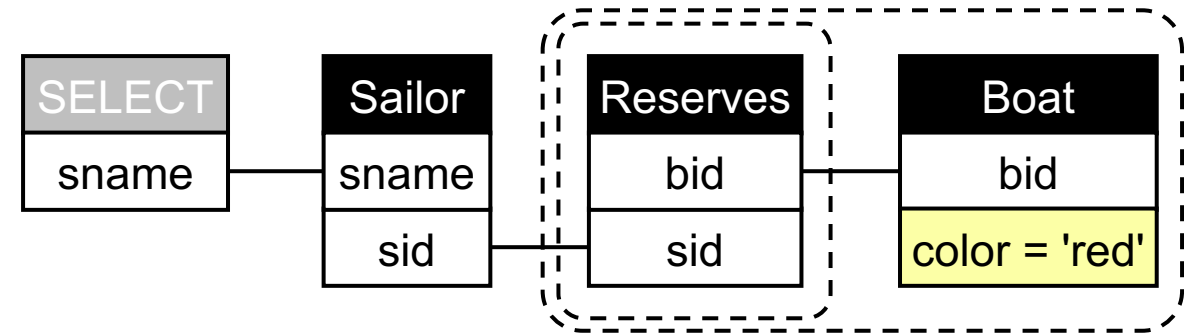


340



Q:

```
SELECT DISTINCT S.sname
FROM Sailor S
WHERE NOT EXISTS
  (SELECT B.bid
   FROM Boat B
   WHERE B.color = 'red'
   AND NOT EXISTS
     (SELECT R.bid
      FROM Reserves R
      WHERE R.bid = B.bid
      AND R.sid = S.sid))
```



{S.sname |  $\exists S \in \text{Sailor} . (\nexists B \in \text{Boat} . (B.\text{color} = \text{'red'} \wedge \nexists R \in \text{Reserves} . (B.\text{bid} = R.\text{bid} \wedge R.\text{sid} = S.\text{sid})))$ }

# Nested query 5

Sailor (sid, sname, rating, age)  
Reserves (sid, bid, day)  
Boat (bid, bname, color)

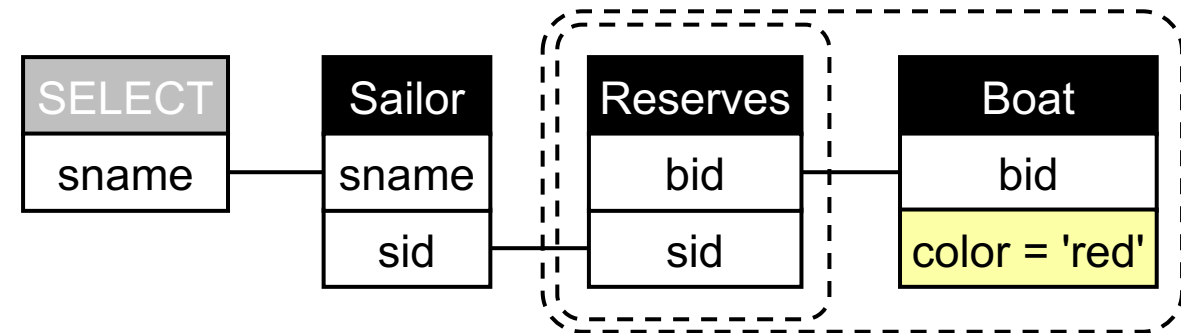


340

= Find sailors who have reserved **all red** boats

Q: Find sailors so there is **no red** boat that is **not** reserved by the sailor.

```
SELECT DISTINCT S.sname
FROM Sailor S
WHERE NOT EXISTS
  (SELECT B.bid
   FROM Boat B
   WHERE B.color = 'red'
   AND NOT EXISTS
     (SELECT R.bid
      FROM Reserves R
      WHERE R.bid = B.bid
      AND R.sid = S.sid))
```



*I don't know of a way to write that query with IN instead of EXISTS and without an explicit cross product between sailors and red boats. (More on that in a moment and also later when we discuss this query in relational algebra.)*

{S.sname |  $\exists S \in \text{Sailor}. (\nexists B \in \text{Boat}. (B.\text{color} = \text{'red'} \wedge \nexists R \in \text{Reserves}. (B.\text{bid} = R.\text{bid} \wedge R.\text{sid} = S.\text{sid})))$ }

# Nested query 5 (universal)

Sailor (sid, sname, rating, age)  
Reserves (sid, bid, day)  
Boat (bid, bname, color)

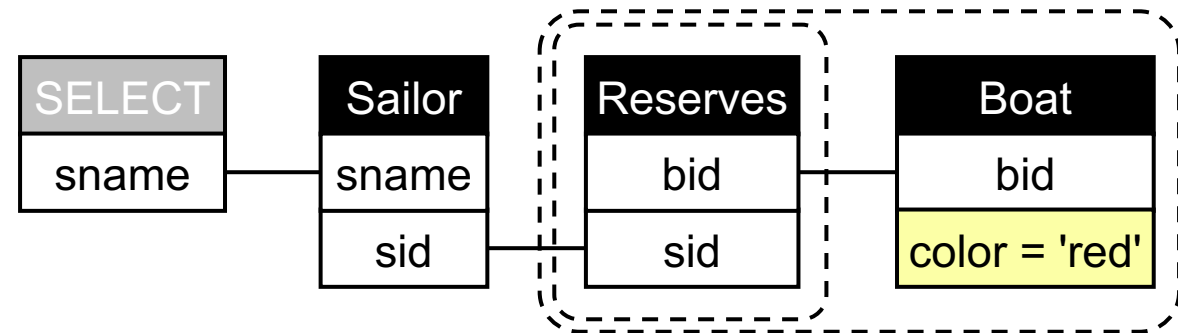


340

= Find sailors who have reserved **all red** boats

Q: Find sailors so there is **no red** boat that is **not** reserved by the sailor.

```
SELECT DISTINCT S.sname
FROM Sailor S
WHERE NOT EXISTS
  (SELECT B.bid
   FROM Boat B
   WHERE B.color = 'red'
   AND NOT EXISTS
     (SELECT R.bid
      FROM Reserves R
      WHERE R.bid = B.bid
      AND R.sid = S.sid))
```



I don't know of a way to write that query with IN instead of EXISTS and without an explicit cross product between sailors and red boats. (More on that in a moment and also later when we discuss this query in relational algebra.)

{S.sname |  $\exists S \in \text{Sailor}. (\forall B \in \text{Boat}. (B.\text{color} = \text{'red'} \rightarrow \exists R \in \text{Reserves}. (B.\text{bid} = R.\text{bid} \wedge R.\text{sid} = S.\text{sid})))$ }

{S.sname |  $\exists S \in \text{Sailor}. (\nexists B \in \text{Boat}. (B.\text{color} = \text{'red'} \wedge \nexists R \in \text{Reserves}. (B.\text{bid} = R.\text{bid} \wedge R.\text{sid} = S.\text{sid})))$ }

# Nested query 5 (w/o correlation)

Sailor (sid, sname, rating, age)  
Reserves (sid, bid, day)  
Boat (bid, bname, color)

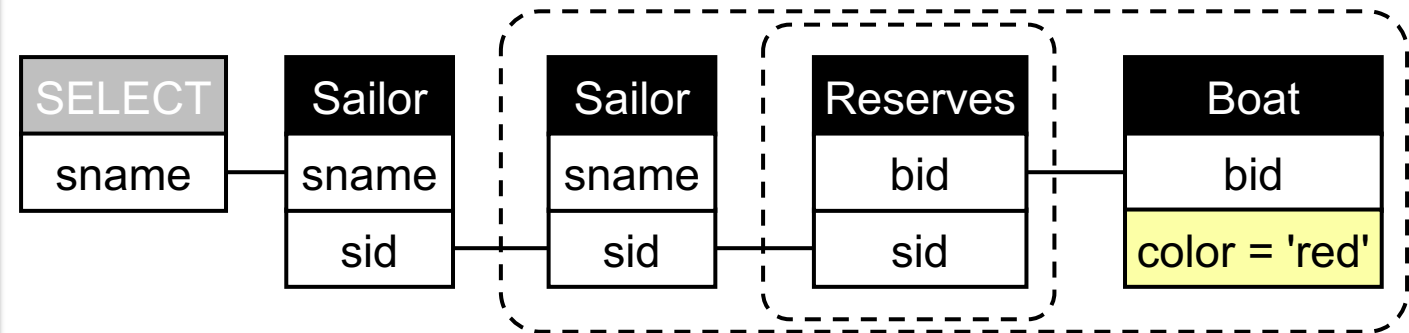


340

= Find sailors who have reserved **all red** boats

Q: Find sailors so there is **no red** boat that is **not** reserved by the sailor.

```
SELECT DISTINCT S.sname
FROM Sailor S
WHERE S.sid NOT IN
  (SELECT S2.sid
   FROM Sailor S2, Boat B
   WHERE B.color = 'red'
   AND (S2.sid, B.bid) NOT IN
     (SELECT R.sid, R.bid
      FROM Reserves R))
```



$\{S.sname \mid \exists S \in \text{Sailor}. (\forall S2 \in \text{Sailor}, \forall B \in \text{Boat}. (B.color = 'red' \wedge S2.sid = S.sid \rightarrow \exists R \in \text{Reserves}. (B.bid = R.bid \wedge R.sid = S2.sid)))\}$   
 $\{S.sname \mid \exists S \in \text{Sailor}. (\nexists S2 \in \text{Sailor}, \nexists B \in \text{Boat}. (B.color = 'red' \wedge S2.sid = S.sid \wedge \nexists R \in \text{Reserves}. (B.bid = R.bid \wedge R.sid = S.sid)))\}$

# Nested query 5 (w/o correlation)

Sailor (sid, sname, rating, age)  
Reserves (sid, bid, day)  
Boat (bid, bname, color)

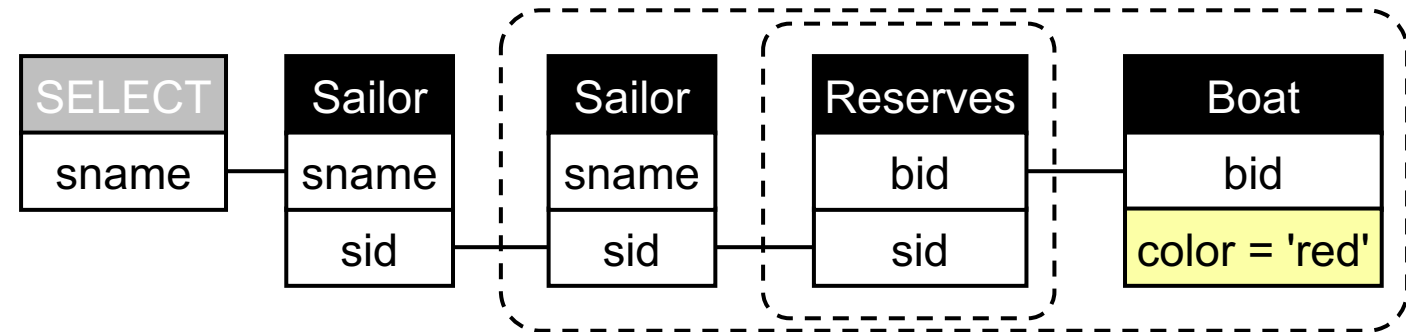


340

= Find sailors who have reserved **all red** boats

Q: Find sailors so there is **no red** boat that is **not** reserved by the sailor.

```
SELECT DISTINCT S.sname
FROM Sailor S
WHERE NOT EXISTS
  (SELECT *
   FROM Sailor S2, Boat B
   WHERE B.color = 'red'
   AND S.sid = S2.sid
   AND NOT EXISTS
     (SELECT *
      FROM Reserves R
      WHERE B.bid = R.bid
      AND S2.sid = R.sid))
```



$\{S.sname \mid \exists S \in \text{Sailor}. (\forall S2 \in \text{Sailor}, \forall B \in \text{Boat}. (B.color = 'red' \wedge S2.sid = S.sid \rightarrow \exists R \in \text{Reserves}. (B.bid = R.bid \wedge R.sid = S2.sid)))\}$

$\{S.sname \mid \exists S \in \text{Sailor}. (\nexists S2 \in \text{Sailor}, \nexists B \in \text{Boat}. (B.color = 'red' \wedge S2.sid = S.sid \wedge \nexists R \in \text{Reserves}. (B.bid = R.bid \wedge R.sid = S.sid)))\}$



# Towards SQL patterns

Sailor (sid, sname, rating, age)  
Reserves (sid, bid, day)  
Boat (bid, bname, color)

	Sailors who <b>have not</b> reserved a red boat	Sailors who reserved <b>only</b> red boats	Sailors who reserved <b>all</b> red boats
SQL	<pre>SELECT DISTINCT S.sname FROM Sailor S WHERE NOT EXISTS(   SELECT *   FROM Reserves R, Boat B   WHERE R.sid = S.sid   AND R.bid = B.bid   AND B.color = 'red')</pre>	<pre>SELECT DISTINCT S.sname FROM Sailor S WHERE NOT EXISTS(   SELECT *   FROM Reserves R   WHERE R.sid = S.sid   AND NOT EXISTS(     SELECT *     FROM Boat B     WHERE R.bid = B.bid     AND B.color = 'red'))</pre>	<pre>SELECT DISTINCT S.sname FROM Sailor S WHERE NOT EXISTS(   SELECT *   FROM Boat B   WHERE B.color = 'red'   AND NOT EXISTS(     SELECT *     FROM Reserves R     WHERE R.bid = B.bid     AND R.sid = S.sid))</pre>

# Towards SQL patterns

Sailor (sid, sname, rating, age)  
 Reserves (sid, bid, day)  
 Boat (bid, bname, color)

	Sailors who <b>have not</b> reserved a red boat	Sailors who reserved <b>only</b> red boats	Sailors who reserved <b>all</b> red boats
SQL	<pre>SELECT DISTINCT S.sname FROM Sailor S WHERE NOT EXISTS(   SELECT *   FROM Reserves R, Boat B   WHERE R.sid = S.sid   AND R.bid = B.bid   AND B.color = 'red')</pre>	<pre>SELECT DISTINCT S.sname FROM Sailor S WHERE NOT EXISTS(   SELECT *   FROM Reserves R   WHERE R.sid = S.sid   AND NOT EXISTS(     SELECT *     FROM Boat B     WHERE R.bid = B.bid     AND B.color = 'red'))</pre>	<pre>SELECT DISTINCT S.sname FROM Sailor S WHERE NOT EXISTS(   SELECT *   FROM Boat B   WHERE B.color = 'red'   AND NOT EXISTS(     SELECT *     FROM Reserves R     WHERE R.bid = B.bid     AND R.sid = S.sid))</pre>
RD			

Sailor (sid, sname, rating, age)  
Reserves (sid, bid, day)  
Boat (bid, bname, color)

	<b>not</b>	<b>only</b>	<b>all</b>
Sailors renting boats	have not reserved a red boat	reserved only red boats	reserved all red boats

Sailor (sid, sname, rating, age)  
 Reserves (sid, bid, day)  
 Boat (bid, bname, color)

Student (sid, sname)  
 Takes (sid, cid, semester)  
 Course (cid, cname, department)

	<b>not</b>	<b>only</b>	<b>all</b>
Sailors renting boats	have not reserved a red boat	reserved only red boats	reserved all red boats
Students taking classes	took no art class	took only art classes	took all art classes

Sailor (sid, sname, rating, age)  
 Reserves (sid, bid, day)  
 Boat (bid, bname, color)

Student (sid, sname)  
 Takes (sid, cid, semester)  
 Course (cid, cname, department)

Actor (aid, aname)  
 Plays (aid, mid, role)  
 Movie (mid, mname, director)

	not	only	all
Sailors renting boats	have not reserved a red boat	reserved only red boats	reserved all red boats
Students taking classes	took no art class	took only art classes	took all art classes
Actors playing in movies	did not play in a Hitchcock movie	played only Hitchcock movies	played in all Hitchcock movies

Sailor (sid, sname, rating, age)  
 Reserves (sid, bid, day)  
 Boat (bid, bname, color)

Student (sid, sname)  
 Takes (sid, cid, semester)  
 Course (cid, cname, department)

Actor (aid, aname)  
 Plays (aid, mid, role)  
 Movie (mid, mname, director)

	not	only	all
Sailors	<pre>SELECT DISTINCT S.sname FROM Sailor S WHERE NOT EXISTS(   SELECT *   FROM Reserves R, Boat B   WHERE R.sid = S.sid   AND R.bid = B.bid   AND B.color = 'red')</pre>	<pre>SELECT DISTINCT S.sname FROM Sailor S WHERE NOT EXISTS(   SELECT *   FROM Reserves R   WHERE R.sid = S.sid   AND NOT EXISTS(     SELECT *     FROM Boat B     WHERE R.bid = B.bid     AND B.color = 'red'))</pre>	<pre>SELECT DISTINCT S.sname FROM Sailor S WHERE NOT EXISTS(   SELECT *   FROM Boat B   WHERE B.color = 'red'   AND NOT EXISTS(     SELECT *     FROM Reserves R     WHERE R.bid = B.bid     AND R.sid = S.sid))</pre>
Students	<pre>SELECT DISTINCT S.sname FROM Student S WHERE NOT EXISTS(   SELECT *   FROM Takes T, Class C   WHERE T.sid = S.sid   AND T.cid = C.cid   AND C.department = 'art')</pre>	<pre>SELECT DISTINCT S.sname FROM Student S WHERE NOT EXISTS(   SELECT *   FROM Takes T   WHERE T.sid = S.sid   AND NOT EXISTS(     SELECT *     FROM Class C     WHERE T.cid = C.cid     AND C.department= 'art'))</pre>	<pre>SELECT DISTINCT S.sname FROM Student S WHERE NOT EXISTS(   SELECT *   FROM Class C   WHERE C.department = 'art'   AND NOT EXISTS(     SELECT *     FROM Takes T     WHERE T.cid = C.cid     AND T.sid = S.sid))</pre>
Actors	<pre>SELECT DISTINCT A.aname FROM Actor A WHERE NOT EXISTS(   SELECT *   FROM Plays P, Movie M   WHERE P.aid = A.aid   AND P.mid = M.mid   AND M.director= 'Hitchcock')</pre>	<pre>SELECT DISTINCT A.aname FROM Actor A WHERE NOT EXISTS(   SELECT *   FROM Plays P   WHERE P.aid = A.aid   AND NOT EXISTS(     SELECT *     FROM Movie M     WHERE P.mid = M.mid     AND M.director= 'Hitchcock'))</pre>	<pre>SELECT DISTINCT A.aname FROM Actor A WHERE NOT EXISTS(   SELECT *   FROM Movie M   WHERE M.director= 'Hitchcock'   AND NOT EXISTS(     SELECT *     FROM Plays P     WHERE P.mid = M.mid     AND P.aid = A.aid))</pre>

Sailor (sid, sname, rating, age)  
 Reserves (sid, bid, day)  
 Boat (bid, bname, color)

Student (sid, sname)  
 Takes (sid, cid, semester)  
 Course (cid, cname, department)

Actor (aid, aname)  
 Plays (aid, mid, role)  
 Movie (mid, mname, director)

	not	only	all
Sailors			
Students			
Actors			

# Logical SQL Patterns

**Logical patterns** are the building blocks of most SQL queries.

Patterns are very hard to extract from the SQL text.

A pattern can appear across different database schemas.

Think of queries like:

- Find sailors who reserved all red boats
- Find students who took all art classes
- Find actors who played in all movies by Hitchcock



# What does this query return?

Likes(drinker,beer)

```
SELECT L1.drinker
FROM Likes L1
WHERE not exists
  (SELECT *
   FROM Likes L2
   WHERE L1.drinker <> L2.drinker
   AND not exists
     (SELECT *
      FROM Likes L3
      WHERE L3.drinker = L2.drinker
      AND not exists
        (SELECT *
         FROM Likes L4
         WHERE L4.drinker = L1.drinker
         AND L4.beer = L3.beer)))
AND not exists
  (SELECT *
   FROM Likes L5
   WHERE L5.drinker = L1.drinker
   AND not exists
     (SELECT *
      FROM Likes L6
      WHERE L6.drinker = L2.drinker
      AND L6.beer = L5.beer)))
```

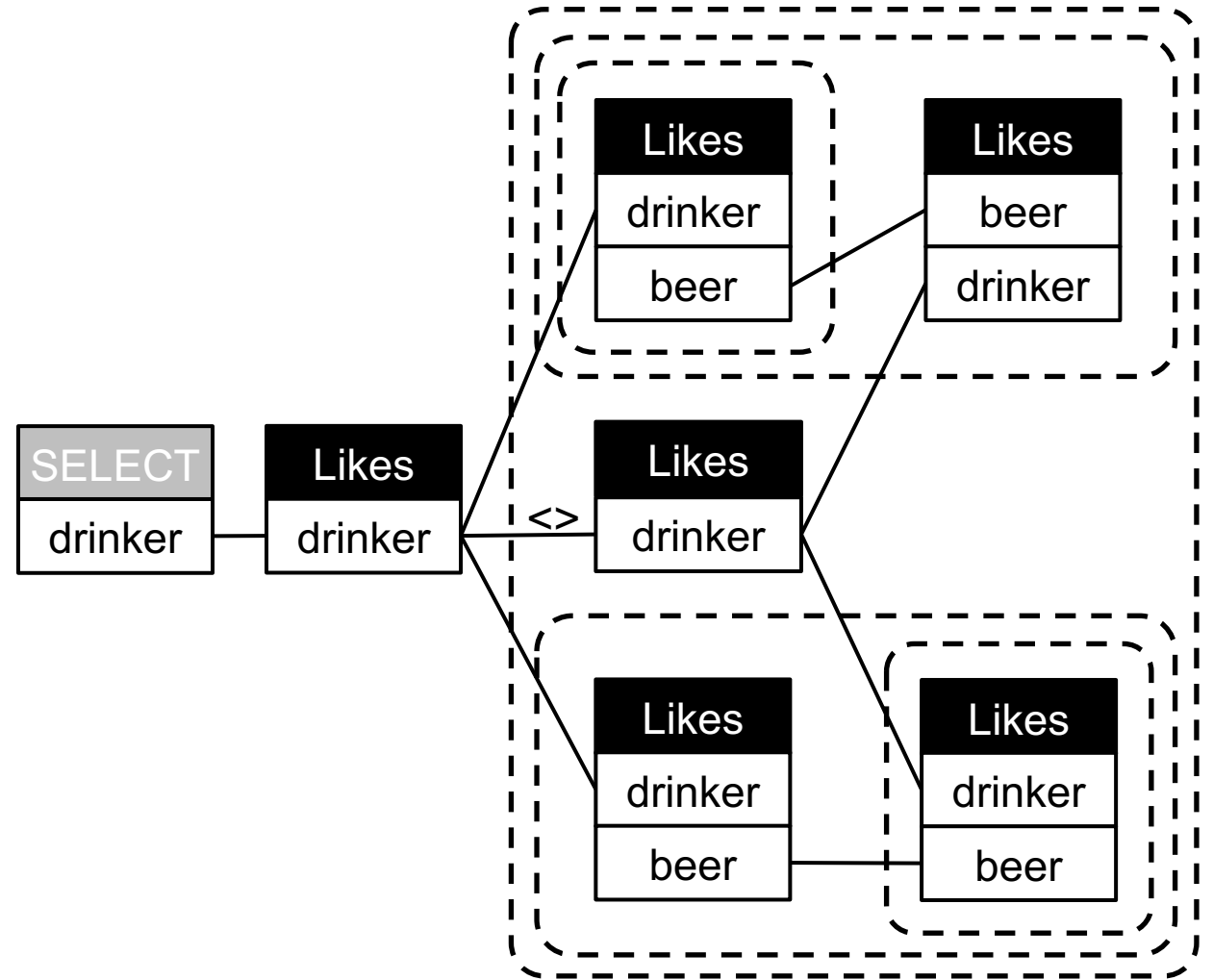
# What does this query return?

```
SELECT L1.drinker
FROM Likes L1
WHERE not exists
  (SELECT *
   FROM Likes L2
   WHERE L1.drinker <> L2.drinker
   AND not exists
     (SELECT *
      FROM Likes L3
      WHERE L3.drinker = L2.drinker
      AND not exists
        (SELECT *
         FROM Likes L4
         WHERE L4.drinker = L1.drinker
         AND L4.beer = L3.beer)))
AND not exists
  (SELECT *
   FROM Likes L5
   WHERE L5.drinker = L1.drinker
   AND not exists
     (SELECT *
      FROM Likes L6
      WHERE L6.drinker = L2.drinker
      AND L6.beer = L5.beer)))
```



## Relational Diagrams scoping

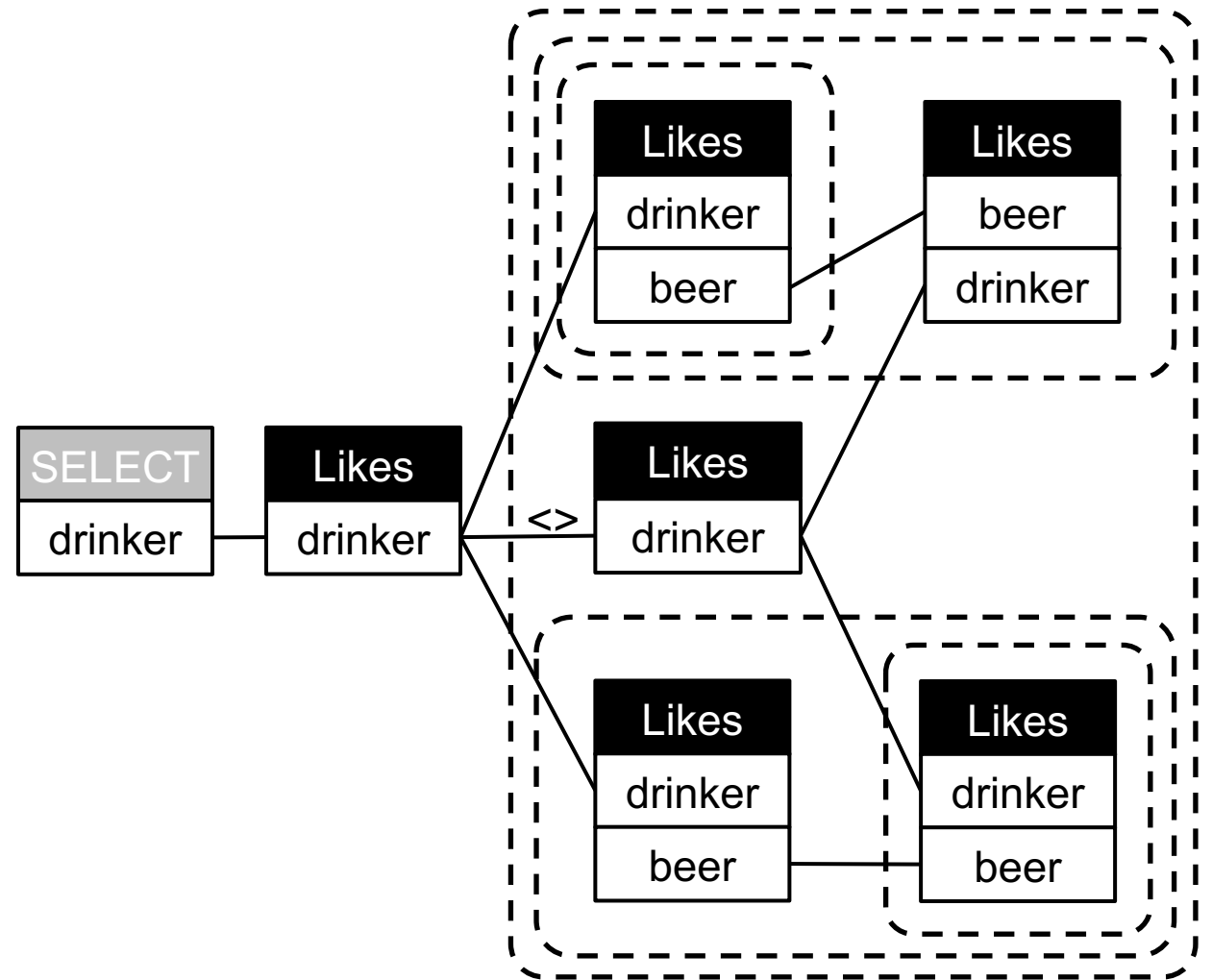
Likes(drinker,beer)



# Q: Finder drinkers with a unique beer taste

Likes(drinker,beer)

```
SELECT L1.drinker
FROM Likes L1
WHERE not exists
  (SELECT *
   FROM Likes L2
   WHERE L1.drinker <> L2.drinker
   AND not exists
     (SELECT *
      FROM Likes L3
      WHERE L3.drinker = L2.drinker
      AND not exists
        (SELECT *
         FROM Likes L4
         WHERE L4.drinker = L1.drinker
         AND L4.beer = L3.beer)))
AND not exists
  (SELECT *
   FROM Likes L5
   WHERE L5. drinker = L1. drinker
   AND not exists
     (SELECT *
      FROM Likes L6
      WHERE L6.drinker = L2.drinker
      AND L6.beer= L5.beer)))
```

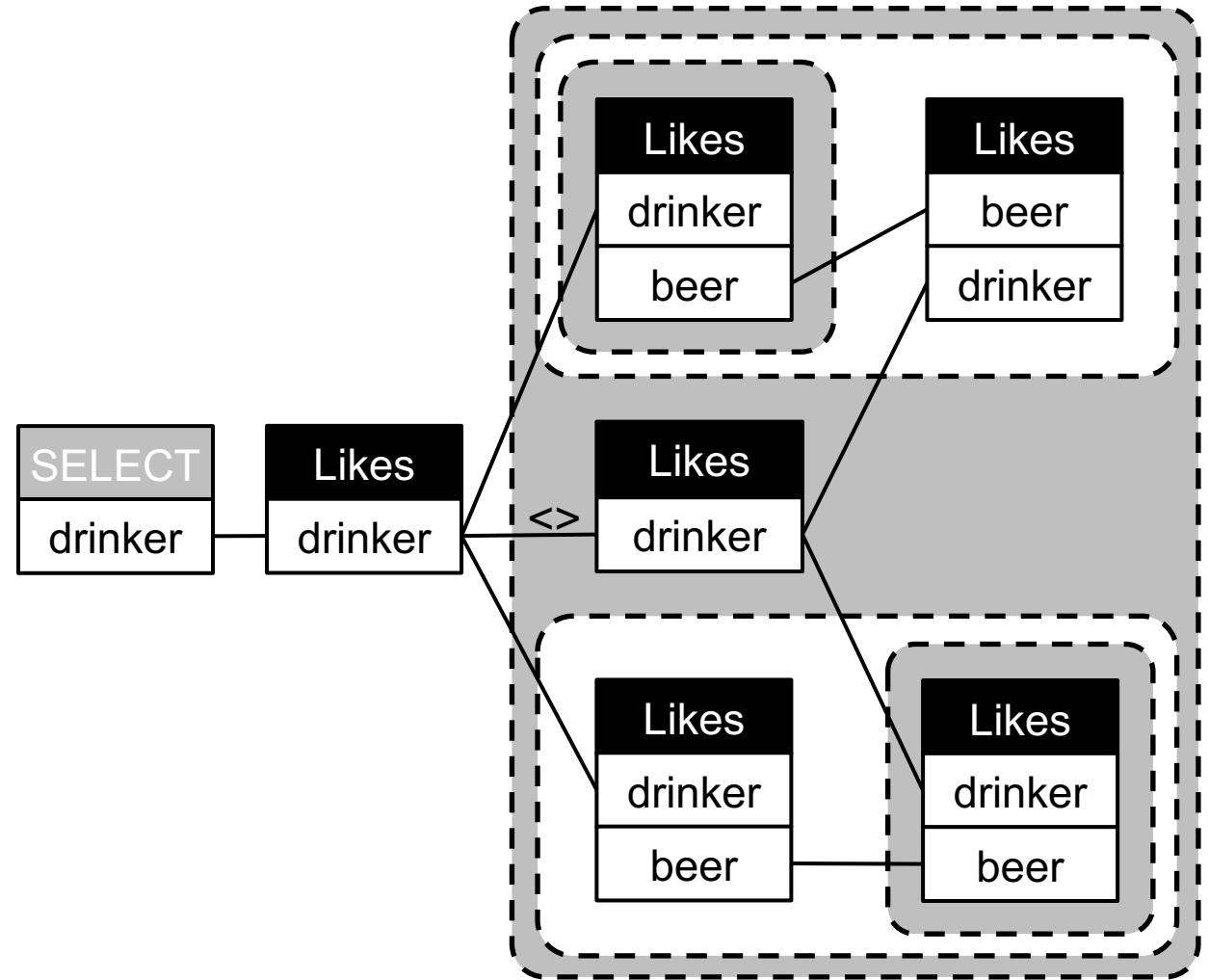


Relational Diagrams scoping

# Q: Finder drinkers with a unique beer taste

Likes(drinker,beer)

```
SELECT L1.drinker
FROM Likes L1
WHERE not exists
  (SELECT *
   FROM Likes L2
   WHERE L1.drinker <> L2.drinker
   AND not exists
     (SELECT *
      FROM Likes L3
      WHERE L3.drinker = L2.drinker
      AND not exists
        (SELECT *
         FROM Likes L4
         WHERE L4.drinker = L1.drinker
         AND L4.beer = L3.beer)))
AND not exists
  (SELECT *
   FROM Likes L5
   WHERE L5.drinker = L1.drinker
   AND not exists
     (SELECT *
      FROM Likes L6
      WHERE L6.drinker = L2.drinker
      AND L6.beer = L5.beer)))
```

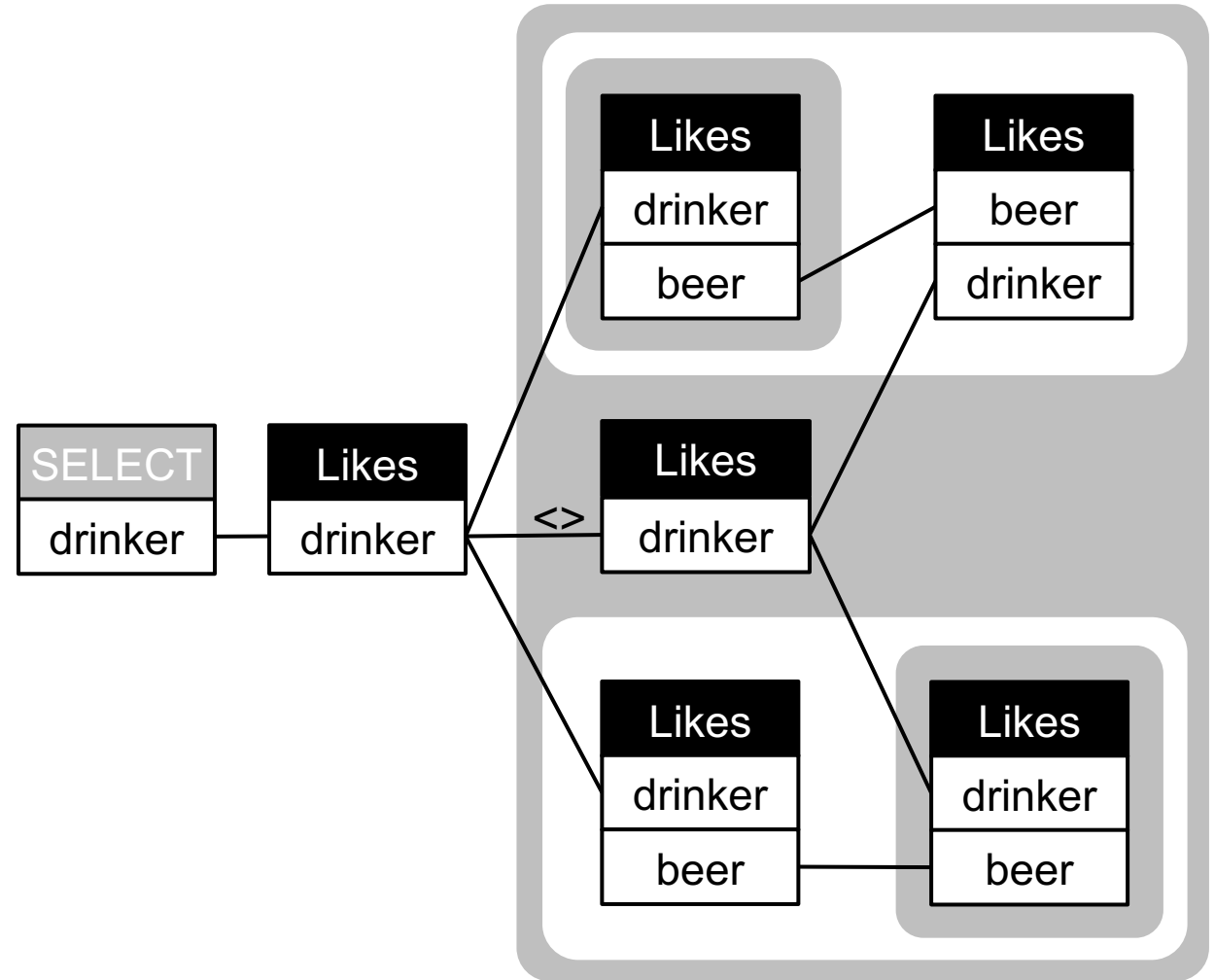


Relational Diagrams scoping

# Q: Finder drinkers with a unique beer taste

Likes(drinker,beer)

```
SELECT L1.drinker
FROM Likes L1
WHERE not exists
  (SELECT *
   FROM Likes L2
   WHERE L1.drinker <> L2.drinker
   AND not exists
     (SELECT *
      FROM Likes L3
      WHERE L3.drinker = L2.drinker
      AND not exists
        (SELECT *
         FROM Likes L4
         WHERE L4.drinker = L1.drinker
         AND L4.beer = L3.beer)))
AND not exists
  (SELECT *
   FROM Likes L5
   WHERE L5.drinker = L1.drinker
   AND not exists
     (SELECT *
      FROM Likes L6
      WHERE L6.drinker = L2.drinker
      AND L6.beer = L5.beer)))
```

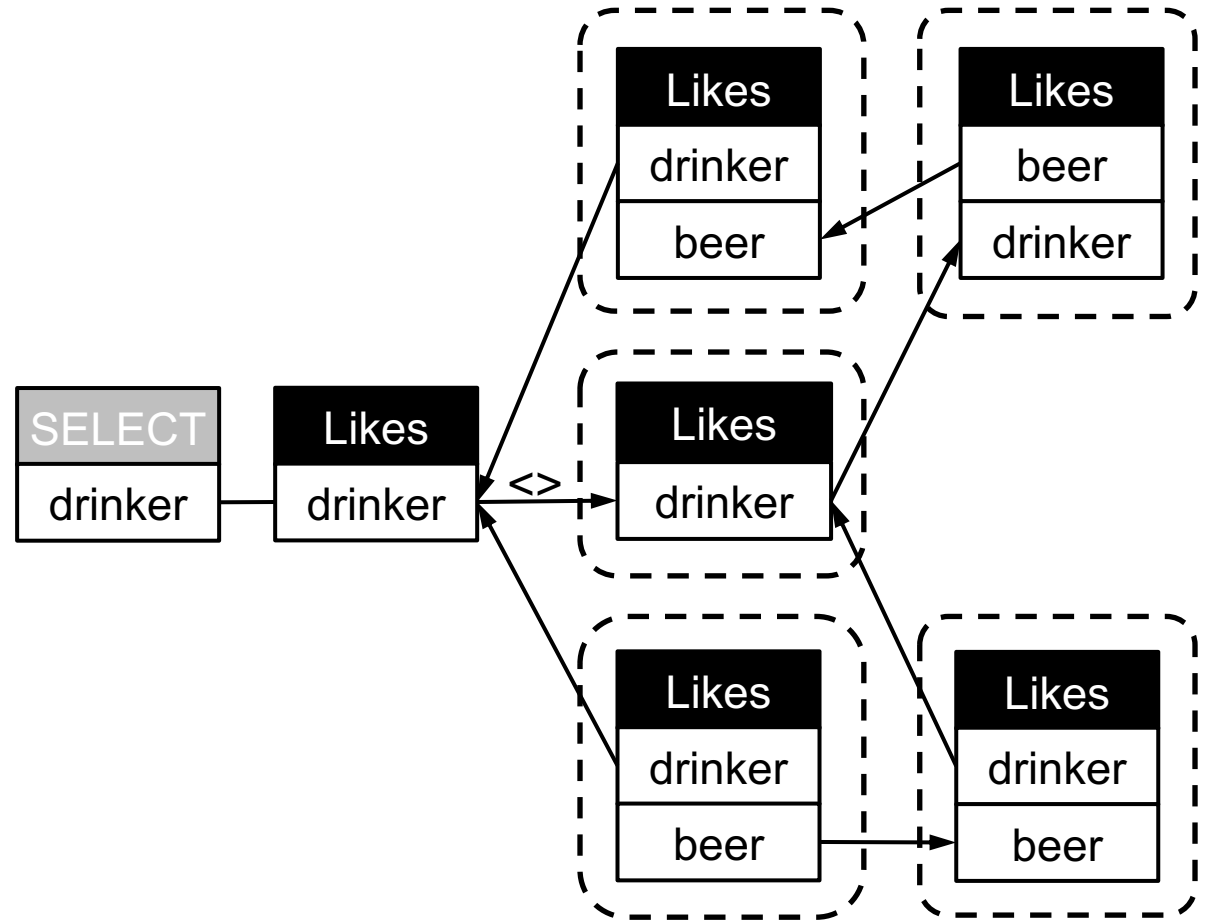


Relational Diagrams scoping

# Q: Finder drinkers with a unique beer taste

Likes(drinker,beer)

```
SELECT L1.drinker
FROM Likes L1
WHERE not exists
  (SELECT *
   FROM Likes L2
   WHERE L1.drinker <> L2.drinker
   AND not exists
     (SELECT *
      FROM Likes L3
      WHERE L3.drinker = L2.drinker
      AND not exists
        (SELECT *
         FROM Likes L4
         WHERE L4.drinker = L1.drinker
         AND L4.beer = L3.beer))
     AND not exists
       (SELECT *
        FROM Likes L5
        WHERE L5.drinker = L1.drinker
        AND not exists
          (SELECT *
           FROM Likes L6
           WHERE L6.drinker = L2.drinker
           AND L6.beer = L5.beer)))
```



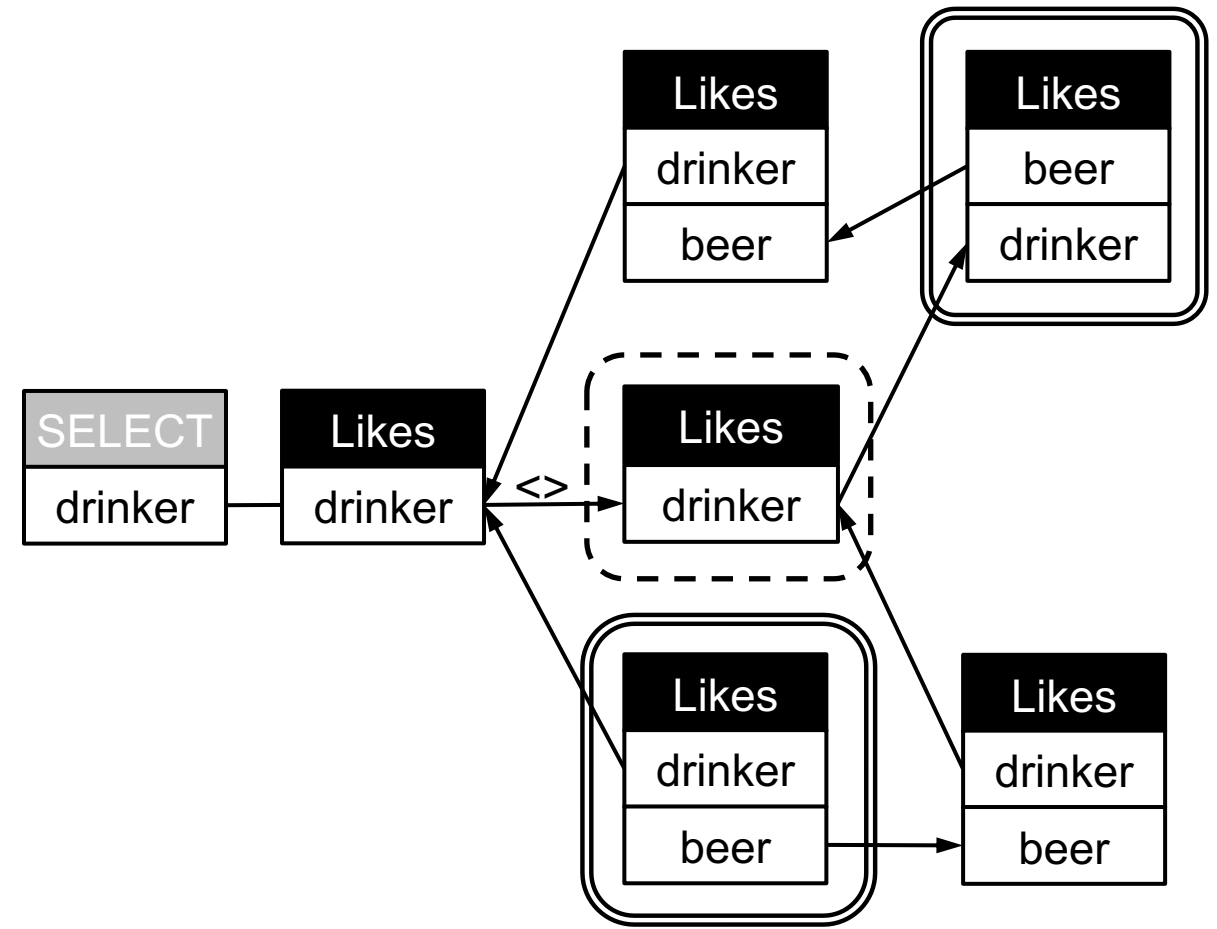
QueryVis scoping

Relational Diagrams scoping

# Q: Finder drinkers with a unique beer taste

Likes(drinker,beer)

```
SELECT L1.drinker
FROM Likes L1
WHERE not exists
  (SELECT *
   FROM Likes L2
   WHERE L1.drinker <> L2.drinker
   AND not exists
     (SELECT *
      FROM Likes L3
      WHERE L3.drinker = L2.drinker
      AND not exists
        (SELECT *
         FROM Likes L4
         WHERE L4.drinker = L1.drinker
         AND L4.beer = L3.beer)))
AND not exists
  (SELECT *
   FROM Likes L5
   WHERE L5.drinker = L1.drinker
   AND not exists
     (SELECT *
      FROM Likes L6
      WHERE L6.drinker = L2.drinker
      AND L6.beer = L5.beer)))
```



QueryVis scoping

Relational Diagrams scoping

<https://demo.queryvis.com>

# QueryViz

Input: Schema

Input Query

Output: Visualization

**Your Input**

Specify or choose a pre-defined schema help

Employee and Department

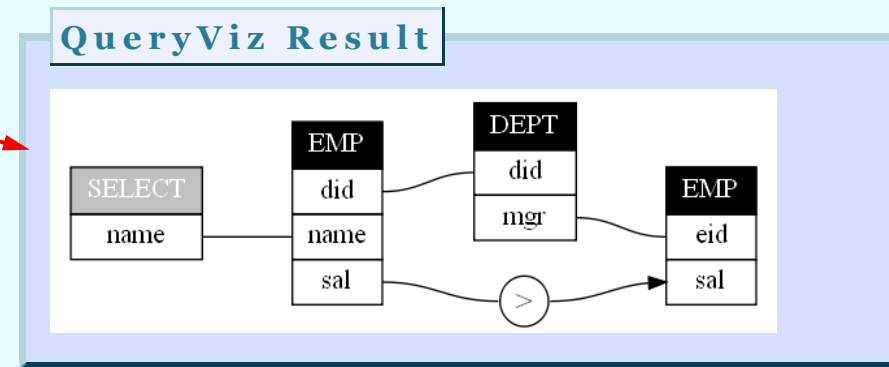
```
EMP(eid,name,sal,did)
DEPT(did,dname,mgr)
```

Specify or choose an SQL Query help

Query 8

```
SELECT e1.name
FROM EMP e1, EMP e2, DEPT d
WHERE e1.did = d.did
AND d.mgr = e2.eid
AND e1.sal > e2.sal
```

Submit



<https://queryvis.com/>

<http://www.youtube.com/watch?v=kVFnQRGAQIs>

Source: Danaparamita, Gatterbauer: QueryViz: Helping users understand SQL queries and their patterns. EDBT 2011. <https://doi.org/10.14778/3402755.3402805>

See also: Gatterbauer, Dunne, Jagadish, Riedewald: Principles of Query Visualization. IEEE Debull 2023. <http://sites.computer.org/debull/A22sept/p47.pdf>

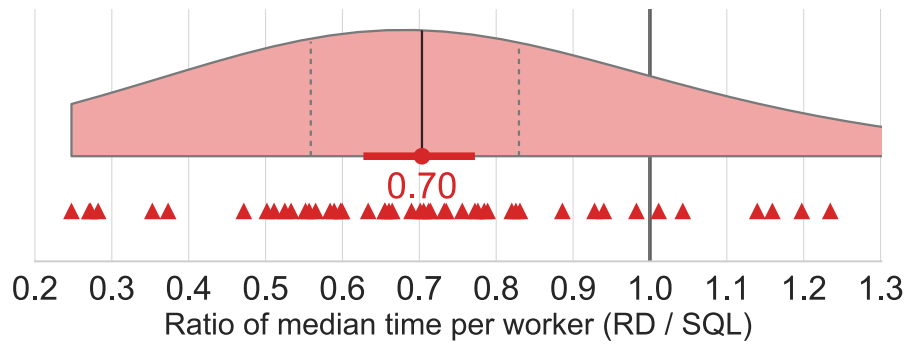
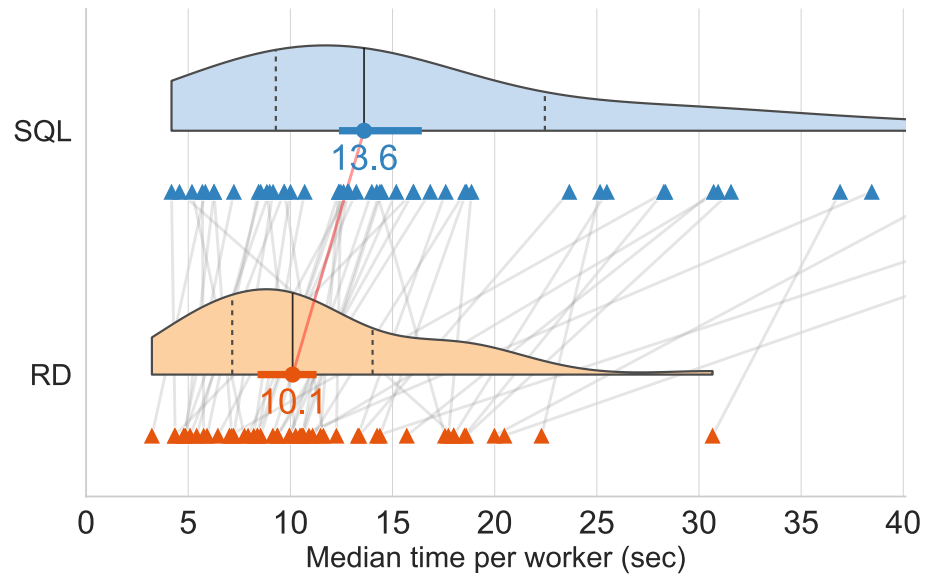
Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>



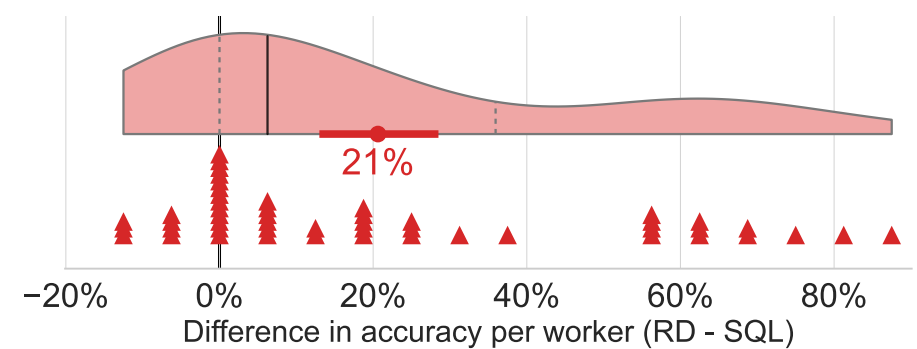
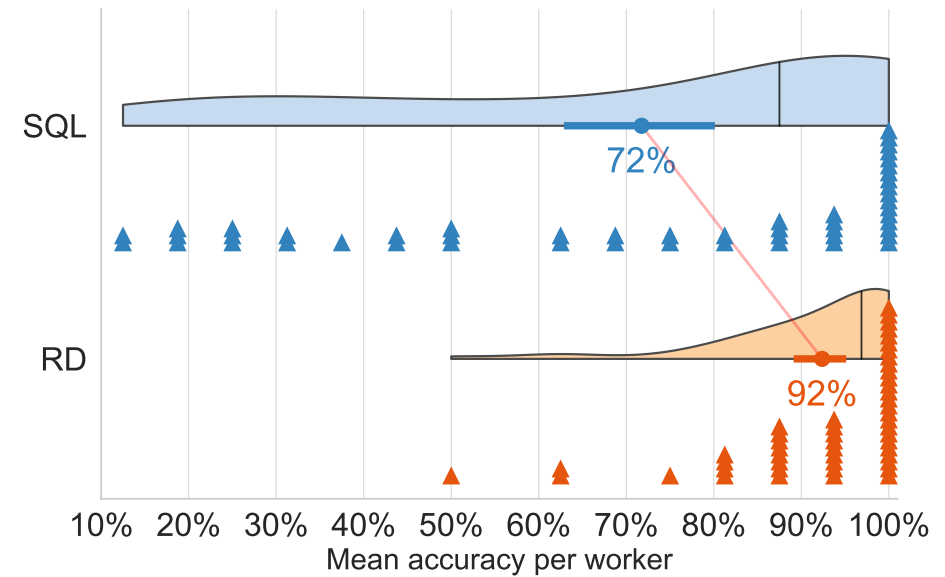
# Preregistered, randomized user study on AMT

n = 50 participants, preregistration: <https://osf.io/4zpsk>

## Speed



## Accuracy



# Preregistered, randomized user study on AMT

n = 50 participants, preregistration: <https://osf.io/4zpsk>

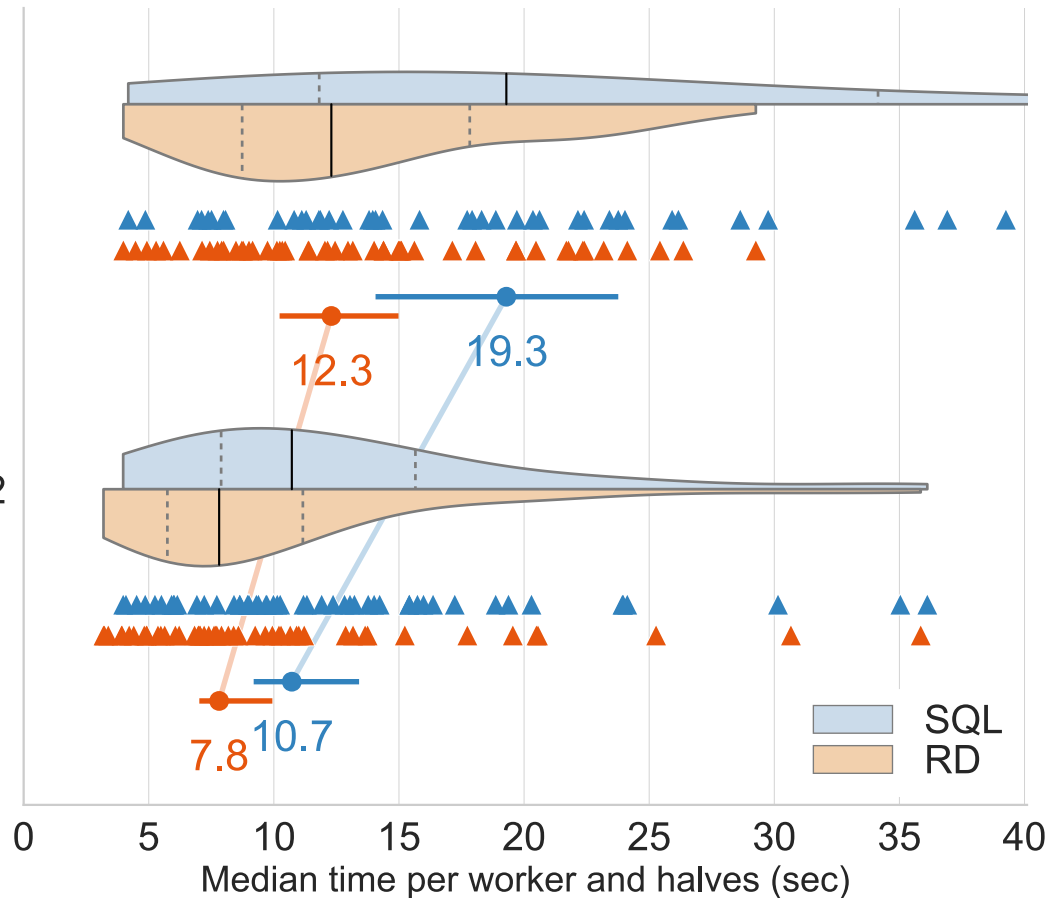
## Learning

#1 = first 16 questions

#2 = second 16 questions

H1

H2



# Outline: T1-U1: SQL

- SQL

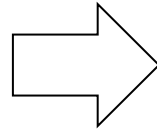
- Schema, keys, referential integrity
- Joins
- Aggregates and grouping
- Nested queries (Subqueries)
- **Union and Theta Joins**
- Nulls & Outer joins
- Top-k
- [Recursion: moved to T1-U4: Datalog]

# Union

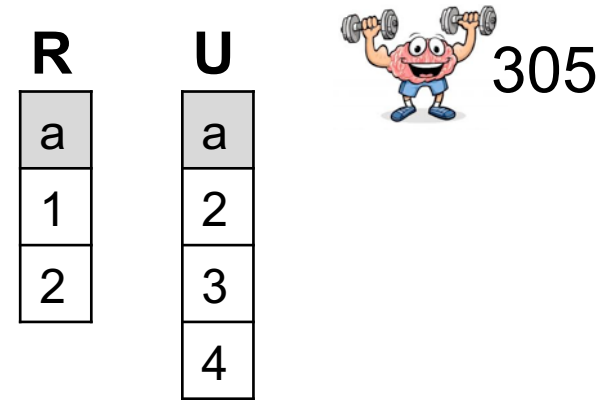
R	U
a	a
1	2
2	3
	4



```
SELECT a
FROM R
UNION
SELECT a
FROM U
```

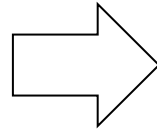


# Union uses set semantics

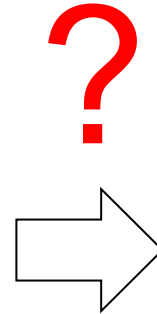


What if we wanted also the duplicates

```
SELECT a
FROM R
UNION
SELECT a
FROM U
```

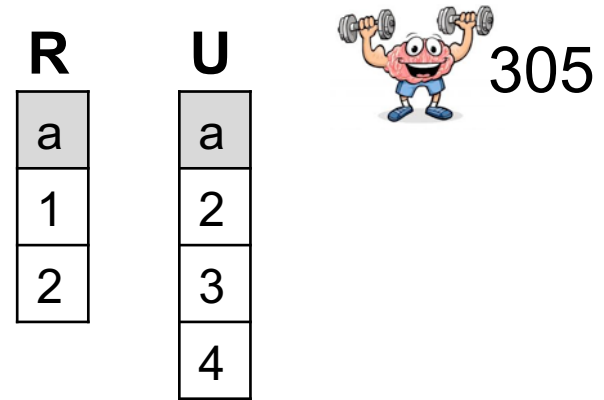


a
1
2
3
4

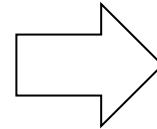


a
1
2
2
3
4

# Union ALL uses bag semantics



```
SELECT a
FROM R
UNION ALL
SELECT a
FROM U
```



a
1
2
2
3
4

# Theta joins

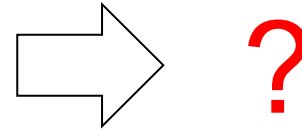
*What do these queries compute?*

R
a
1
2

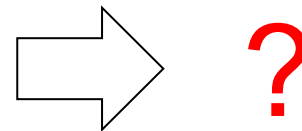
U
a
2
3
4



```
SELECT R.a, U.a as b
FROM   R, U
WHERE  R.a < U.a
```



```
SELECT R.a, U.a as b
FROM   R, U
WHERE  R.a >= U.a
```

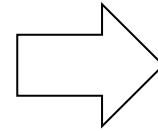


A **Theta-join** allows for arbitrary comparison relationships (such as  $\geq$ ).  
An **equijoin** is a theta join using the equality operator.

# Theta joins

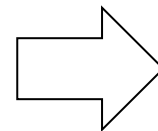
*What do these queries compute?*

```
SELECT R.a, U.a as b
FROM R, U
WHERE R.a < U.a
```



a	b
1	2
1	3
1	4
2	3
2	4

```
SELECT R.a, U.a as b
FROM R, U
WHERE R.a >= U.a
```



?

**R**

a
1
2

**U**

a
2
3
4



A **Theta-join** allows for arbitrary comparison relationships (such as  $\geq$ ).  
An **equijoin** is a theta join using the equality operator.

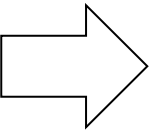


# Theta joins

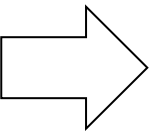
*What do these queries compute?*

```
SELECT R.a, U.a as b
FROM R, U
WHERE R.a < U.a
```


```
SELECT R.a, U.a as b
FROM R, U
WHERE R.a >= U.a
```



a	b
1	2
1	3
1	4
2	3
2	4



a	b
2	2



R
a
1
2

U
a
2
3
4



*Think about these two queries as a partition of the Cartesian product*

A **Theta-join** allows for arbitrary comparison relationships (such as  $\geq$ ). An **equijoin** is a theta join using the equality operator.

# Outline: T1-U1: SQL

- SQL

- Schema, keys, referential integrity
- Joins
- Aggregates and grouping
- Nested queries (Subqueries)
- Union and Theta Joins
- Nulls & Outer joins
- Top-k
- [Recursion: moved to T1-U4: Datalog]

# 3-valued logic example



- Three logicians walk into a bar. The bartender asks: "Do all of you want a drink?"
- The 1st logician says: "I don't know."
- The 2nd logician says: "I don't know."
- The 3rd logician says: "Yes!"

What is going on here ?

# Nulls in SQL

- Whenever we don't have a value, we can put a NULL
- Can mean many things, e.g.:



# Nulls in SQL

- Whenever we don't have a value, we can put a NULL

- Can mean many things, e.g.:

- Value exists but is unknown
- Value not applicable (yet)

*A new class without a grade*

Student	Class	Semester	grade
Alice	cs3200	Fall 2023	B+
Bob	cs3200	Spring 2024	null

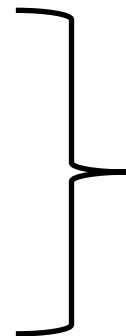
- The schema specifies for each attribute if it can be NULL (nullable attribute) or not ("**NOT NULL**")
- Lots of ongoing research on NULLs
- Next: How does SQL cope with tables that have NULLs ?

# Null Values

- In SQL there are three Boolean values ("**ternary logic**")
  - FALSE, TRUE, UNKNOWN
- If  $x = \text{NULL}$  then
  - **Boolean conditions** are also NULL. E.g:  $x = \text{'Joe'}$
  - **Arithmetic operations** produce NULL. E.g:  $4 * (3 - x) / 7$
  - But **aggregates** ignore NULL values (exception:  $\text{count}(*))$

- Logical reasoning:

- FALSE = 0
- TRUE = 1
- **UNKNOWN = 0.5**



$x \text{ AND } y = \min(x, y)$   
 $x \text{ OR } y = \max(x, y)$   
NOT  $x = (1 - x)$

Handwritten examples in red:  
1 0.5 → 0.5  
0 0.5 → 0

# Join Illustration



## English

eText	<u>eid</u>
One	1
Two	2
Three	3
Four	4
Five	5
Six	6

## French

<u>fid</u>	fText
1	Un
3	Trois
4	Quatre
5	Cinq
6	Siz
7	Sept
8	Huit

An "inner join":

```
SELECT *  
FROM English, French  
WHERE eid = fid
```

Same as (alternative join syntax):

```
SELECT *  
FROM English JOIN French  
ON eid = fid
```

etext	eid	fid	ftext
One	1	1	Un
Three	3	3	Trois
Four	4	4	Quatre
Five	5	5	Cinq
Six	6	6	Siz

← shortform of "INNER JOIN"

# Join Illustration



## English

eText	eid
One	1
Two	2
Three	3
Four	4
Five	5
Six	6

## French

fid	fText
1	Un
3	Trois
4	Quatre
5	Cinq
6	Six
7	Sept
8	Huit

How do we get a join with the full data ?

```
SELECT *  
FROM English JOIN French  
ON eid = fid
```

Null also sometimes just shown as empty

etext	eid	fid	ftext
One	1	1	Un
Two	2	NULL	NULL
Three	3	3	Trois
Four	4	4	Quatre
Five	5	5	Cinq
Six	6	6	Six
NULL	NULL	7	Sept
NULL	NULL	8	Huit



# Join Illustration



## English

eText	eid
One	1
Two	2
Three	3
Four	4
Five	5
Six	6

## French

fid	fText
1	Un
3	Trois
4	Quatre
5	Cinq
6	Siz
7	Sept
8	Huit

shortform of:  
"FULL OUTER JOIN"

```
SELECT *  
FROM English FULL JOIN French  
ON English.eid = French.fid
```

```
SELECT *  
FROM English JOIN French  
ON eid = fid
```

Null also sometimes  
just shown as empty

etext	eid	fid	ftext
One	1	1	Un
Two	2	NULL	NULL
Three	3	3	Trois
Four	4	4	Quatre
Five	5	5	Cinq
Six	6	6	Siz
NULL	NULL	7	Sept
NULL	NULL	8	Huit

# Join Illustration



## English

eText	<u>eid</u>
One	1
Two	2
Three	3
Four	4
Five	5
Six	6

## French

<u>fid</u>	fText
1	Un
3	Trois
4	Quatre
5	Cinq
6	Siz
<del>7</del>	<del>Sept</del>
<del>8</del>	<del>Huit</del>

```
SELECT *  
FROM English LEFT JOIN French  
ON English.eid = French.fid
```

etext	eid	fid	ftext
One	1	1	Un
Two	2	NULL	NULL
Three	3	3	Trois
Four	4	4	Quatre
Five	5	5	Cinq
Six	6	6	Siz

# Join Illustration

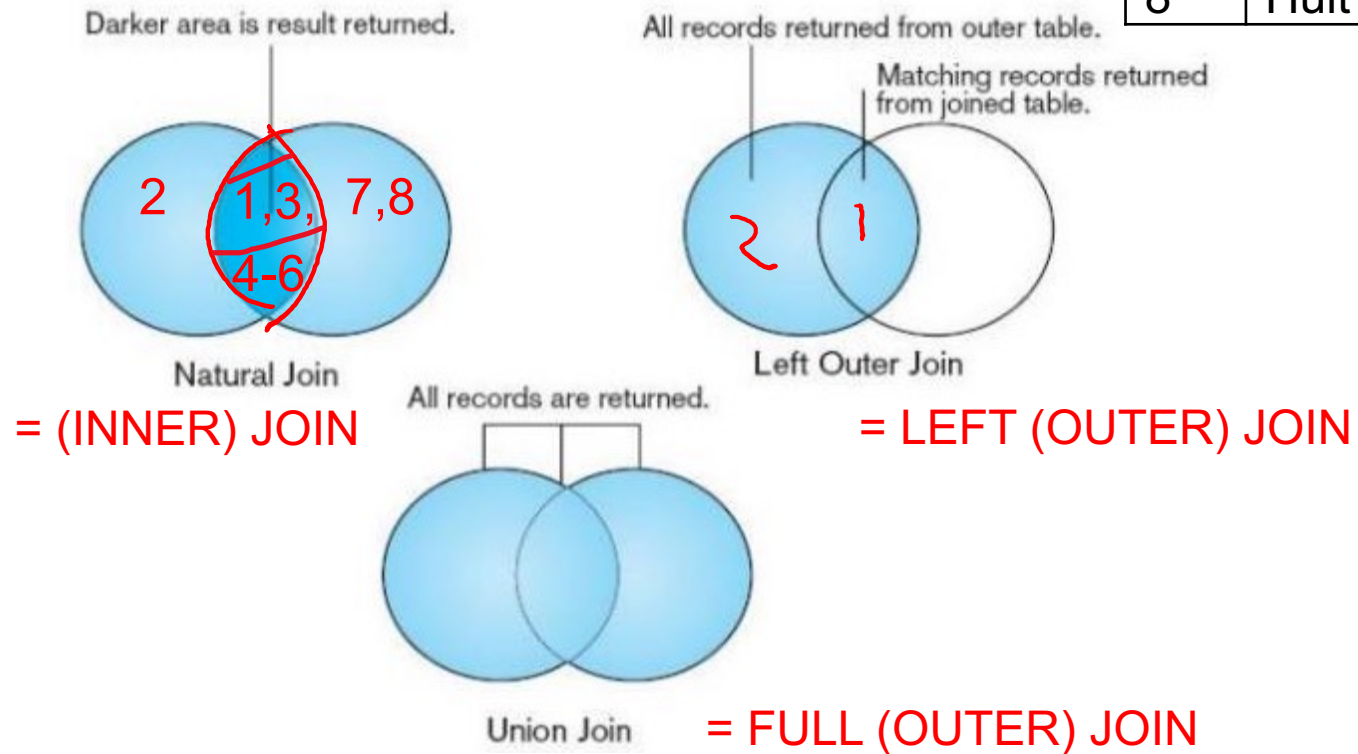


## English

eText	eid
One	1
Two	2
Three	3
Four	4
Five	5
Six	6

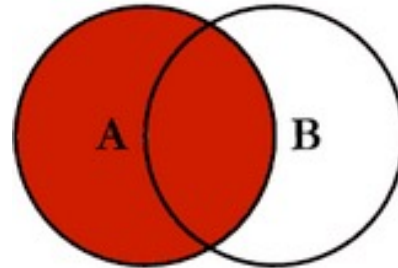
## French

fid	fText
1	Un
3	Trois
4	Quatre
5	Cinq
6	Siz
7	Sept
8	Huit

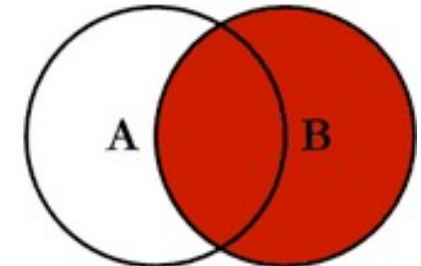


# Detailed Illustration with Examples (follow the link)

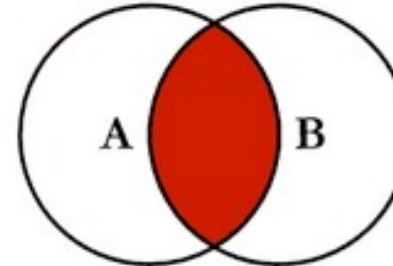
## SQL JOINS



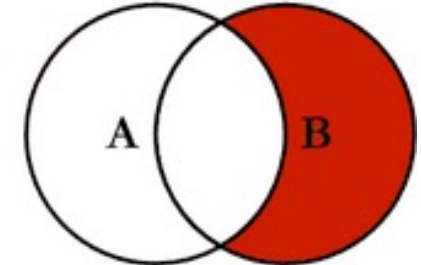
```
SELECT <select_list>  
FROM TableA A  
LEFT JOIN TableB B  
ON A.Key = B.Key
```



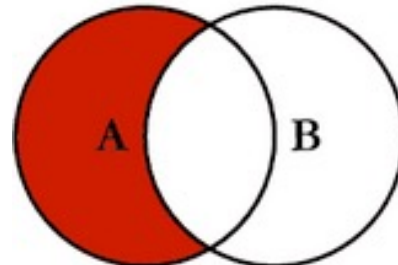
```
SELECT <select_list>  
FROM TableA A  
RIGHT JOIN TableB B  
ON A.Key = B.Key
```



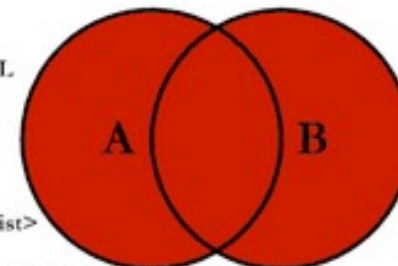
```
SELECT <select_list>  
FROM TableA A  
INNER JOIN TableB B  
ON A.Key = B.Key
```



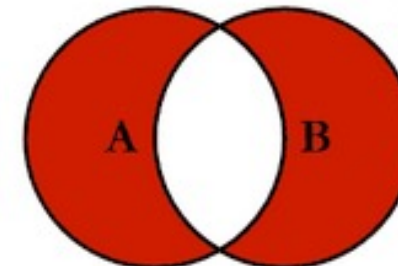
```
SELECT <select_list>  
FROM TableA A  
RIGHT JOIN TableB B  
ON A.Key = B.Key  
WHERE A.Key IS NULL
```



```
SELECT <select_list>  
FROM TableA A  
LEFT JOIN TableB B  
ON A.Key = B.Key  
WHERE B.Key IS NULL
```



```
SELECT <select_list>  
FROM TableA A  
FULL OUTER JOIN TableB B  
ON A.Key = B.Key
```



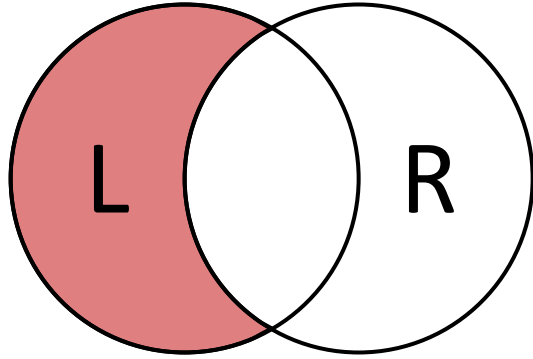
```
SELECT <select_list>  
FROM TableA A  
FULL OUTER JOIN TableB B  
ON A.Key = B.Key  
WHERE A.Key IS NULL  
OR B.Key IS NULL
```

© C.L. Moffitt, 2008

*also called "anti-join"*

Check this web page for illustrating examples

# Let's practice anti-joins



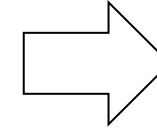
```
SELECT <select_list>
FROM L
LEFT JOIN R
ON L.key = R.key
WHERE R.key IS NULL
```

English

eText	<u>eid</u>
One	1
Two	2
Three	3
Four	4
Five	5
Six	6

French

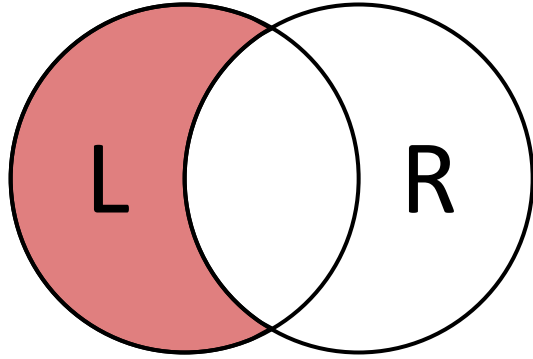
<u>fid</u>	fText
1	Un
3	Trois
4	Quatre
5	Cinq
6	Six
7	Sept
8	Huit



Results



# Let's practice anti-joins



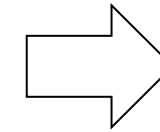
```
SELECT <select_list>
FROM L
LEFT JOIN R
ON L.key = R.key
WHERE R.key IS NULL
```

English

eText	<u>eid</u>
One	1
Two	2
Three	3
Four	4
Five	5
Six	6

French

<u>fid</u>	fText
1	Un
3	Trois
4	Quatre
5	Cinq
6	Six
7	Sept
8	Huit



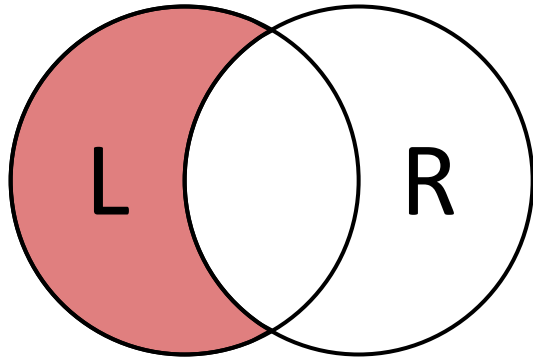
Results

eText	<u>eid</u>
Two	2

How to write in SQL?



# Let's practice anti-joins



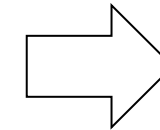
```
SELECT <select_list>
FROM L
LEFT JOIN R
ON L.key = R.key
WHERE R.key IS NULL
```

English

eText	eid
One	1
Two	2
Three	3
Four	4
Five	5
Six	6

French

fid	fText
1	Un
3	Trois
4	Quatre
5	Cinq
6	Six
7	Sept
8	Huit



Results

eText	eid
Two	2

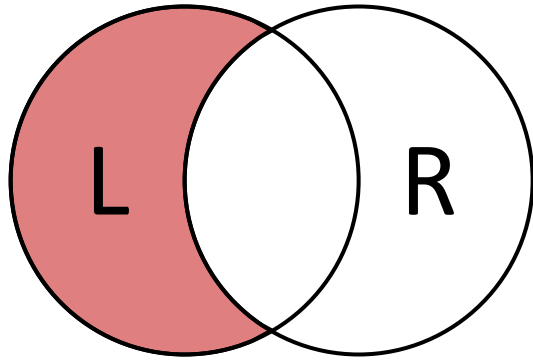
How to write in SQL?

```
SELECT eText, eid
FROM English
LEFT JOIN French
ON eid = fid
WHERE fid IS NULL
```

Any alternative?



# Let's practice anti-joins



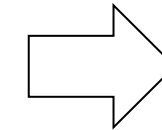
```
SELECT <select_list>
FROM L
LEFT JOIN R
ON L.key = R.key
WHERE R.key IS NULL
```

English

eText	<u>eid</u>
One	1
Two	2
Three	3
Four	4
Five	5
Six	6

French

<u>fid</u>	fText
1	Un
3	Trois
4	Quatre
5	Cinq
6	Six
7	Sept
8	Huit



Results

eText	<u>eid</u>
Two	2

How to write in SQL?

```
SELECT eText, eid
FROM English
LEFT JOIN French
ON eid = fid
WHERE fid IS NULL
```

Any alternative?

```
SELECT *
FROM English
WHERE eid NOT IN
(SELECT fid
FROM French)
```



# Semi-joins: kind of the anti-anti-joins...

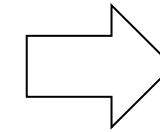


English

eText	<u>eid</u>
One	1
Two	2
Three	3
Four	4
Five	5
Six	6

French

<u>fid</u>	fText
1	Un
3	Trois
4	Quatre
5	Cinq
6	Six
7	Sept
8	Huit



Results

eText	<u>eid</u>
One	1
Three	3
Four	4
Five	5
Six	6

What do we have to change to these queries to get the tuples in English that have a partner in French?

?

```
SELECT eText, eid
FROM English
LEFT JOIN French
ON eid = fid
WHERE fid IS NULL
```

```
SELECT *
FROM English
WHERE eid NOT IN
(SELECT fid
FROM French)
```

# Semi-joins: kind of the anti-anti-joins...

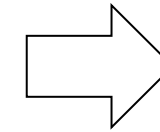


English

eText	<u>eid</u>
One	1
Two	2
Three	3
Four	4
Five	5
Six	6

French

<u>fid</u>	fText
1	Un
3	Trois
4	Quatre
5	Cinq
6	Six
7	Sept
8	Huit



Results

eText	<u>eid</u>
One	1
Three	3
Four	4
Five	5
Six	6

What do we have to change to these queries to get the tuples in English that have a partner in French?

What if fid is not a key?

?

```
SELECT eText, eid
FROM English
LEFT JOIN French
ON eid = fid
WHERE fid IS NOT NULL
```

```
SELECT *
FROM English
WHERE eid IN
  (SELECT fid
   FROM French)
```

# Semi-joins: kind of the anti-anti-joins...

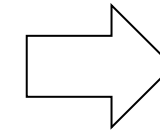


English

eText	<u>eid</u>
One	1
Two	2
Three	3
Four	4
Five	5
Six	6

French

<u>fid</u>	fText
1	Un
3	Trois
4	Quatre
5	Cinq
6	Six
7	Sept
8	Huit



Results

eText	<u>eid</u>
One	1
Three	3
Four	4
Five	5
Six	6

What do we have to change to these queries to get the tuples in English that have a partner in French?

What if fid is not a key?

**DISTINCT**

```
SELECT * eText, eid
FROM English
LEFT JOIN French
ON eid = fid
WHERE fid IS NOT NULL
```

These queries express "semi-joins"

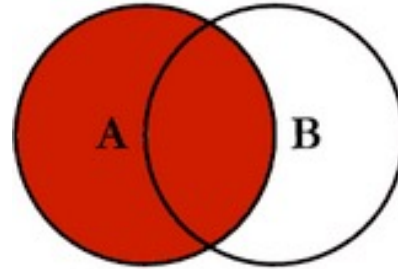
```
SELECT *
FROM English
WHERE eid IN
  (SELECT fid
   FROM French)
```

# Semi-joins: kind of the anti-anti-joins...

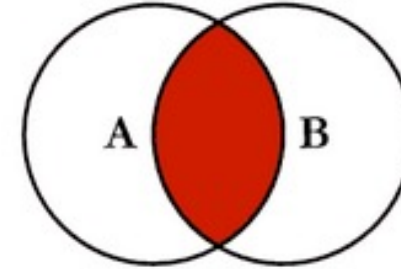


361

Which of these two VENN diagrams from earlier correspond to a semi-join?



```
SELECT <select_list>
FROM TableA A
LEFT JOIN TableB B
ON A.Key = B.Key
```



```
SELECT <select_list>
FROM TableA A
INNER JOIN TableB B
ON A.Key = B.Key
```

Get the tuples in English that have a partner in French?

**DISTINCT**

```
SELECT †eText, eid
FROM English
LEFT JOIN French
ON eid = fid
WHERE fid IS NOT NULL
```

These queries express "semi-joins"

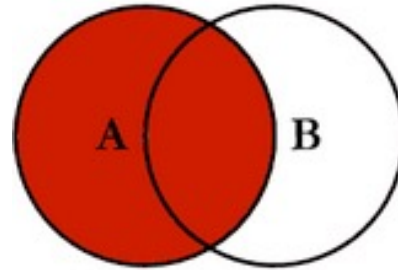
```
SELECT *
FROM English
WHERE eid IN
  (SELECT fid
   FROM French)
```

# Semi-joins: kind of the anti-anti-joins...

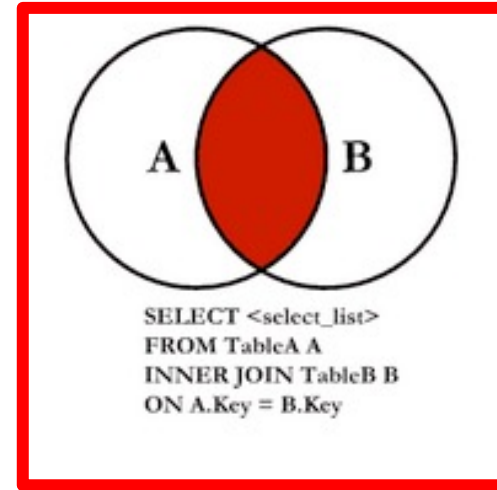


361

Which of these two VENN diagrams from earlier correspond to a semi-join?



```
SELECT <select_list>
FROM TableA A
LEFT JOIN TableB B
ON A.Key = B.Key
```



```
SELECT <select_list>
FROM TableA A
INNER JOIN TableB B
ON A.Key = B.Key
```

Only preserve tuples from A that also appear in B...

Get the tuples in English that have a partner in French?

**DISTINCT**

```
SELECT †eText, eid
FROM English
LEFT JOIN French
ON eid = fid
WHERE fid IS NOT NULL
```

These queries express "semi-joins"

```
SELECT *
FROM English
WHERE eid IN
  (SELECT fid
   FROM French)
```

# Another look at Outer Joins



```
SELECT *
FROM English FULL JOIN French
ON English.eid = French.fid
```

English

eText	eid
One	1
Two	2
Three	3
Four	4
Five	5
Six	6

French

fid	fText
1	Un
3	Trois
4	Quatre
5	Cinq
6	Six
7	Sept
8	Huit

*FULL JOIN can be written as union of inner join with anti-joins*

?

etext	eid	fid	ftext
One	1	1	Un
Two	2	NULL	NULL
Three	3	3	Trois
Four	4	4	Quatre
Five	5	5	Cinq
Six	6	6	Six
NULL	NULL	7	Sept
NULL	NULL	8	Huit

# Another look at Outer Joins



```
SELECT *
FROM English FULL JOIN French
ON English.eid = French.fid
```

```
SELECT etext, eid, fid, ftext
FROM English INNER JOIN French
ON English.eid = French.fid
```

```
UNION ALL
SELECT etext, eid, NULL, NULL
FROM English
WHERE NOT EXISTS(
  SELECT *
  FROM French
  WHERE eid=fid)
```

```
UNION ALL
SELECT NULL, NULL, fid, ftext
FROM French
WHERE NOT EXISTS(
  SELECT *
  FROM English
  WHERE eid=fid)
```

English

eText	eid
One	1
Two	2
Three	3
Four	4
Five	5
Six	6

French

fid	fText
1	Un
3	Trois
4	Quatre
5	Cinq
6	Siz
7	Sept
8	Huit

*anti-join*

etext	eid	fid	ftext
One	1	1	Un
Two	2	NULL	NULL
Three	3	3	Trois
Four	4	4	Quatre
Five	5	5	Cinq
Six	6	6	Siz
NULL	NULL	7	Sept
NULL	NULL	8	Huit

# Are these two queries equivalent?

Q<sub>2</sub>: Find all companies that make only products with price < 25

```
SELECT DISTINCT C.cname
FROM   Company C
WHERE  C.cid NOT IN ( SELECT P.cid
                     FROM   Product P
                     WHERE  P.price >= 25)
```

Is the following query identical to the one above?

```
SELECT C.name
FROM   Company C, Product P
WHERE  C.cid=P.cid
GROUP BY cname
HAVING MAX(P.price) < 25
```



# Are these two queries equivalent?

Q<sub>2</sub>: Find all companies that make only products with price < 25

```
SELECT DISTINCT C.cname
FROM   Company C
WHERE  C.cid NOT IN ( SELECT P.cid
                     FROM   Product P
                     WHERE  P.price >= 25)
```



Almost, but not really.  
The upper query would  
return a company that  
has no product, the  
one below would not.

Is the following query identical to the one above?

```
SELECT C.name
FROM   Company C, Product P
WHERE  C.cid=P.cid
GROUP BY cname
HAVING MAX(P.price) < 25
```

# Outer Joins, Coalesce, and non-associativity

# Coalesce function



M	N
a	a
1	2
2	3

```
SELECT M.a, N.a, COALESCE(M.a, N.a) as b
FROM M
FULL JOIN N
ON M.a = N.a
```

COALESCE: takes first non-NULL value,

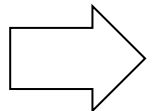
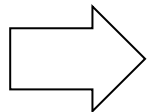
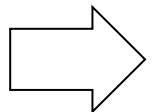
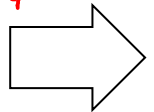
(<sup>1</sup>N <sup>2</sup>N <sup>3</sup>9 <sup>4</sup>8)

```
SELECT COALESCE(1, NULL)
```

```
SELECT COALESCE(NULL, 3)
```

```
SELECT COALESCE(1, 2)
```

```
SELECT COALESCE(NULL, NULL)
```



?

SQL example available at: <https://github.com/northeastern-datalab/cs3200-activities/tree/master/sql>

Also see use of COALESCE across programming languages: [https://en.wikipedia.org/wiki/Null\\_coalescing\\_operator](https://en.wikipedia.org/wiki/Null_coalescing_operator)

Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>

# Coalesce function

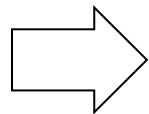


M	N
a	a
1	2
2	3

```
SELECT M.a, N.a, COALESCE(M.a, N.a) as b
FROM M
FULL JOIN N
ON M.a = N.a
```

Result

M.a	N.a	b
-----	-----	---

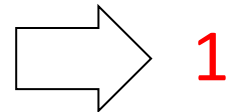


?

?

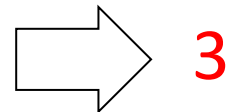
COALESCE: takes first non-NULL value,  
 $C(x,y,z) = C(x,C(y,z)) = C(C(x,y),z)$

```
SELECT COALESCE(1, NULL)
```



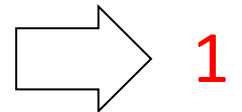
1

```
SELECT COALESCE(NULL, 3)
```



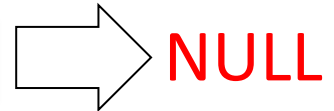
3

```
SELECT COALESCE(1, 2)
```



1

```
SELECT COALESCE(NULL, NULL)
```



NULL

SQL example available at: <https://github.com/northeastern-datalab/cs3200-activities/tree/master/sql>

Also see use of COALESCE across programming languages: [https://en.wikipedia.org/wiki/Null\\_coalescing\\_operator](https://en.wikipedia.org/wiki/Null_coalescing_operator)

Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>

# Coalesce function



M	N
a	a
1	2
2	3

```
SELECT M.a, N.a, COALESCE(M.a, N.a) as b
FROM M
FULL JOIN N
ON M.a = N.a
```

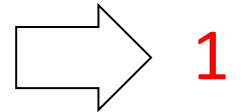


## Result

M.a	N.a	b
1	NULL	1
2	2	2
NULL	3	3

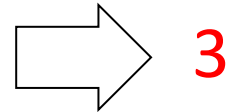
COALESCE: takes first non-NULL value,  
 $C(x,y,z) = C(x,C(y,z)) = C(C(x,y),z)$

```
SELECT COALESCE(1, NULL)
```



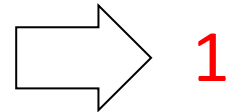
1

```
SELECT COALESCE(NULL, 3)
```



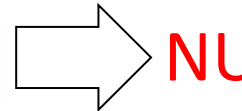
3

```
SELECT COALESCE(1, 2)
```



1

```
SELECT COALESCE(NULL, NULL)
```



NULL

SQL example available at: <https://github.com/northeastern-datalab/cs3200-activities/tree/master/sql>

Also see use of COALESCE across programming languages: [https://en.wikipedia.org/wiki/Null\\_coalescing\\_operator](https://en.wikipedia.org/wiki/Null_coalescing_operator)

Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>

# Coalesce, Natural Outer Join, Union



M	N
a	a
1	2
2	3

```
SELECT *  
FROM M  
NATURAL FULL JOIN N
```

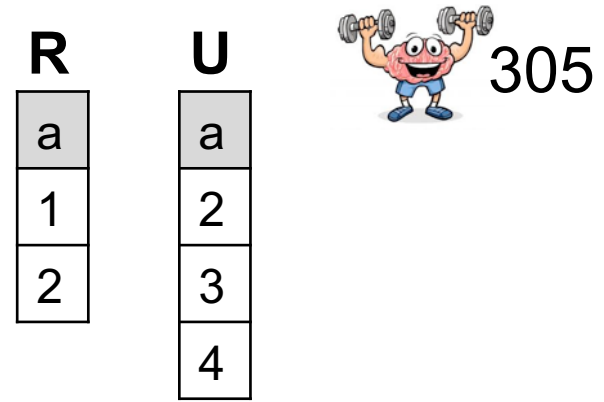
Result

a
1
2
3

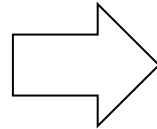
Natural full join models "coalesce"

Join vs. Union – it is actually the same:  
Union is a special case of a join 😊  
(under set semantics)

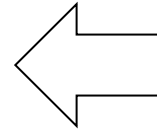
# Do we need "union"?



```
SELECT a
FROM R
UNION
SELECT a
FROM U
```



a
1
2
3
4



# Do we need "union"?



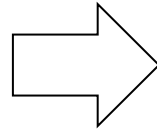
R
a
1
2

U
a
2
3
4

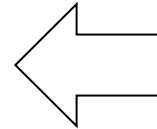
Coalesce (R.A, U.A)

1	U
2	2
U	3
U	4

```
SELECT a
FROM R
UNION
SELECT a
FROM U
```



a
1
2
3
4



```
SELECT COALESCE(R.a,U.a) as a
FROM R FULL JOIN U
on R.a = U.a
```



# Quick recap: Commutativity & Associativity

Multiplication

$$3 \cdot 2 \cdot 4 = 24$$

?

Multiplication is  
associative 😊

Matrix multiplication

$$\begin{bmatrix} 1 & 1 \\ 2 & 3 \end{bmatrix} \cdot \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \cdot \begin{bmatrix} 3 \\ 1 \end{bmatrix} =$$

# Quick recap: Commutativity & Associativity

## Multiplication

$$(3 \cdot 2) \cdot 4 = 24$$

*Order of operations can be exchanged:*

$$3 \cdot (2 \cdot 4) = 24$$

Multiplication is associative 😊

?

and commutative 😊

## Matrix multiplication

$$\begin{bmatrix} 1 & 1 \\ 2 & 3 \end{bmatrix} \cdot \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \cdot \begin{bmatrix} 3 \\ 1 \end{bmatrix} =$$

# Quick recap: Commutativity & Associativity

## Multiplication

$$(3 \cdot 2) \cdot 4 = 24$$

Order of operations can be exchanged:

$$3 \cdot (2 \cdot 4) = 24$$

Multiplication is associative 😊

Order of operands can be exchanged:

$$4 \cdot 2$$

and commutative 😊

## Matrix multiplication

$$\begin{bmatrix} 1 & 1 \\ 2 & 3 \end{bmatrix} \cdot \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \cdot \begin{bmatrix} 3 \\ 1 \end{bmatrix} =$$

?

# Quick recap: Commutativity & Associativity

## Multiplication

$$(3 \cdot 2) \cdot 4 = 24$$

Order of operations can be exchanged:

$$3 \cdot (2 \cdot 4) = 24$$

Multiplication is associative 😊

Order of operands can be exchanged:

$$4 \cdot 2$$

and commutative 😊

## Matrix multiplication

$$\begin{pmatrix} 1 & 1 \\ 2 & 3 \end{pmatrix} \cdot \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \cdot \begin{pmatrix} 3 \\ 1 \end{pmatrix} = \begin{pmatrix} 18 \\ 49 \end{pmatrix}$$

$$\begin{pmatrix} 1 & 1 \\ 2 & 3 \end{pmatrix} \cdot \left( \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \cdot \begin{pmatrix} 3 \\ 1 \end{pmatrix} \right) = \begin{pmatrix} 18 \\ 49 \end{pmatrix}$$

Matrix multipl. is associative 😊

$$\begin{pmatrix} 3 \\ 1 \end{pmatrix} \cdot \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}$$

#col ≠ #row

... but \*not\* commutative 😞

It turns out this is mainly a problem of syntax, not semantics. Einstein notation (and similar more recent extensions like "EINSUM") solves that. See e.g. Laue et al. *A Simple and Efficient Tensor Calculus*. AAAI 2020. <https://arxiv.org/abs/2010.03313>. Alternatively, think about the relational join operator as a commutative notation for sparse matrix multiplication (also Cartesian product under named or unnamed perspective)

# The power of associativity

Option 1:  $\left( \begin{array}{|c|c|} \hline 1 & 1 \\ \hline 2 & 3 \\ \hline \end{array} \cdot \begin{array}{|c|c|} \hline 1 & 2 \\ \hline 3 & 4 \\ \hline \end{array} \right) \cdot \begin{array}{|c|} \hline 3 \\ \hline 1 \\ \hline \end{array} = \begin{array}{|c|} \hline 18 \\ \hline 49 \\ \hline \end{array}$

Option 2:  $\begin{array}{|c|c|} \hline 1 & 1 \\ \hline 2 & 3 \\ \hline \end{array} \cdot \left( \begin{array}{|c|c|} \hline 1 & 2 \\ \hline 3 & 4 \\ \hline \end{array} \cdot \begin{array}{|c|} \hline 3 \\ \hline 1 \\ \hline \end{array} \right) = \begin{array}{|c|} \hline 18 \\ \hline 49 \\ \hline \end{array}$

Which option would you choose to evaluate this matrix multiplication ?

# The power of associativity

Option 1:  $\left( \begin{array}{|c|c|} \hline 1 & 1 \\ \hline 2 & 3 \\ \hline \end{array} \cdot \begin{array}{|c|c|} \hline 1 & 2 \\ \hline 3 & 4 \\ \hline \end{array} \right) \cdot \begin{array}{|c|} \hline 3 \\ \hline 1 \\ \hline \end{array} = \begin{array}{|c|} \hline 18 \\ \hline 49 \\ \hline \end{array}$

$\begin{array}{|c|c|} \hline 4 & 6 \\ \hline 11 & 16 \\ \hline \end{array}$

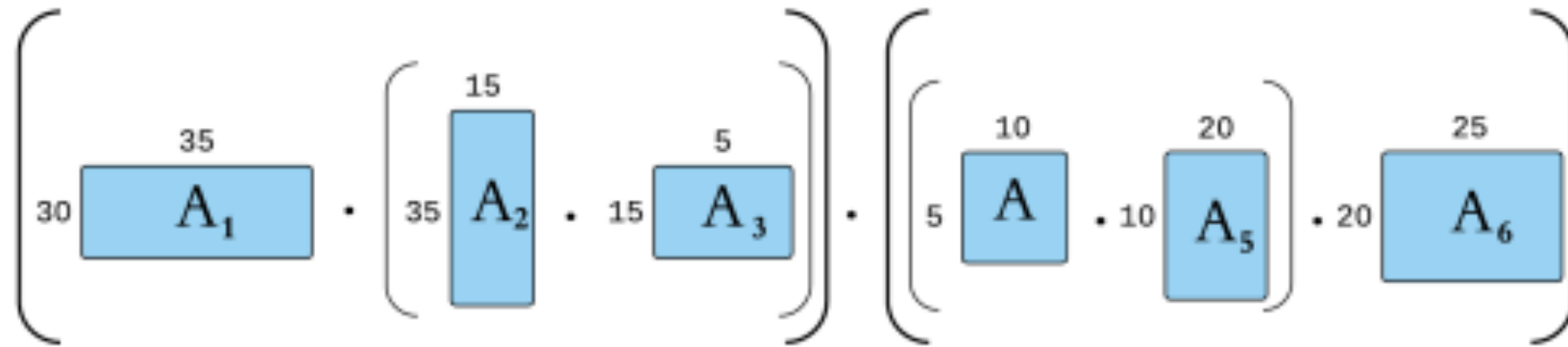
Option 2:  $\begin{array}{|c|c|} \hline 1 & 1 \\ \hline 2 & 3 \\ \hline \end{array} \cdot \left( \begin{array}{|c|c|} \hline 1 & 2 \\ \hline 3 & 4 \\ \hline \end{array} \cdot \begin{array}{|c|} \hline 3 \\ \hline 1 \\ \hline \end{array} \right) = \begin{array}{|c|} \hline 18 \\ \hline 49 \\ \hline \end{array}$

$\begin{array}{|c|} \hline 5 \\ \hline 13 \\ \hline \end{array}$

*All variants give the same result. But some are faster.  
Intuition: we like to have small intermediate result sizes!*

# Matrix chain multiplication

Given  $n$  matrices, what is the optimal sequence to multiply them? **?**



This is an example  
optimal factorization. **?**  
What is its cost?

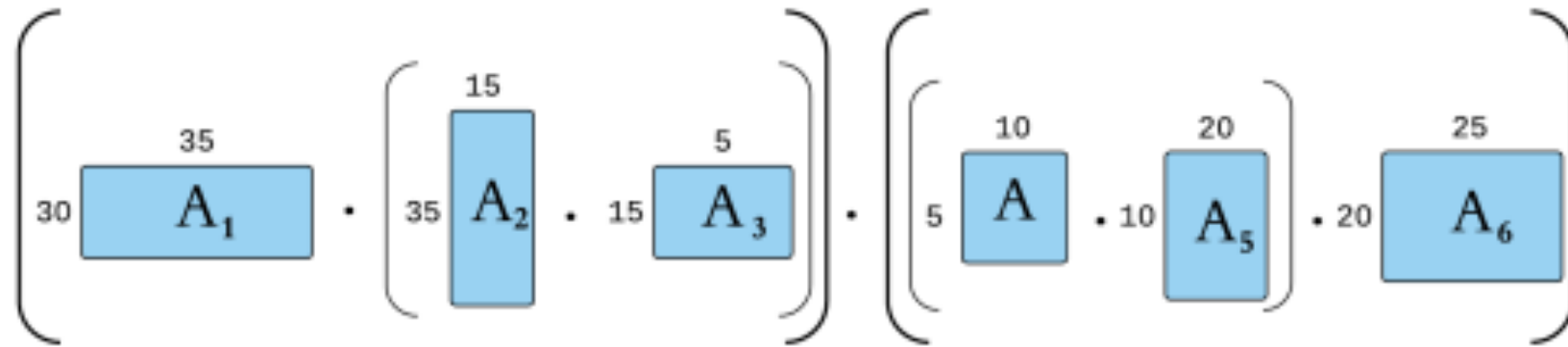
See also [https://en.wikipedia.org/wiki/Catalan\\_number](https://en.wikipedia.org/wiki/Catalan_number), [https://en.wikipedia.org/wiki/Matrix\\_chain\\_multiplication](https://en.wikipedia.org/wiki/Matrix_chain_multiplication), [https://en.wikipedia.org/wiki/Matrix\\_multiplication#Associativity](https://en.wikipedia.org/wiki/Matrix_multiplication#Associativity)

Source figure: <https://bruceoutdoors.wordpress.com/2015/11/24/matrix-chain-multiplication-with-c-code-part-3-extracting-the-sequence/>

Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>

# Matrix chain multiplication

Given  $n$  matrices, what is the optimal sequence to multiply them? ?



This is an example optimal factorization. ?  
What is its cost?

$$\text{MinCost: } (30 \cdot 35 \cdot 5 + (35 \cdot 15 \cdot 5)) + 30 \cdot 5 \cdot 25 + (5 \cdot 10 \cdot 20) + 5 \cdot 20 \cdot 25$$

Nave method: all possible way to place closed parentheses: "Catalan numbers"

Via Dynamic programming:  $O(n^3)$

Best known:  $O(n \log n)$

$C_n$  is the number of different ways  $n + 1$  factors can be completely parenthesized (or the number of ways of associating  $n$  applications of a binary operator, as in the matrix chain multiplication problem). For  $n = 3$ , for example, we have the following five different parenthesizations of four factors:

$$((ab)c)d \quad (a(bc))d \quad (ab)(cd) \quad a((bc)d) \quad a(b(cd))$$

See also [https://en.wikipedia.org/wiki/Catalan\\_number](https://en.wikipedia.org/wiki/Catalan_number), [https://en.wikipedia.org/wiki/Matrix\\_chain\\_multiplication](https://en.wikipedia.org/wiki/Matrix_chain_multiplication), [https://en.wikipedia.org/wiki/Matrix\\_multiplication#Associativity](https://en.wikipedia.org/wiki/Matrix_multiplication#Associativity)

Source figure: <https://bruceoutdoors.wordpress.com/2015/11/24/matrix-chain-multiplication-with-c-code-part-3-extracting-the-sequence/>

Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>



# Commutativity & Associativity



333

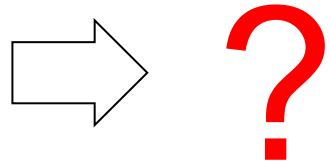
Outer joins

R		S		T	
A	B	B	C	A	C
1	2	2	3	4	5

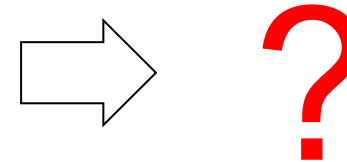
```
SELECT A, B, C
FROM (R
NATURAL FULL JOIN S)
NATURAL FULL JOIN T
```

```
SELECT A, B, C
FROM R
NATURAL FULL JOIN (S
NATURAL FULL JOIN T)
```

Result



Result



# Commutativity & Associativity



333

Outer joins

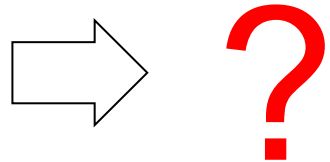
R		S		T	
A	B	B	C	A	C
1	2	2	3	4	5

A	B	C
1	2	3

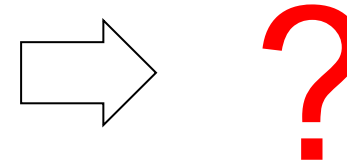
```
SELECT A, B, C
FROM (R
NATURAL FULL JOIN S)
NATURAL FULL JOIN T
```

```
SELECT A, B, C
FROM R
NATURAL FULL JOIN (S
NATURAL FULL JOIN T)
```

Result



Result



# Commutativity & Associativity



333

Outer joins

R		S		T	
A	B	B	C	A	C
1	2	2	3	4	5

A	B	C
1	2	3

```
SELECT A, B, C
FROM (R
NATURAL FULL JOIN S)
NATURAL FULL JOIN T
```

```
SELECT A, B, C
FROM R
NATURAL FULL JOIN (S
NATURAL FULL JOIN T)
```

Result

A	B	C
1	2	3
4	NULL	5

Result

?

# Commutativity & Associativity



333

Outer joins

R		S		T	
A	B	B	C	A	C
1	2	2	3	4	5

A	B	C
1	2	3

A	B	C
1	2	3
4		5

```
SELECT A, B, C
FROM (R
NATURAL FULL JOIN S)
NATURAL FULL JOIN T
```

```
SELECT A, B, C
FROM R
NATURAL FULL JOIN (S
NATURAL FULL JOIN T)
```

Result

A	B	C
1	2	3
4	NULL	5

Result

?

# Commutativity & Associativity



Outer joins

R		S		T	
A	B	B	C	A	C
1	2	2	3	4	5

A	B	C
1	2	3

A	B	C
1	2	3
4		5

```
SELECT A, B, C
FROM (R
NATURAL FULL JOIN S)
NATURAL FULL JOIN T
```

```
SELECT A, B, C
FROM R
NATURAL FULL JOIN (S
NATURAL FULL JOIN T)
```

Result

A	B	C
1	2	3
4	NULL	5

Result

A	B	C
1	2	NULL
NULL	2	3
4	NULL	5

Thus outer joins are not associative! (but they are commutative)

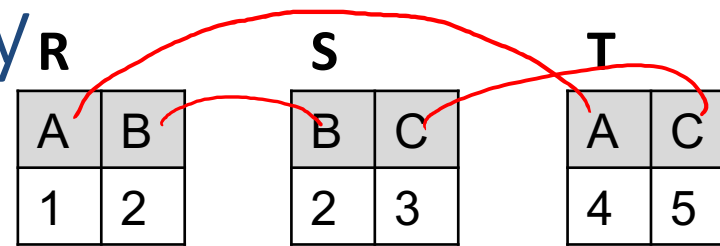
SQL example available at: <https://github.com/northeastern-datalab/cs3200-activities/tree/master/sql>

Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>

# Commutativity & Associativity



## Outer joins



```
SELECT R.a RA, T.a TA, coalesce(R.a, T.a) a,
       R.b RB, S.b SB, coalesce(R.b, S.b) b,
       S.c SC, T.c TC, coalesce(S.c, T.c) c
FROM (R
FULL JOIN S on R.B=S.B)
FULL JOIN T on S.C=T.C AND R.A = T.A
```

```
SELECT R.a RA, T.a TA, coalesce(R.a, T.a) a,
       R.b RB, S.b SB, coalesce(R.b, S.b) b,
       S.c SC, T.c TC, coalesce(S.c, T.c) c
FROM R
FULL JOIN (S
FULL JOIN T on S.C=T.C) on R.B=S.B AND R.A = T.A
```

Data Output	ra	ta	a	rb	sb	b	sc	tc	c
1	1	[null]	1	2	2	2	3	[null]	3
2	[null]	4	4	[null]	[null]	[null]	[null]	5	5

Data Output	ra	ta	a	rb	sb	b	sc	tc	c
1	[null]	[null]	[null]	[null]	2	2	3	[null]	3
2	[null]	4	4	[null]	[null]	[null]	[null]	5	5
3	1	[null]	1	2	[null]	2	[null]	[null]	[null]

# Example: Data Sources on Tourist Information



## Climates

Country	Climate
Canada	diverse
Bahamas	tropical
UK	temperate

## Accommodations

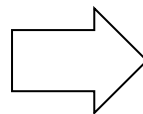
Country	City	Hotel	Stars
Canada	Toronto	Plaza	4
Canada	London	Ramada	3
Bahamas	Nassau	Hilton	

## Sites

Country	City	Site
Canada	London	Air show
Canada		Mount Logan
UK	London	Buckingham
UK	London	Hyde Park

```
SELECT *  
FROM (Accommodations  
NATURAL FULL JOIN Climates)  
NATURAL FULL JOIN Sites
```

Result



# Example: Data Sources on Tourist Information



## Climates

Country	Climate
Canada	diverse
Bahamas	tropical
UK	temperate

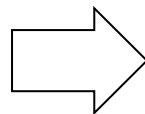
## Accommodations

Country	City	Hotel	Stars
Canada	Toronto	Plaza	4
Canada	London	Ramada	3
Bahamas	Nassau	Hilton	

## Sites

Country	City	Site
Canada	London	Air show
Canada		Mount Logan
UK	London	Buckingham
UK	London	Hyde Park

```
SELECT *
FROM (Accommodations
NATURAL FULL JOIN Climates)
NATURAL FULL JOIN Sites
```



## Result

Country	City	Climate	Hotel	Stars	Site
Canada	Toronto	diverse	Plaza	4	
Canada	London	diverse	Ramada	3	Air Show
Canada					Mount Logan
UK	London				Buckingham
UK	London				Hyde Park
UK		temperate			
Bahamas	Nassau	Tropical	Hilton		

SQL example available at: <https://github.com/northeastern-datalab/cs3200-activities/tree/master/sql>

Example from: Cohen, Fadida, Kanza, Kimelfeld, Sagiv. "Full Disjunctions: Polynomial-Delay Iterators in Action", VLDB 2006. <http://vldb.org/conf/2006/p739-cohen.pdf>

Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>



# Example: Data Sources on Tourist Information



## Climates

Country	Climate
Canada	diverse
Bahamas	tropical
UK	temperate

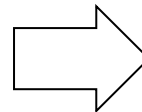
## Accommodations

Country	City	Hotel	Stars
Canada	Toronto	Plaza	4
Canada	London	Ramada	3
Bahamas	Nassau	Hilton	

## Sites

Country	City	Site
Canada	London	Air show
Canada		Mount Logan
UK	London	Buckingham
UK	London	Hyde Park

```
SELECT *  
FROM (Accommodations  
NATURAL FULL JOIN Climates)  
NATURAL FULL JOIN Sites
```



## Result

Country	City	Climate	Hotel	Stars	Site
Canada	Toronto	diverse	Plaza	4	
Canada	London	diverse	Ramada	3	Air Show
Canada					Mount Logan
UK	London				Buckingham
UK	London				Hyde Park
UK		temperate			
Bahamas	Nassau	Tropical	Hilton		

SQL example available at: <https://github.com/northeastern-datalab/cs3200-activities/tree/master/sql>

Example from: Cohen, Fadida, Kanza, Kimelfeld, Sagiv. "Full Disjunctions: Polynomial-Delay Iterators in Action", VLDB 2006. <http://vldb.org/conf/2006/p739-cohen.pdf>

Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>

# Example: Data Sources on Tourist Information



## Climates

Country	Climate
Canada	diverse
Bahamas	tropical
<del>UK</del>	<del>temperate</del>

## Accommodations

Country	City	Hotel	Stars
Canada	Toronto	Plaza	4
Canada	London	Ramada	3
Bahamas	Nassau	Hilton	

## Sites

Country	City	Site
Canada	London	Air show
Canada		Mount Logan
UK	London	Buckingham
UK	London	Hyde Park

```
SELECT *
FROM (Accommodations
NATURAL FULL JOIN Climates)
NATURAL FULL JOIN Sites
```

## Result

Country	City	Climate	Hotel	Stars	Site
Canada	Toronto	diverse	Plaza	4	
Canada	London	diverse	Ramada	3	Air Show
Canada					Mount Logan
UK	London				Buckingham
UK	London				Hyde Park
UK		temperate			
Bahamas	Nassau	Tropical	Hilton		

SQL example available at: <https://github.com/northeastern-datalab/cs3200-activities/tree/master/sql>

Example from: Cohen, Fadida, Kanza, Kimelfeld, Sagiv. "Full Disjunctions: Polynomial-Delay Iterators in Action", VLDB 2006. <http://vldb.org/conf/2006/p739-cohen.pdf>

Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>

# Example: Data Sources on Tourist Information



## Climates

Country	Climate
Canada	diverse
Bahamas	tropical
UK	temperate

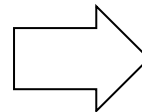
## Accommodations

Country	City	Hotel	Stars
Canada	Toronto	Plaza	4
Canada	London	Ramada	3
Bahamas	Nassau	Hilton	

## Sites

Country	City	Site
Canada	London	Air show
Canada		Mount Logan
UK	London	Buckingham
UK	London	Hyde Park

```
SELECT *
FROM Accommodations
NATURAL FULL JOIN (Climates
NATURAL FULL JOIN Sites)
```



## Result

Country	City	Climate	Hotel	Stars	Site
Canada	Toronto		Plaza	4	
Canada	London	diverse	Ramada	3	Air Show
Canada		diverse			Mount Logan
UK	London	temperate			Buckingham
UK	London	temperate			Hyde Park
Bahamas		Tropical			
Bahamas	Nassau		Hilton		

SQL example available at: <https://github.com/northeastern-datalab/cs3200-activities/tree/master/sql>

Example from: Cohen, Fadida, Kanza, Kimelfeld, Sagiv. "Full Disjunctions: Polynomial-Delay Iterators in Action", VLDB 2006. <http://vldb.org/conf/2006/p739-cohen.pdf>

Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>

# Example: Data Sources on Tourist Information



## Climates

Country	Climate
Canada	diverse
Bahamas	tropical
UK	temperate

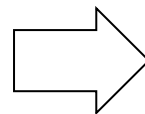
## Accommodations

Country	City	Hotel	Stars
Canada	Toronto	Plaza	4
Canada	London	Ramada	3
Bahamas	Nassau	Hilton	

## Sites

Country	City	Site
Canada	London	Air show
Canada		Mount Logan
UK	London	Buckingham
UK	London	Hyde Park

```
SELECT *  
FROM Accommodations  
NATURAL FULL JOIN (Climates  
NATURAL FULL JOIN Sites)
```



## Result

Country	City	Climate	Hotel	Stars	Site
Canada	Toronto		Plaza	4	
Canada	London	diverse	Ramada	3	Air Show
Canada		diverse			Mount Logan
UK	London	temperate			Buckingham
UK	London	temperate			Hyde Park
Bahamas		Tropical			
Bahamas	Nassau		Hilton		

SQL example available at: <https://github.com/northeastern-datalab/cs3200-activities/tree/master/sql>

Example from: Cohen, Fadida, Kanza, Kimelfeld, Sagiv. "Full Disjunctions: Polynomial-Delay Iterators in Action", VLDB 2006. <http://vldb.org/conf/2006/p739-cohen.pdf>

Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>

# Full disjunction



## Climates

Country	Climate
Canada	diverse
Bahamas	tropical
UK	temperate

## Accommodations

Country	City	Hotel	Stars
Canada	Toronto	Plaza	4
Canada	London	Ramada	3
Bahamas	Nassau	Hilton	

## Sites

Country	City	Site
Canada	London	Air show
Canada		Mount Logan
UK	London	Buckingham
UK	London	Hyde Park

```
SELECT *  
FROM FULL DISJUNCTION(Climates,  
(Accommodations, Sites))
```

FD: variation of the join operator that maximally combines join consistent tuples from connected relations, while preserving all information in the relations.

Not available in SQL! We may discuss later in class in more detail (or skip this year)

## Result

Country	City	Climate	Hotel	Stars	Site
Canada	Toronto	diverse	Plaza	4	
Canada	London	diverse	Ramada	3	Air Show
Canada		diverse			Mount Logan
UK	London	temperate			Buckingham
UK	London	temperate			Hyde Park
Bahamas	Nassau	tropical	Hilton		

SQL example available at: <https://github.com/northeastern-datalab/cs3200-activities/tree/master/sql>

Example from: Cohen, Fadida, Kanza, Kimelfeld, Sagiv. "Full Disjunctions: Polynomial-Delay Iterators in Action", VLDB 2006. <http://vldb.org/conf/2006/p739-cohen.pdf>

Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>

# Full disjunction: definition

## Climates

Country	Climate
Canada	diverse
Bahamas	tropical
UK	temperate

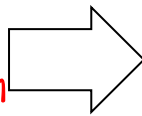
## Accommodations

Country	City	Hotel	Stars
Canada	Toronto	Plaza	4
Canada	London	Ramada	3
Bahamas	Nassau	Hilton	

## Sites

Country	City	Site
Canada	London	Air show
Canada		Mount Logan
UK	London	Buckingham
UK	London	Hyde Park

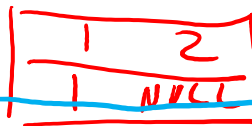
- Two tuples (max one from each relation) are join consistent if they agree on common attributes, e.g.  $t_1/t_2$ ,  $t_3/t_4$ . A set of tuples is join consistent if every pair is join consistent.
- Set of tuples (max one from each relation) is connected if the schema is connected, thus share attributes
- A tuple is in the Full disjunction if it is the inner join from tuples that are connected, join consistent, and there is no superset with both conditions (related to "subsumption").



## Result

Country	City	Climate	Hotel	Stars	Site
Canada	Toronto	diverse	Plaza	4	
Canada	London	diverse	Ramada	3	Air Show
Canada		diverse			Mount Logan
UK	London	temperate			Buckingham
UK	London	temperate			Hyde Park
Bahamas	Nassau	tropical	Hilton		

SUBSUMPTION



# Outline: T1-U1: SQL

- SQL

- Schema, keys, referential integrity
- Joins
- Aggregates and grouping
- Nested queries (Subqueries)
- Union and Theta Joins
- Nulls & Outer joins
- Top-k
- [Recursion: moved to T1-U4: Datalog]

# Sorting & Top- $k$ evaluation with SQL



*R*

<i>A</i>	<i>B</i>
1	0
2	0
3	0
4	1

*S*

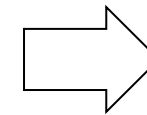
<i>B</i>	<i>C</i>
0	1
0	1
0	1
0	2

*T*

<i>C</i>	<i>D</i>
1	1
1	2
2	3
2	4

```
select A, R.B, S.C, D
from R, S, T
where R.B=S.B and S.C=T.C
```

**Result**

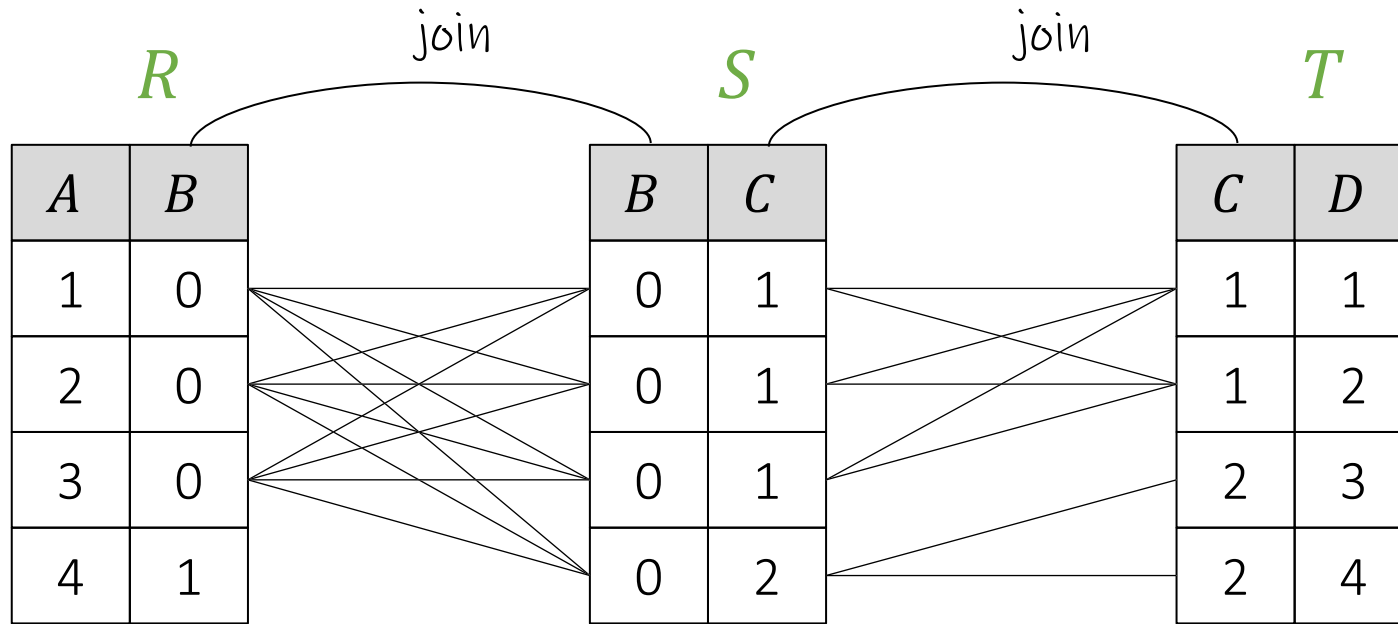


How many results do we get?



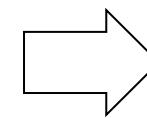


# Sorting & Top-k evaluation with SQL



```
select A, R.B, S.C, D
from R, S, T
where R.B=S.B and S.C=T.C
```

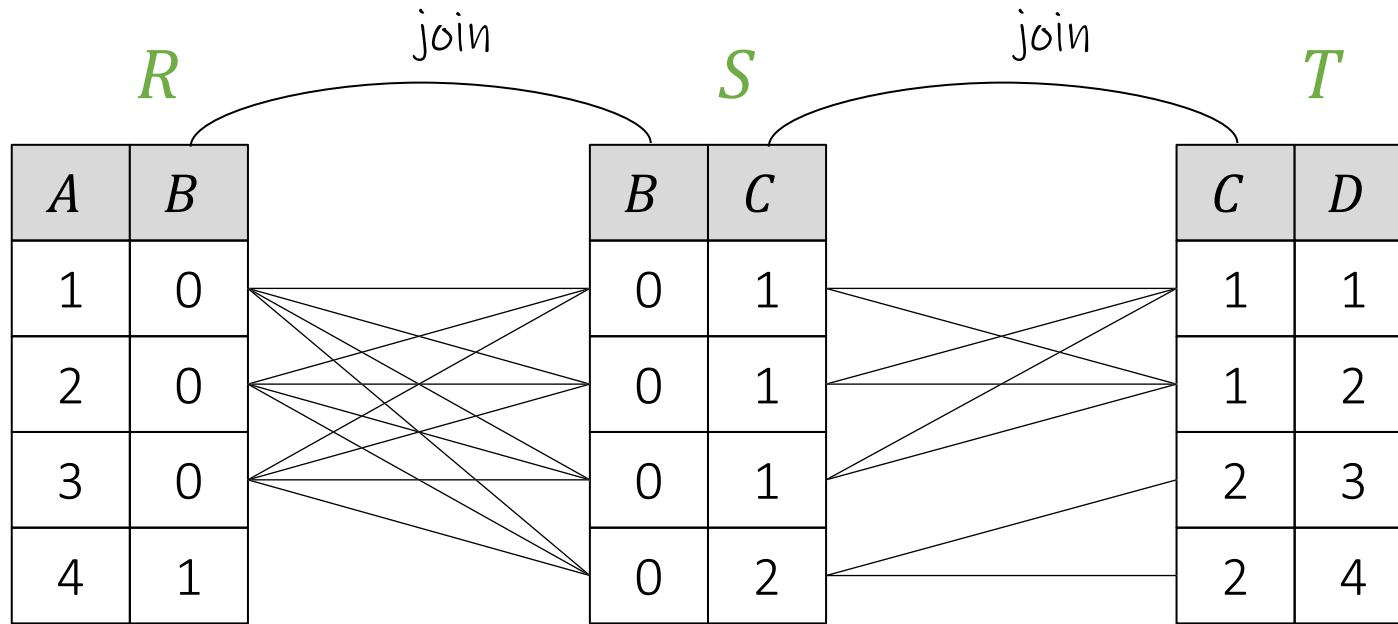
**Result**



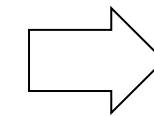
How many results do we get?



# Sorting & Top-k evaluation with SQL



```
select A, R.B, S.C, D
from R, S, T
where R.B=S.B and S.C=T.C
```

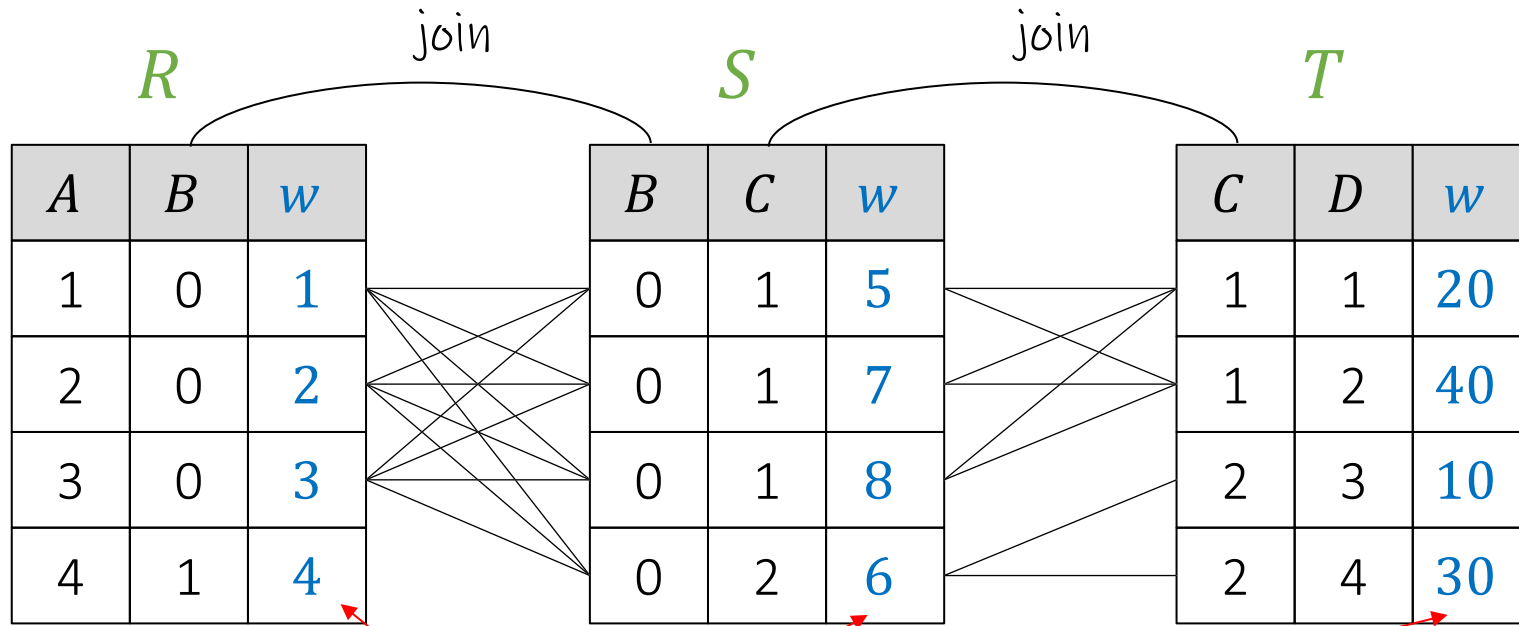


24 total results

## Result

A	B	C	D
1	0	2	3
2	0	2	3
3	0	2	3
1	0	1	1
2	0	1	1
3	0	1	1
1	0	1	1
...	...	...	...

# Sorting & Top-k evaluation with SQL



```
select  A, R.B, S.C, D,  
        R.w + S.w + T.w as weight  
from    R, S, T  
where   R.B=S.B and S.C=T.C  
order  by weight ASC
```

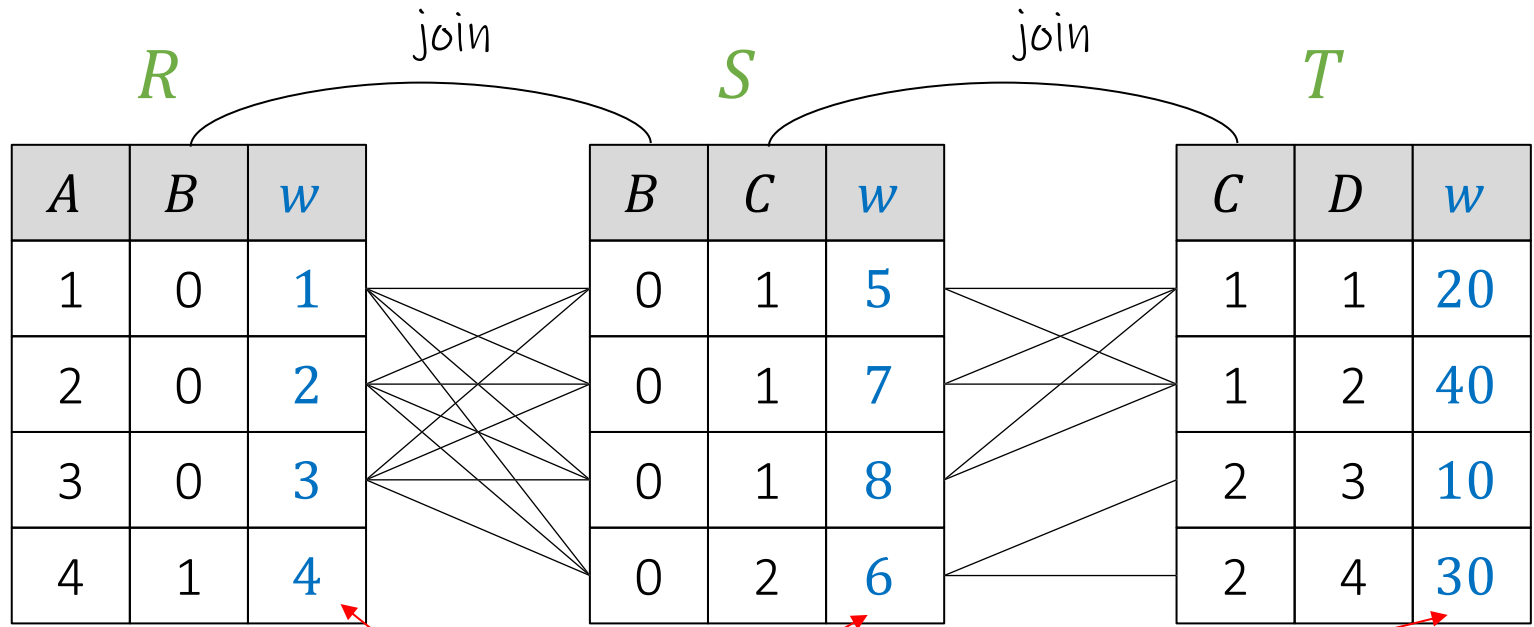
## Result

A	B	C	D	weight
---	---	---	---	--------

What do we get now?

?

# Sorting & Top-k evaluation with SQL

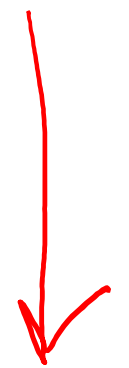


```
select  A, R.B, S.C, D,  
        R.w + S.w + T.w as weight  
from    R, S, T  
where   R.B=S.B and S.C=T.C  
order by weight ASC
```

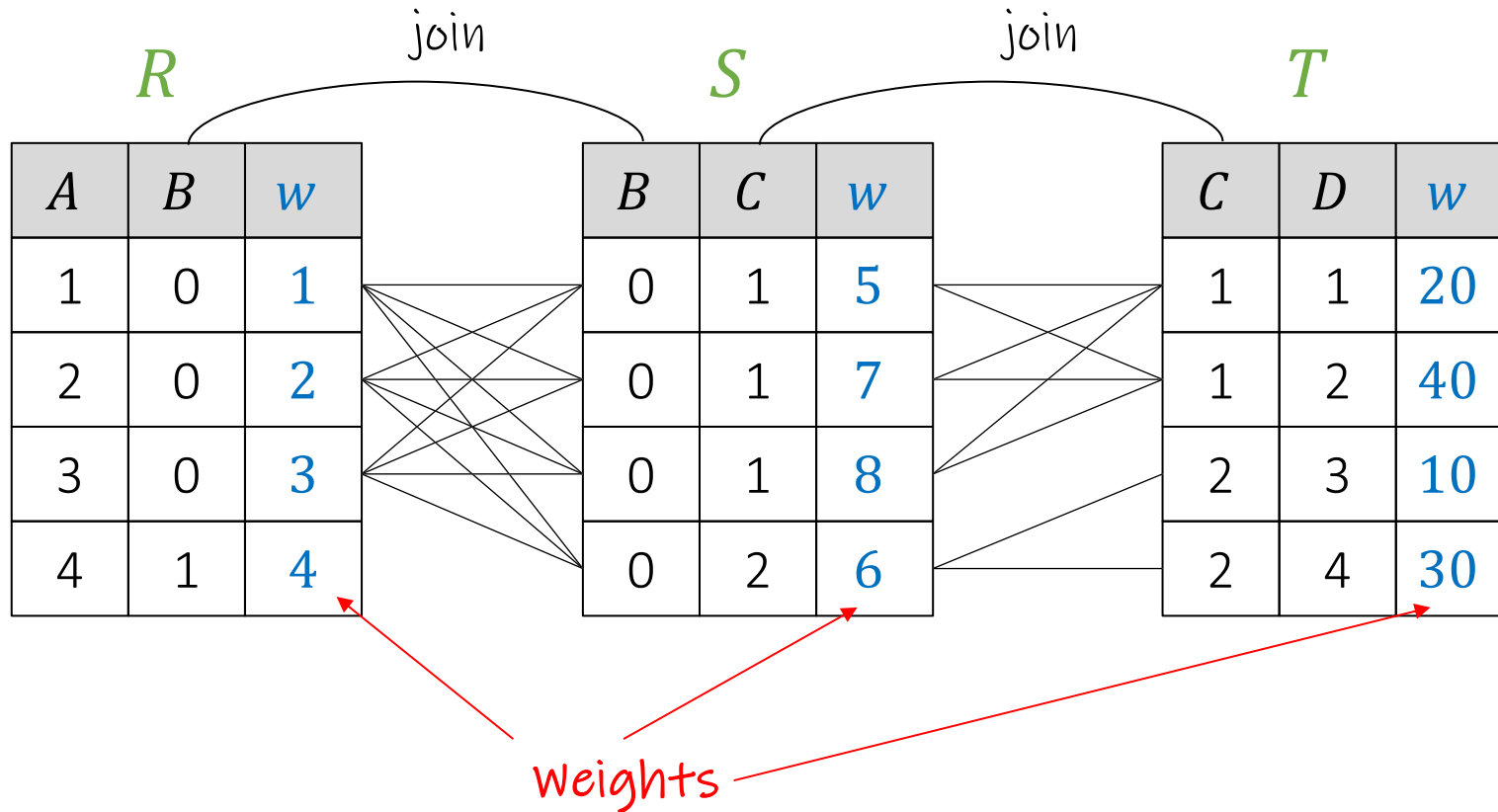
Return all 24 results in order of sum of weights

## Result

A	B	C	D	weight
1	0	2	3	17
2	0	2	3	18
3	0	2	3	19
1	0	1	1	26
2	0	1	1	27
3	0	1	1	28
1	0	1	1	28
...	...	...	...	...



# Sorting & Top-k evaluation with SQL



```
select A, R.B, S.C, D,  
       R.w + S.w + T.w as weight  
from   R, S, T  
where  R.B=S.B and S.C=T.C  
order by weight ASC  
limit  6
```

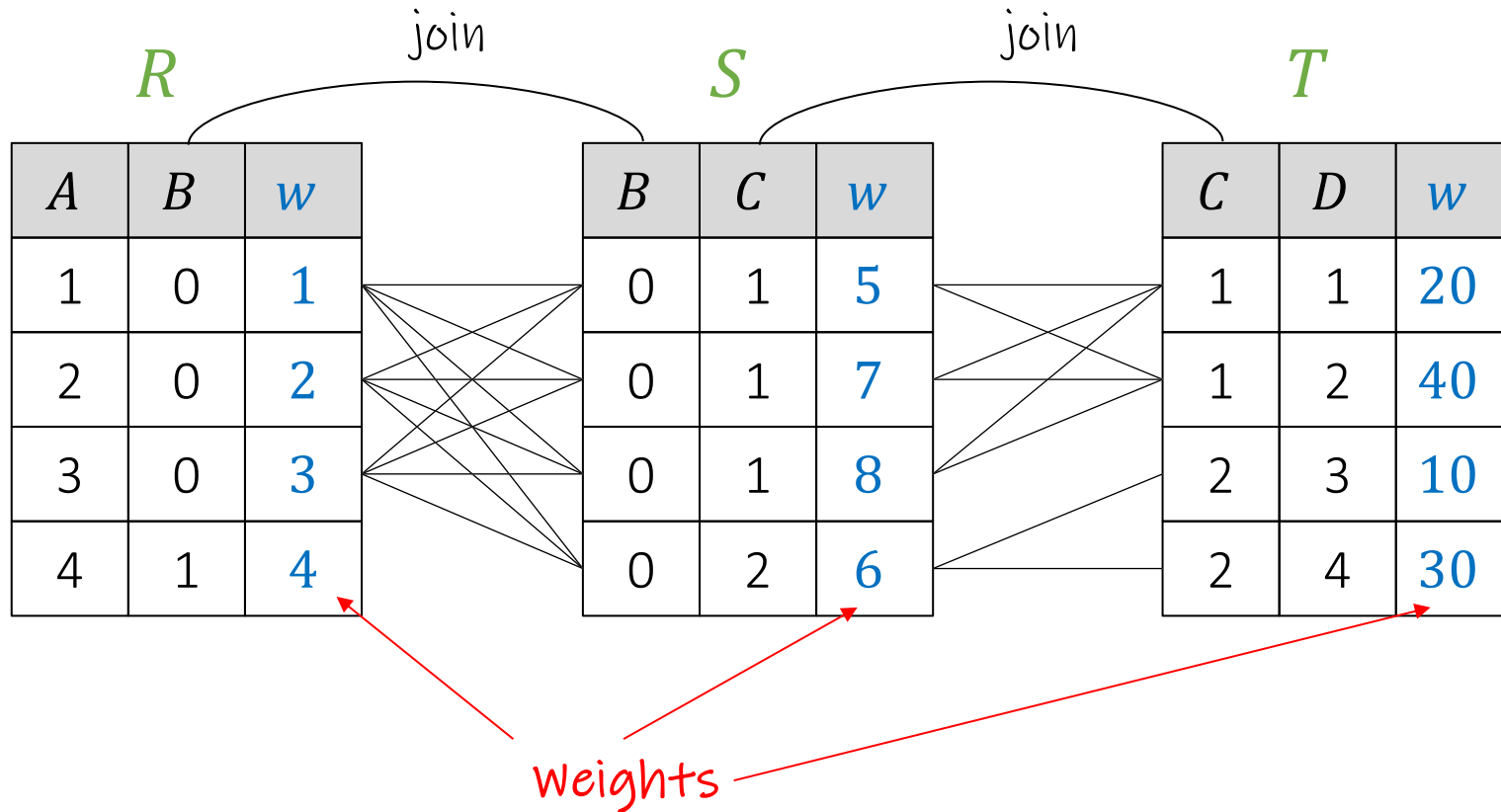
What do we get now?



## Result

A	B	C	D	weight
1	0	2	3	17
2	0	2	3	18
3	0	2	3	19
1	0	1	1	26
2	0	1	1	27
3	0	1	1	28
1	0	1	1	28
...	...	...	...	...

# Sorting & Top-k evaluation with SQL



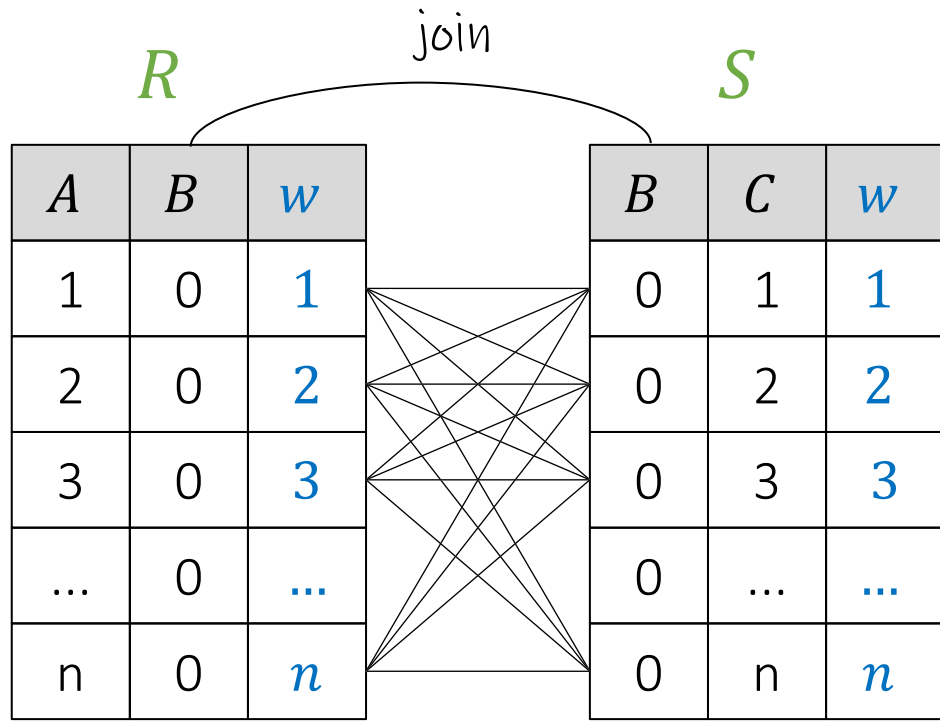
```
select  A, R.B, S.C, D,  
        R.w + S.w + T.w as weight  
from    R, S, T  
where   R.B=S.B and S.C=T.C  
order  by weight ASC  
limit   6
```

What do we get now?

## Result

A	B	C	D	weight
1	0	2	3	17
2	0	2	3	18
3	0	2	3	19
1	0	1	1	26
2	0	1	1	27
3	0	1	1	28
1	0	1	1	28
...	...	...	...	...

# Top- $k$ is evaluated inefficiently by modern DBMS's

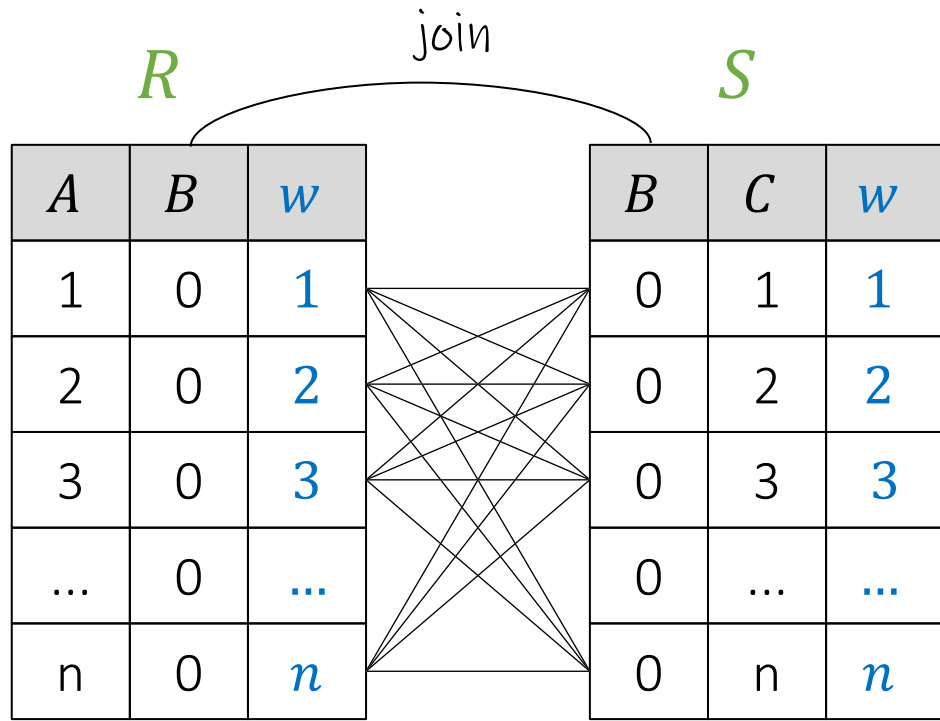


```
select  A, R.B, S.C,  
        R.w + S.w as weight  
from    R, S  
where   R.B=S.B  
order  by weight ASC  
limit   1
```

What will this query return?



# Top-k is evaluated inefficiently by modern DBMS's



```
select A, R.B, S.C,  
       R.w + S.w as weight  
from   R, S  
where  R.B=S.B  
order by weight ASC  
limit  1
```

Handwritten annotations in red: 2, 7, 2, 4, 5

## Result

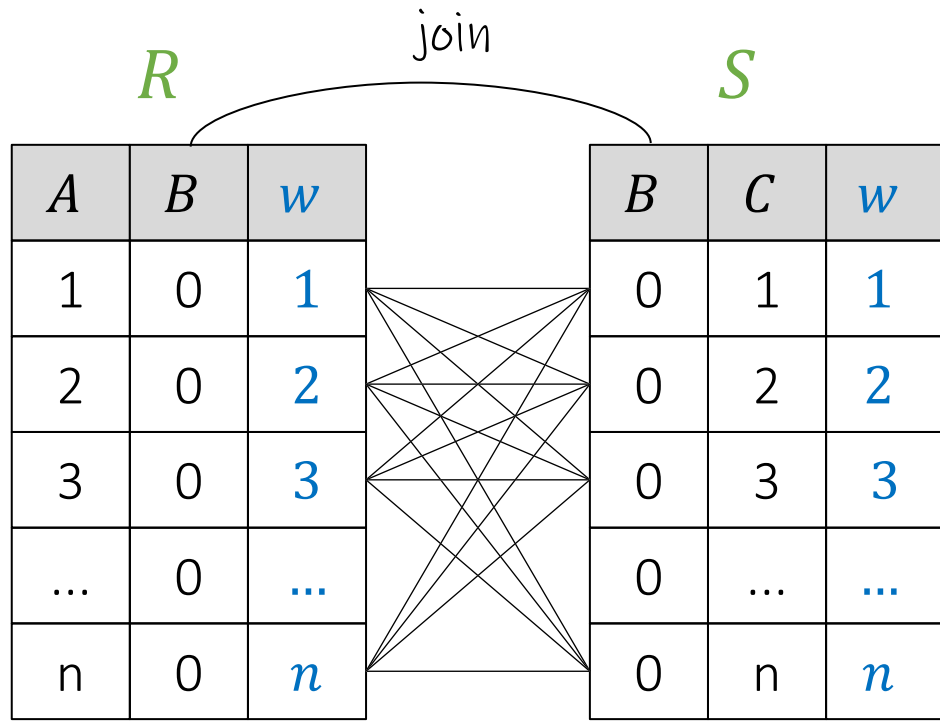
A	B	C	weight
1	0	1	2
1	0	2	3
2	0	1	3
3	0	1	4
2	0	2	4
1	0	3	4
4	0	1	5
...	...	...	...
n	0	n	n*n

Can you see any possible problem of this query as n gets bigger?





# Top- $k$ is evaluated inefficiently by modern DBMS's



```
select  A, R.B, S.C,  
        R.w + S.w as weight  
from    R, S  
where   R.B=S.B  
order  by weight ASC  
limit   1
```

Result

A	B	C	weight
1	0	1	2
1	0	2	3
2	0	1	3
3	0	1	4
2	0	2	4
1	0	3	4
4	0	1	5
...	...	...	...
$n$	0	$n$	$n*n$

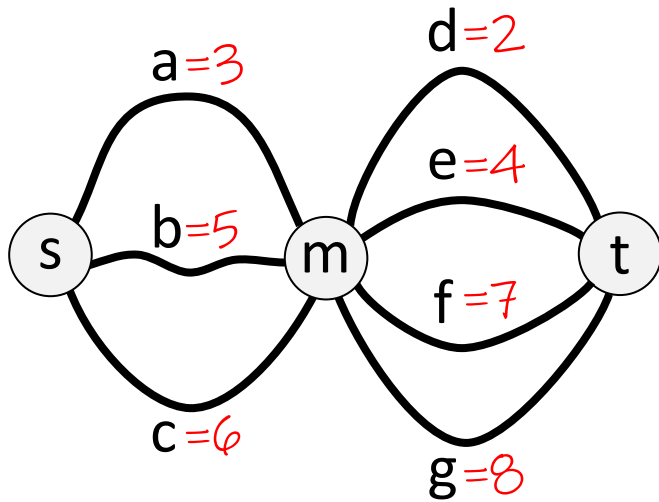
$n^2$  total results. But we are only interested in the top-1.

Problem: Database first calculates all  $n^2$  results before sorting.

Question: is there any way to push the sorting behind the join?



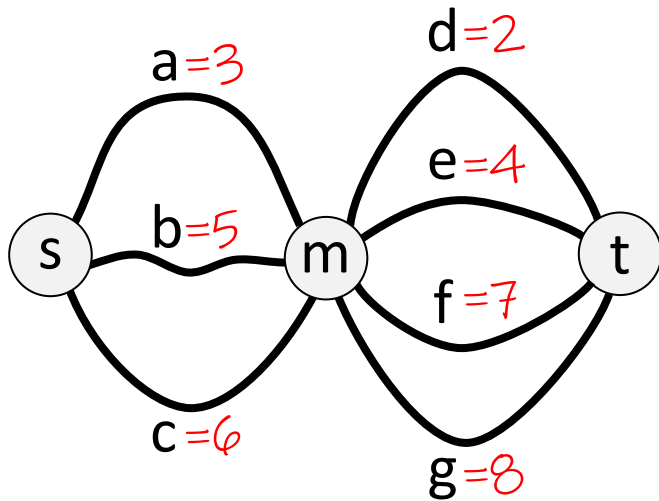
# Digression: Distributivity = efficient factorization



What is the shortest path from s to t?



# Digression: Distributivity = efficient factorization



$$\min [a + d, a + e, a + f, a + g, \dots, c + g]$$

$$\min [3+2, 3+4, 3+7, 3+8, \dots, 6+8]$$

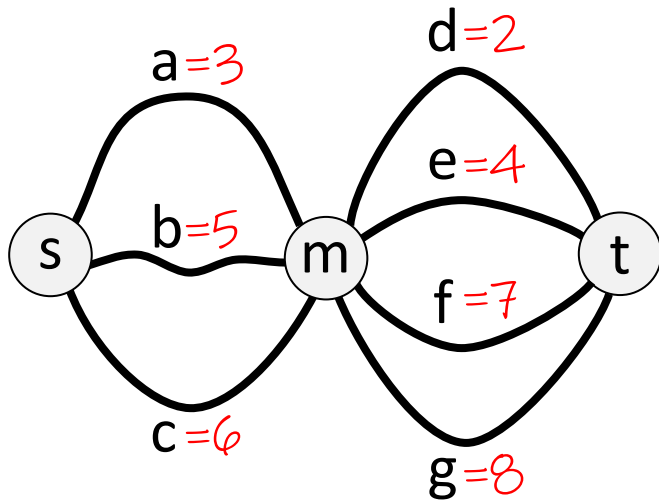
?

What is the shortest path from s to t?

Answer:  $5 = 3 + 2$

# Digression: Distributivity = efficient factorization

Principle of optimality from Dynamic Programming:  
*irrespective of the initial state and decision, an optimal solution continues optimally from the resulting state*



What is the shortest path from s to t?

Answer:  $5 = 3 + 2$

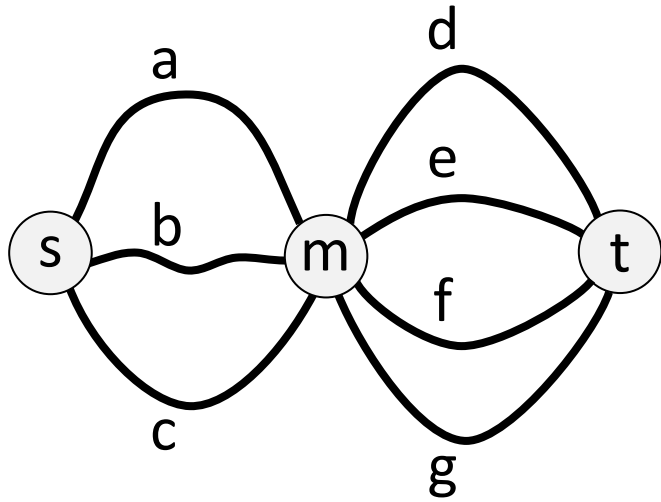
$$\min [a + d, a + e, a + f, a + g, \dots, c + g]$$
$$\min [3+2, 3+4, 3+7, 3+8, \dots, 6+8]$$

$$= \min [a, b, c] + \min [d, e, f, g]$$
$$\min [3, 5, 6] + \min [2, 4, 7, 8]$$

$$\min [x, y] + z = \min [(x+z), (y+z)]$$

(+ distributes over min)

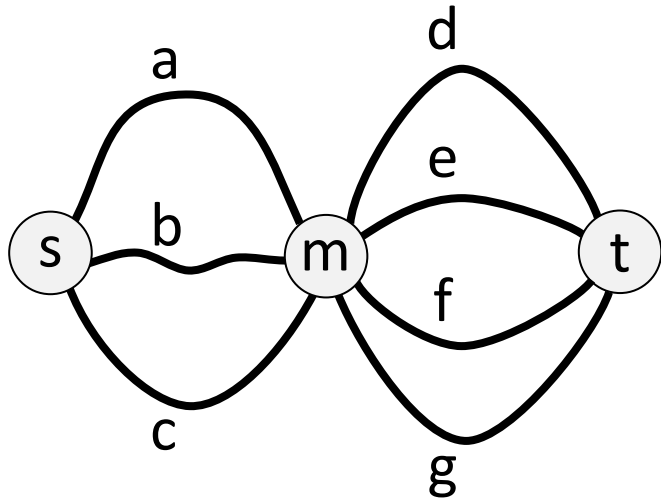
# Digression: Distributivity = efficient factorization



How many paths are there from s to t?



# Digression: Distributivity = efficient factorization

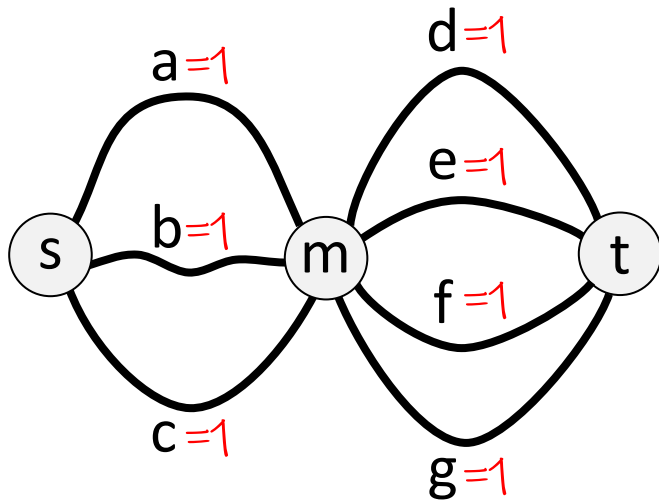


How many paths are there from s to t?

Answer:  $12 = 3 \cdot 4$

# Digression: Distributivity = efficient factorization

The more general algebraic structure behind these two examples are "semirings" (much more on those later in class)



How many paths are there from s to t?

Answer:  $12 = 3 \cdot 4$

$$\text{count}[a \cdot d, a \cdot e, a \cdot f, a \cdot g, \dots, c \cdot g]$$

$$\text{count}[\underbrace{1 \cdot 1, 1 \cdot 1, 1 \cdot 1, 1 \cdot 1, \dots, 1 \cdot 1}_{12}]$$

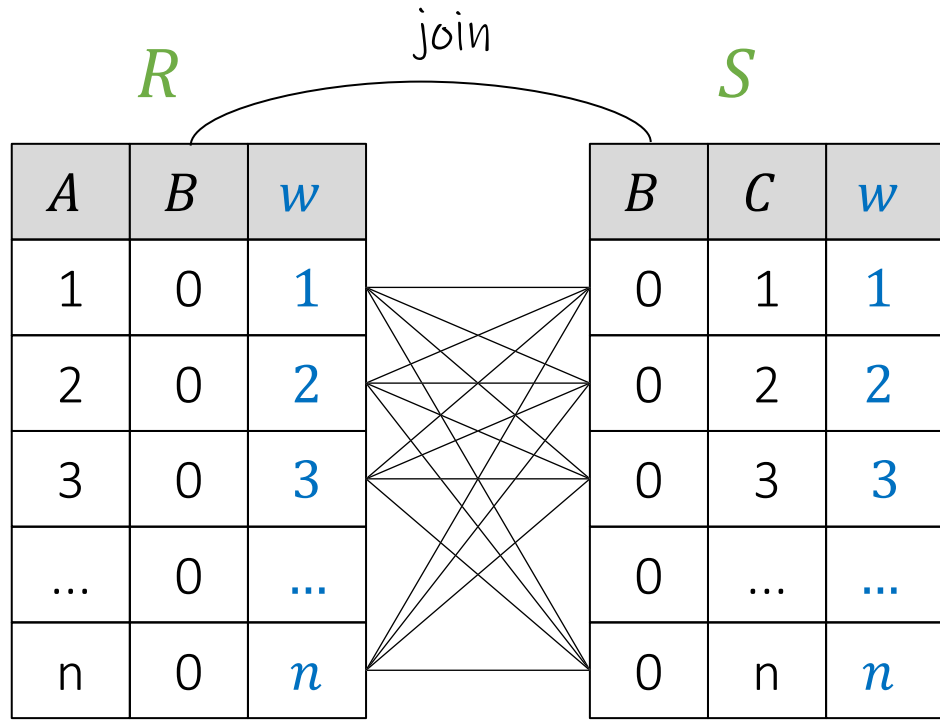
$$= \text{count}[a, b, c] \cdot \text{count}[d, e, f, g]$$

$$\text{count}[1, 1, 1] \cdot \text{count}[1, 1, 1, 1]$$

$$+[x, y] \cdot z = +[x \cdot z, y \cdot z]$$

( $\cdot$  distributes over  $+$ )

# Top- $k$ is evaluated inefficiently by modern DBMS's



-----  
 -- Query 2  
 -----

-----  
 -- Query 1  
 -----

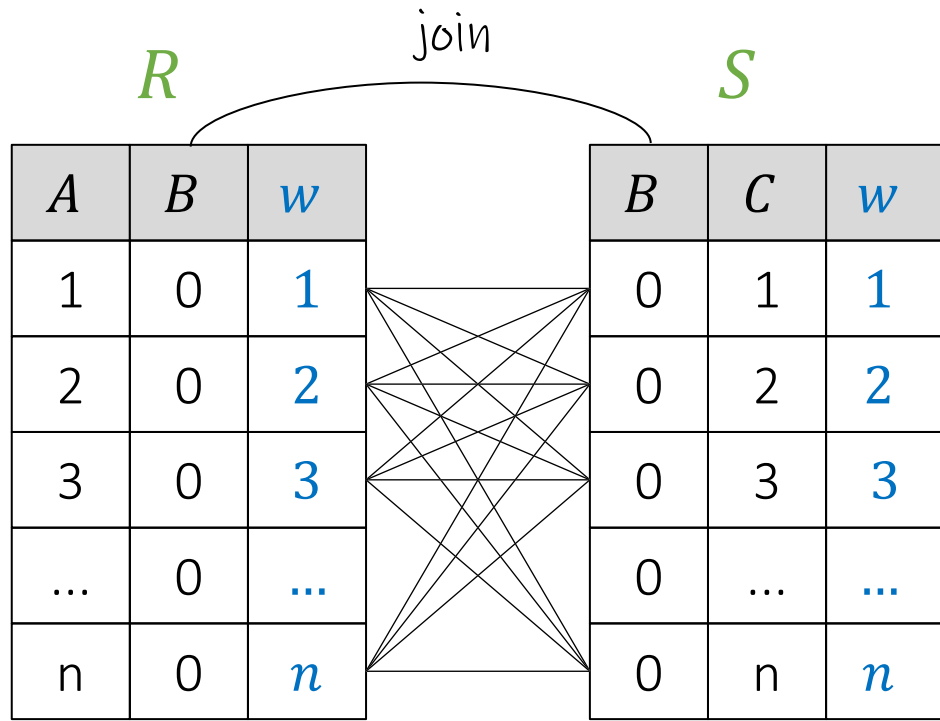


```
SELECT  A, R.B, S.C,
        R.W + S.W as weight
FROM    R, S
WHERE   R.B=S.B
ORDER BY weight ASC
LIMIT  1;
```

$n=1000$ :             $t_{Q1} = 0.88 \text{ sec}$              $t_{Q2} = 2 \text{ msec}$   
 $n=5000$ :             $t_{Q1} = 18.6 \text{ sec}$              $t_{Q2} = 8 \text{ msec}$



# Top-k is evaluated inefficiently by modern DBMS's



Maximal intermediate result size is  $O(n)$  😊  
 What is this algorithm called?



```

-----
-- Query 1
-----
SELECT  A, R.B, S.C,
        R.W + S.W as weight
FROM    R, S
WHERE   R.B=S.B
ORDER BY weight ASC
LIMIT  1;
    
```

```

-----
-- Query 2
-----
SELECT R.A, X.B, S.C, X.W as weight
FROM R, S,
     (SELECT T1.B, W1, W2, W1+W2 W
      FROM
        (SELECT B, MIN(W) W1
         FROM R
         GROUP BY B) T1,
        (SELECT B, MIN(W) W2
         FROM S
         GROUP BY B) T2
      WHERE T1.B = T2.B
      ORDER BY W ASC
      LIMIT 1) X
WHERE X.B = R.B
AND X.W1 = R.W
AND X.B = S.B
AND X.W2 = S.W
LIMIT 1;
    
```

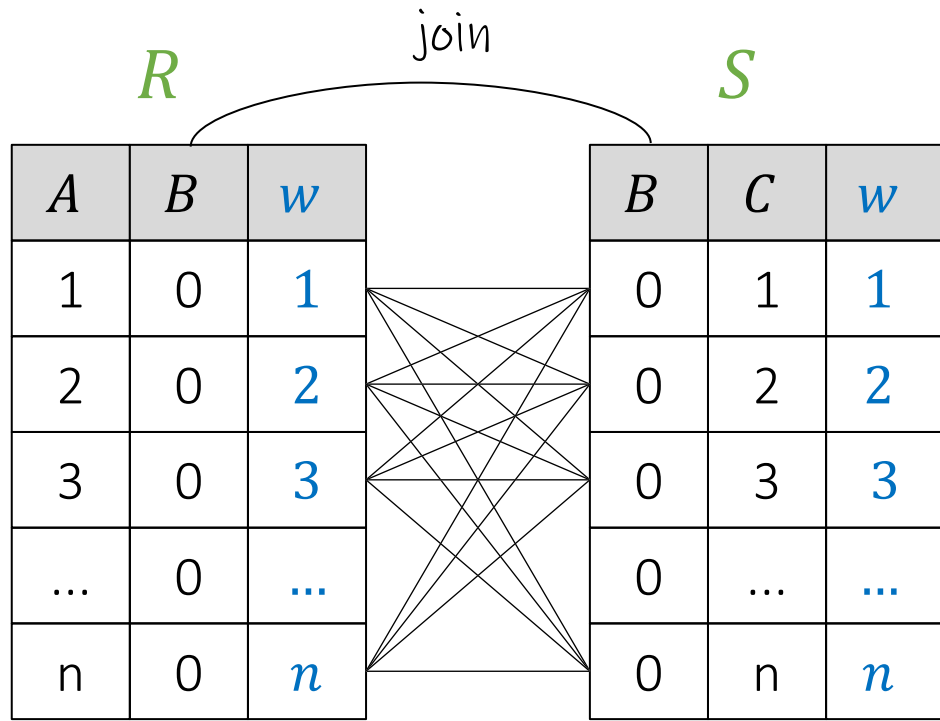
n=1000: t<sub>Q1</sub> = 0.88 sec

t<sub>Q2</sub> = 2 msec

n=5000: t<sub>Q1</sub> = 18.6 sec

t<sub>Q2</sub> = 8 msec

# Top-k is evaluated inefficiently by modern DBMS's



Maximal intermediate result size is  $O(n)$  😊  
 What is this algorithm called?

Dynamic programming

```

-- Query 1
SELECT  A, R.B, S.C,
        R.W + S.W as weight
FROM    R, S
WHERE   R.B=S.B
ORDER BY weight ASC
LIMIT  1;
    
```

n=1000:  $t_{Q1} = 0.88 \text{ sec}$   
 n=5000:  $t_{Q1} = 18.6 \text{ sec}$

$O(n^2)$

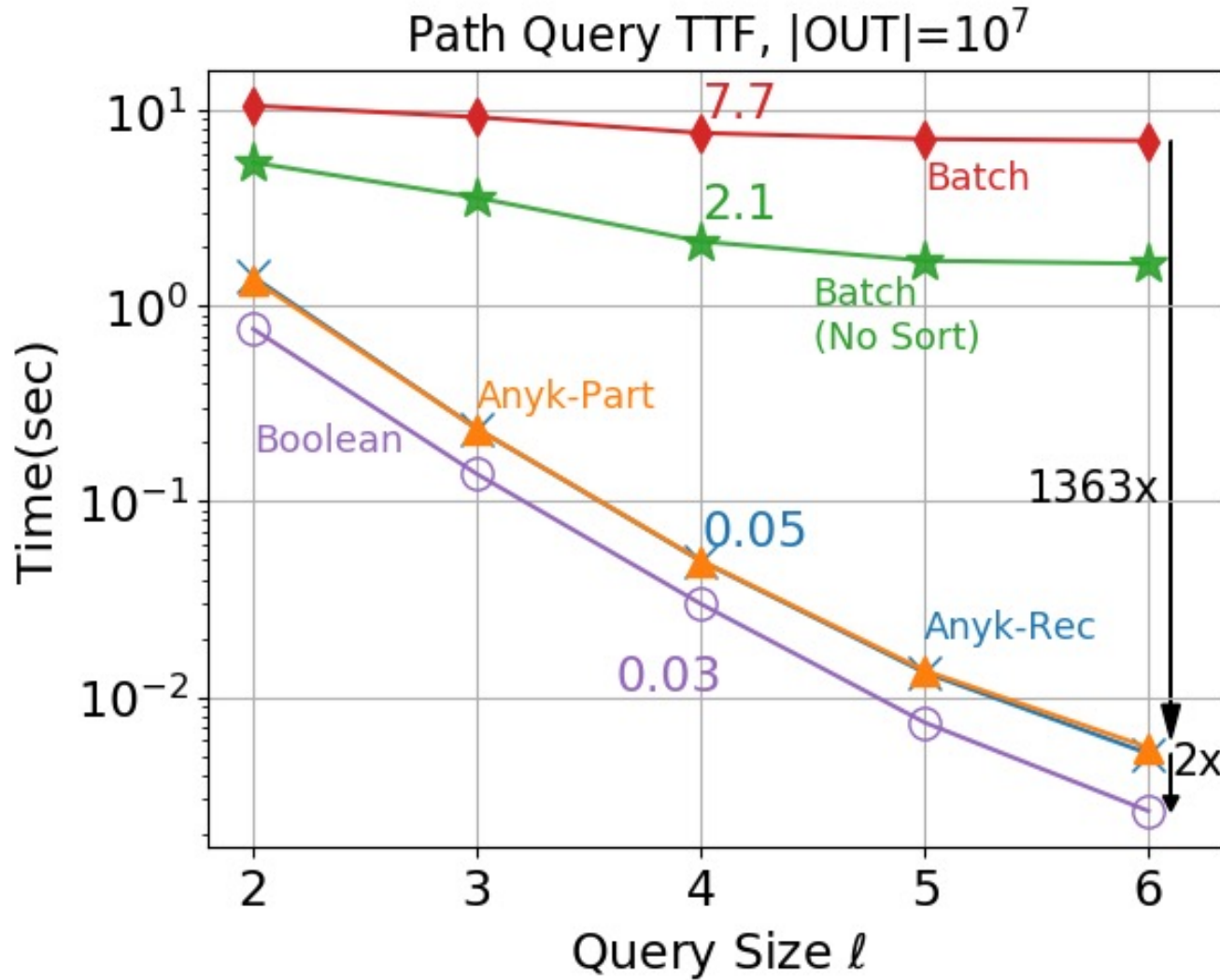
```

-- Query 2
SELECT R.A, X.B, S.C, X.W as weight
FROM R, S,
     (SELECT T1.B, W1, W2, W1+W2 W
      FROM
         (SELECT B, MIN(W) W1
          FROM R
          GROUP BY B) T1,
         (SELECT B, MIN(W) W2
          FROM S
          GROUP BY B) T2
      WHERE T1.B = T2.B
      ORDER BY W ASC
      LIMIT 1) X
WHERE X.B = R.B
AND X.W1 = R.W
AND X.B = S.B
AND X.W2 = S.W
LIMIT 1;
    
```

$t_{Q2} = 2 \text{ msec}$   
 $t_{Q2} = 8 \text{ msec}$

$O(n)$

# Any- $k$ : Faster and more versatile than Top- $k$



Path query with constant size output and increasing query size

<https://northeastern-datalab.github.io/anyk/>

<https://northeastern-datalab.github.io/topk-join-tutorial/>

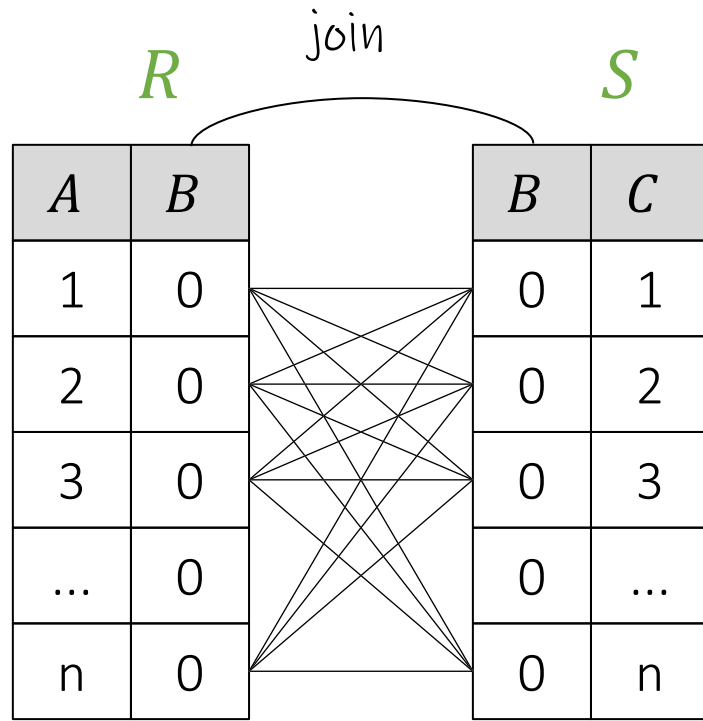
<https://northeastern-datalab.github.io/responsive-dbms-tutorial/>

[https://www.youtube.com/watch?v=KpUQayBuaQI&list=PL\\_72ERGF6DR7kvGNwwjWlbpScKtGjt9R&index=2](https://www.youtube.com/watch?v=KpUQayBuaQI&list=PL_72ERGF6DR7kvGNwwjWlbpScKtGjt9R&index=2)

Tziavelis, Ajwani, Gatterbauer, Riedewald, Yang. Optimal Algorithms for Ranked Enumeration of Answers to Full Conjunctive Queries. PVLDB 2020. <https://doi.org/10.14778/3397230.3397250>

Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>

# Even Grouping Aggregates can be improved



```
-- Query 1
SELECT count(*) CT
INTO record1
FROM R, S
WHERE R.B=S.B;
```

n=1000:  $t_{Q1} = 0.374 \text{ sec}$

n=5000:  $t_{Q1} = 10.021 \text{ sec}$

```
-- Query 2
```

```
SELECT SUM(C) as CT
INTO record2
FROM
  (SELECT T1.B, C1, C2, C1*C2 C
   FROM
     (SELECT B, COUNT(*) C1
      FROM R
      GROUP BY B) T1,
     (SELECT B, COUNT(*) C2
      FROM S
      GROUP BY B) T2
   WHERE T1.B = T2.B) X;
```

$t_{Q2} = 3 \text{ msec}$

$t_{Q2} = 5 \text{ msec}$

# Outline: T1-U1: SQL

- SQL

- Schema, keys, referential integrity
- Joins
- Aggregates and grouping
- Nested queries (Subqueries)
- Union and Theta Joins
- Nulls & Outer joins
- Top-k
- [Recursion: moved to T1-U4: Datalog]