

DATABASE FUNDAMENTALS (RECALL/CRASH COURSE)

Database functionalities

Database Management Systems

- Functionality provided
 - What kind of data can I put in? **Relations/documents/pairs...**
 - How can I get data out of it? **query languages/API**
 - How does it handle concurrent access?
ACID (or less)
 - How long does a given operation take?
Query execution, optimization
- Implementation (internals)
 - How does it cope with scale?
 - for reads? **Smart storage and indexing structures**
 - for writes? **Concurrency control**

Relational Database Management Systems

- Functionality provided
 - What kind of data can I put in? **Relations**
 - How can I get data out of it? **SQL query language**
 - How does it handle concurrent access? **ACID (or less)**
 - How long does a given operation take? **Query optimization**
- Implementation (internals)
 - How does it cope with scale?
 - for reads? **Smart storage and indexing structures**
 - for writes? **Concurrency control**

Fundamental database features

1. Data storage

- Protection against unauthorized access, data loss

2. Ability to at least **add** to and **remove** data to the database

- Also: **updates; active behavior** upon update (triggers)

3. Support for **accessing** the data

- Declarative query languages: say what data you need, not how to find it

DATABASE FUNDAMENTALS (RECALL/CRASH COURSE)

Query processing

How are queries processed?



How are queries processed?

Driver	
name	ID
Julie	1
Damien	2


Car	
driver	license
1	'123AB'
2	'171KZ'

1. Load 

2. SQL 

```
select driver.name
from driver, car
where
driver.ID=car.driver
and
car.license='123AB'
```

Database

 3. Results

name
Julie

How are queries processed?

Driver	
name	ID
Julie	1
Damien	2

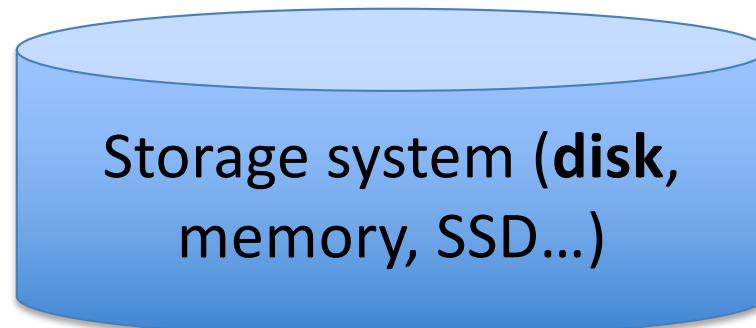
Car	
driver	license
1	'123AB'
2	'171KZ'


1. Load 

2. SQL 

```
select driver.name
from driver, car
where
driver.ID=car.driver
and
car.license='123AB'
```

Database



 3. Results

name
Julie

How are queries processed?

1. Load 

2. SQL 

```
select driver.name  
from driver, car  
where  
driver.ID=car.driver  
and  
car.license='123AB'
```

Database

Driver		Car	
name	ID	driver	license
Julie	1	1	'123AB'
Damien	2	2	'171KZ'

 3. Results

name

Julie

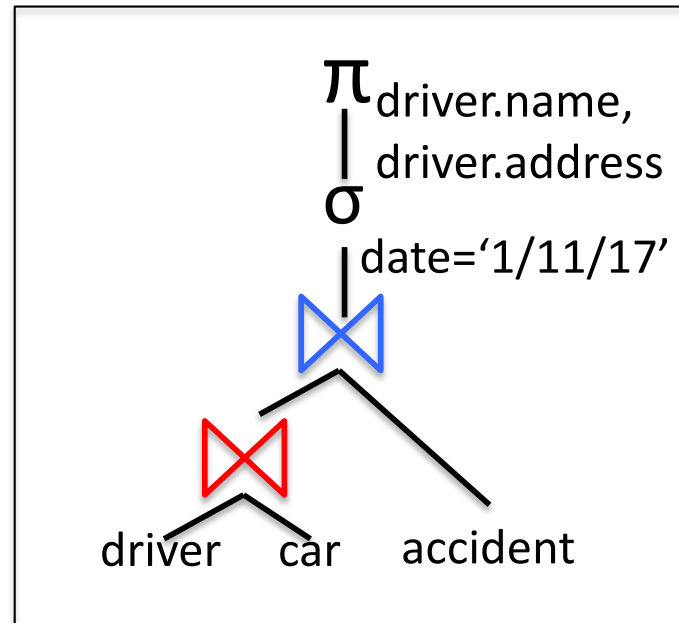
How are queries processed?

SQL

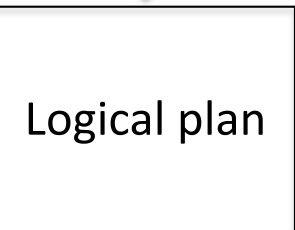


```
select driver.name,  
driver.address  
from driver, car,  
accident  
where  
driver.ID=car.driver  
and  
car.license=accident  
.carLicense and  
accident.date='1/11  
/17'
```

select... from driver, car, accident where...



Query language



.....

Driver		Accident		Car	
name	ID	driver	license		
Julie	1	1	'123AB'		
Damien	2	2	'171KZ'		



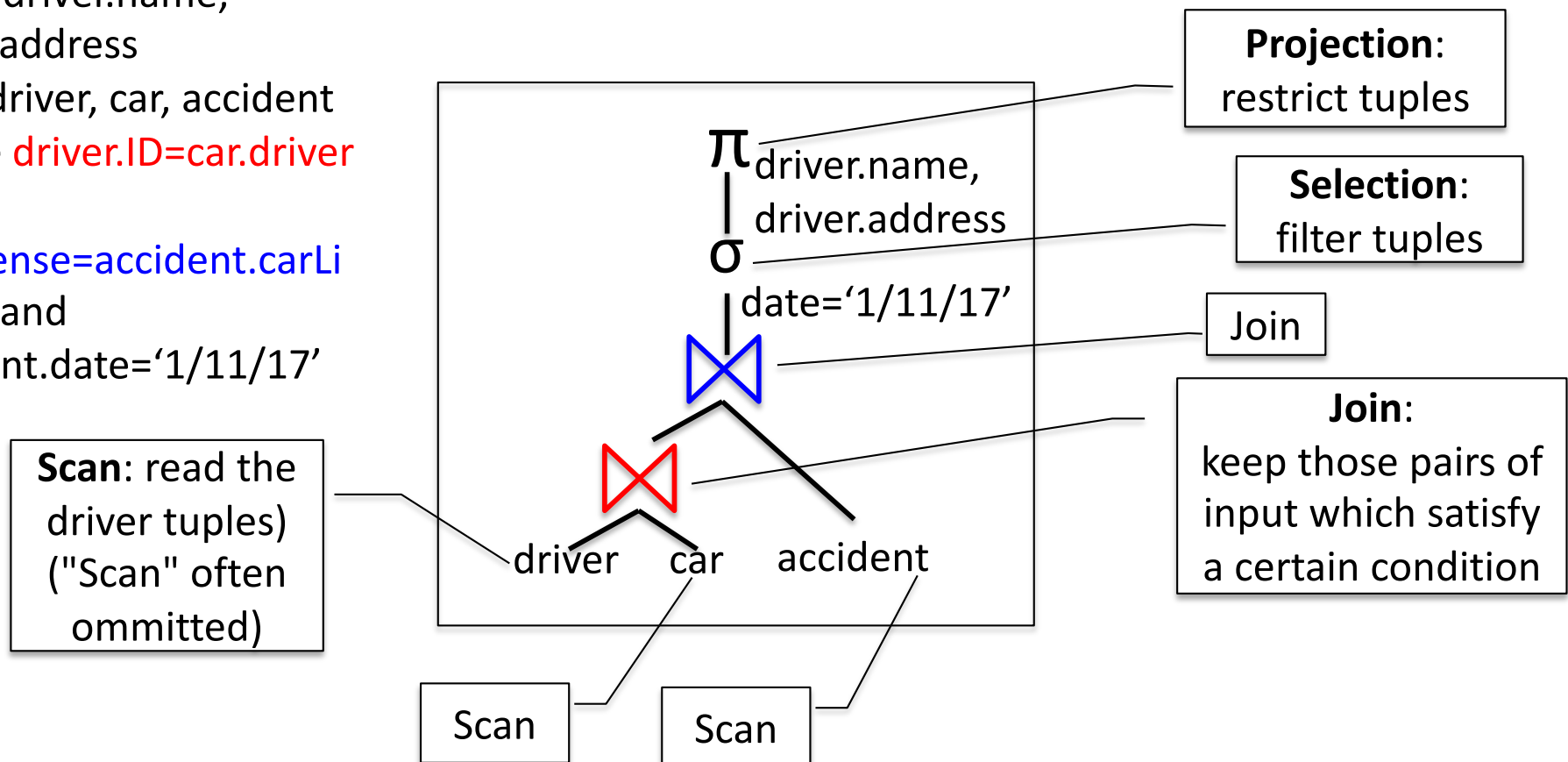
Results

Logical query plans

- Trees made of logical operators, each of which specializes in a certain task

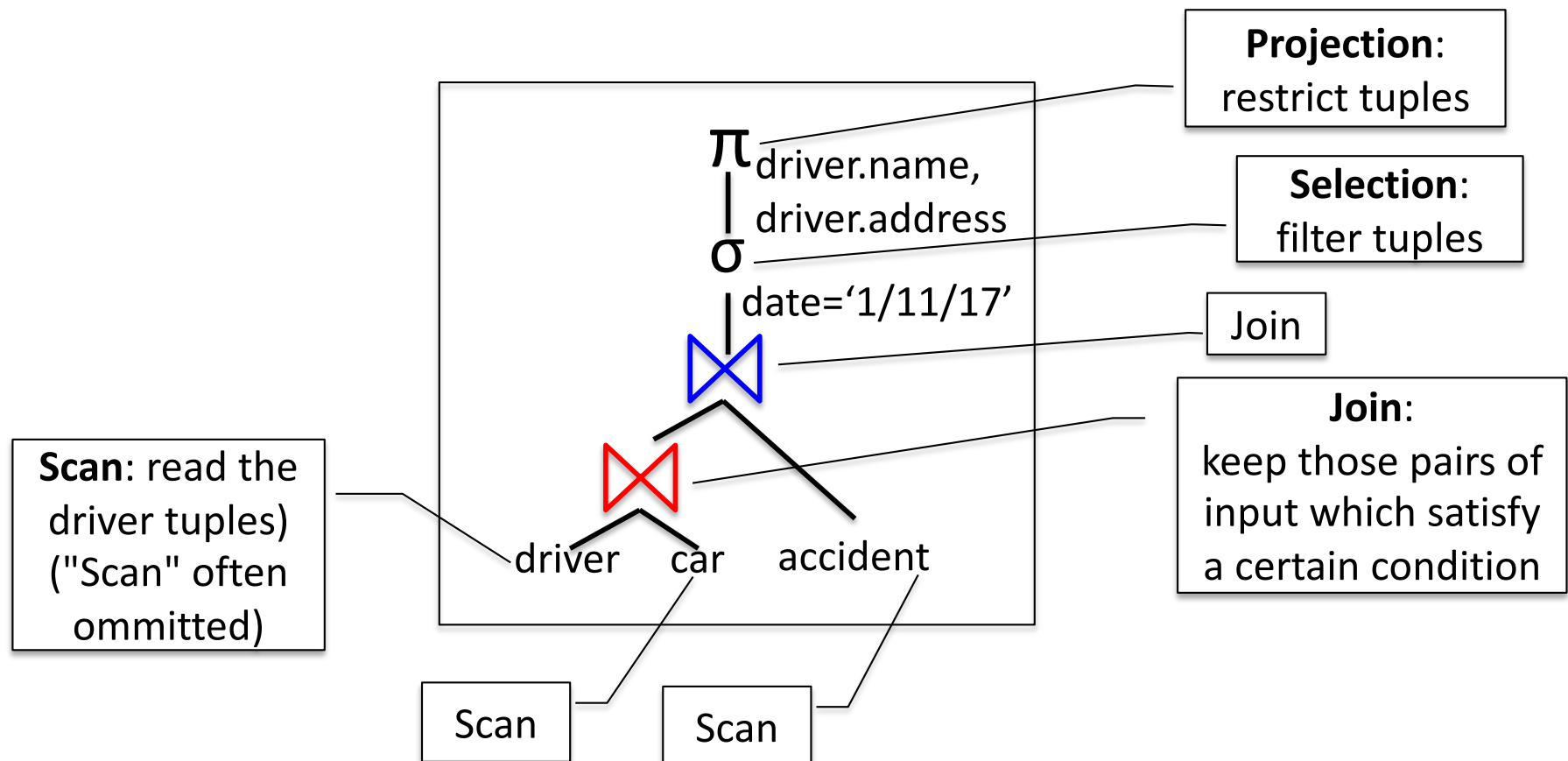
SQL:

```
select driver.name,  
driver.address  
from driver, car, accident  
where driver.ID=car.driver  
and  
car.license=accident.carLi  
cense and  
accident.date='1/11/17'
```



Logical query plans

- Trees made of logical operators, each of which specializes in a certain task
- Logical operators: they are defined by their result, not by an algorithm
- Physical operators (see next) implement actual algorithms

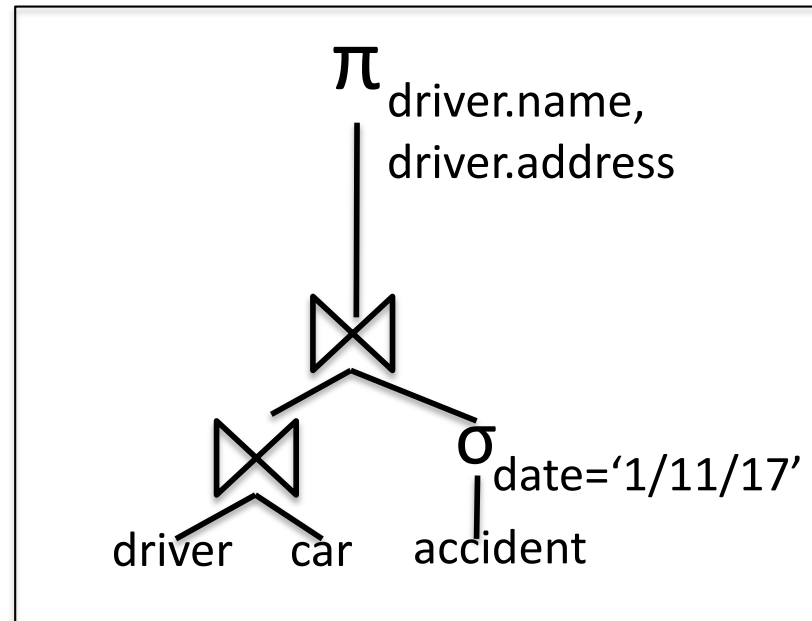


How are queries processed?

SQL

```
select driver.name,  
driver.address  
from driver, car,  
accident  
where  
driver.ID=car.driver  
and  
car.license=accident  
.carLicense and  
accident.date='1/11  
/13'
```

select... from driver, car, accident where...



Query language

Logical plan 1

Logical plan 2

.....

Driver		Accident		Car
name	ID	driver	license	
Julie	1	1	'123AB'	
Damien	2	2	'171KZ'	

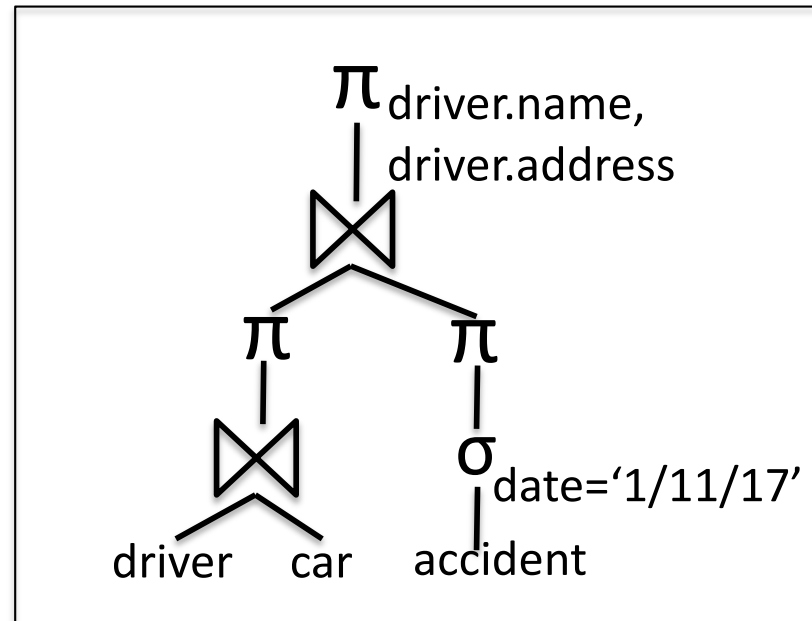
Results

How are queries processed?

SQL

```
select driver.name,  
driver.address  
from driver, car,  
accident  
where  
driver.ID=car.driver  
and  
car.license=accident  
.carLicense and  
accident.date='1/11  
/17'
```

select... from driver, car, accident where...



Driver		Accident		Car
name	ID	driver	license	
Julie	1	1	'123AB'	
Damien	2	2	'171KZ'	

Query language

Logical plan 1

Logical plan 2

Logical plan 3

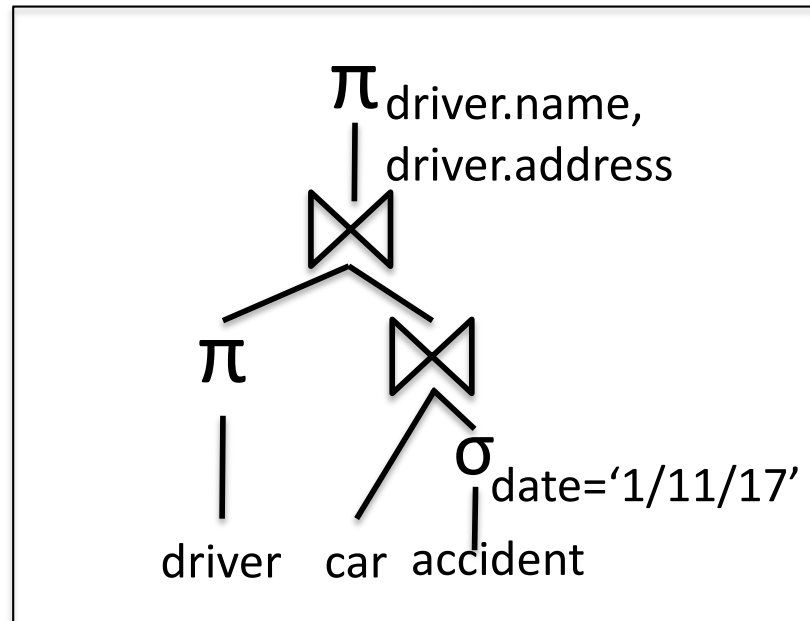
Results

How are queries processed?

SQL

```
select driver.name,  
driver.address  
from driver, car,  
accident  
where  
driver.ID=car.driver  
and  
car.license=accident  
.carLicense and  
accident.date='1/11  
/17'
```

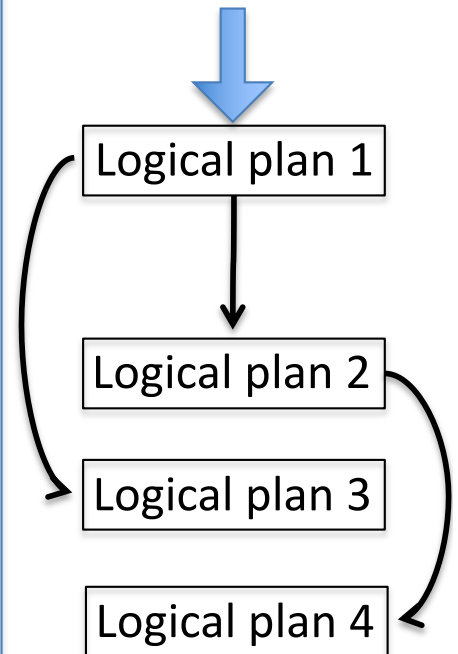
select... from driver, car, accident where...



.....

Driver		Accident		Car	
name	ID	driver	license		
Julie	1	1	'123AB'		
Damien	2	2	'171KZ'		

Query language



Results

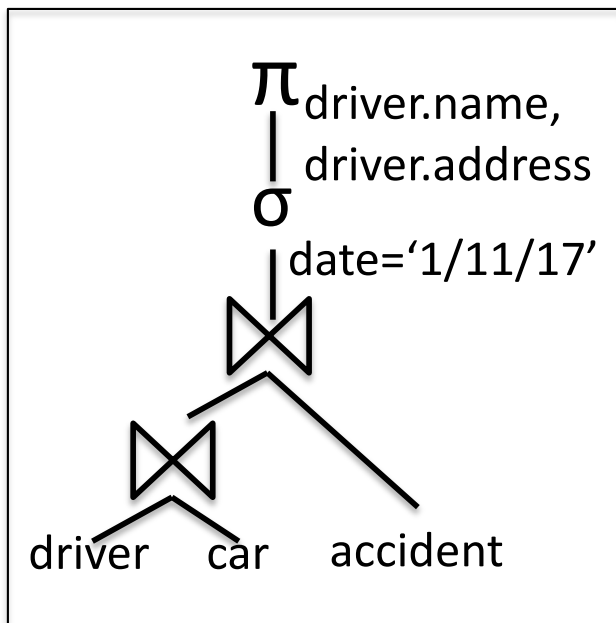
Logical query optimization

- Enumerates logical plans
- All logical plans compute the query result
 - They are **equivalent**
- Some are (much) more **efficient** than others
- **Logical optimization**: moving from a plan to a more efficient one
 - Pushing selections
 - Pushing projections
 - Join reordering: most important source of optimizations

Logical query optimization example

1.000.000 cars, 1.000.000 drivers, 1.000 accidents, 2 cars per accident, 10 accidents on 1/11/17

« Name and address of drivers in accidents on 1/11/2017? »



Cost of an operator: depends on the number of tuples (or tuple pairs) which it must process
e.g. $c_{\text{disk}} \times \text{number of tuples read from disk}$
e.g. $c_{\text{cpu}} \times \text{number of tuples compared}$

Cardinality of an operator's output: how many tuples result from this operator

The cardinality of one operator's output determines the cost of its parent operator

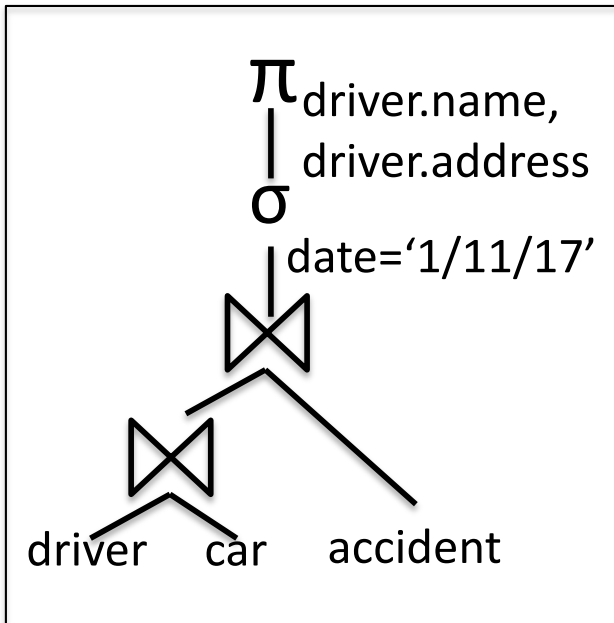
Plan **cost** = the sum of the costs of all operators in a plan

Logical query optimization example

1.000.000 cars, 1.000.000 drivers, 1.000 accidents, 2 cars per accident, 10 accidents on 1/11/17

« Name and address of drivers in accidents on 1/11/2017? »

Pessi-
mistic
(worst-
case)
estim.



cs, cj, cf constant

Scan **costs**: $cs \times (10^6 + 10^6 + 10^3)$

Scan **cardinality** estimations: $10^6, 10^6, 10^3$

Driver-car join **cost** estimation: $cj \times (10^6 \times 10^6 = 10^{12})$

Driver-car join **cardinality** estimation: 10^6

Driver-car-accident join **cost** estim.: $cj \times (10^6 \times 10^3 = 10^9)$

Driver-car-accident join **cardinality** estimation: 2×10^3

Selection **cost** estimation: $cf \times (2 \times 10^3)$

Selection **cardinality** estimation: 10

Projection (similar), negligible

Total **cost** estimation: $cs \times (2 \times 10^6 + 10^3) + cf \times 2 \times 10^3$

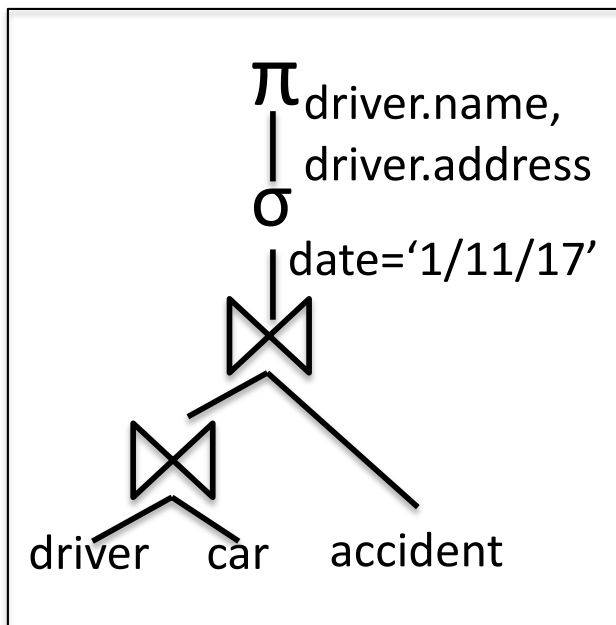
$+ cj \times (10^{12} + 2 \times 10^3) \sim cj \times 10^{12} \sim \mathbf{10^{12}}$

Logical query optimization example

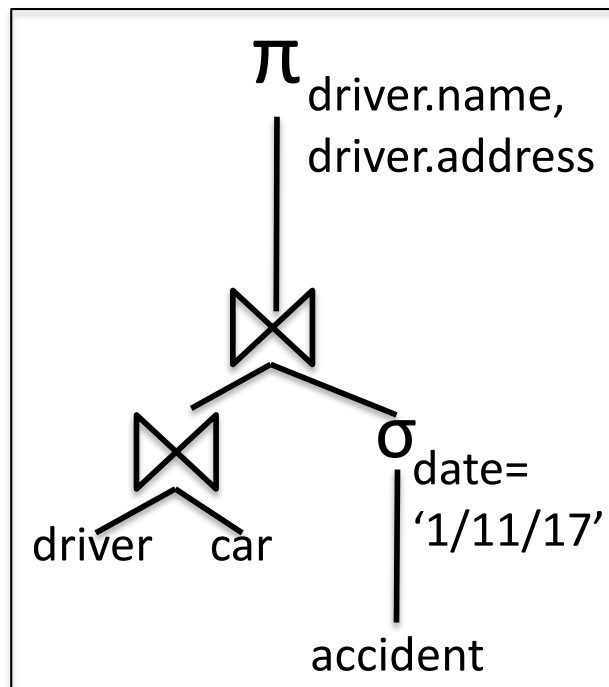
1.000.000 cars, 1.000.000 drivers, 1.000 accidents, 2 cars per accident, 10 accidents on 1/11/17

« Name and address of drivers in accidents on 1/11/2017? »

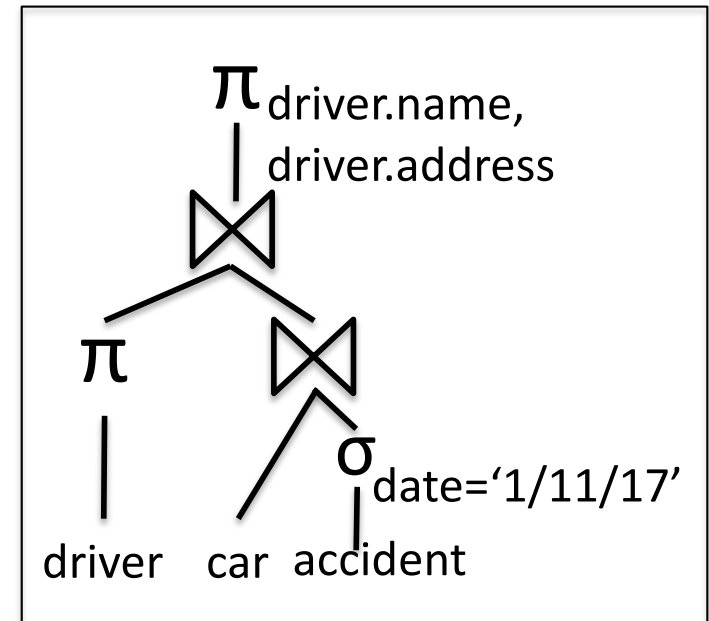
Three plans, same scan costs (neglected below); join costs dominant



$$10^9 + 10^{12} \sim 10^{12}$$



$$10^9 + 10^7 \sim 10^9$$



$$10^7 + 2 \cdot 10^7 \sim 3 \cdot 10^7$$

Logical query optimization example

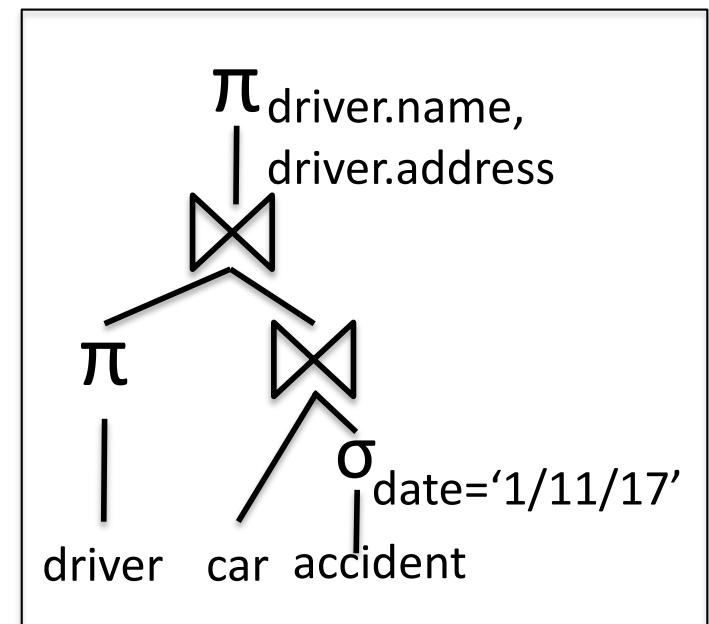
1.000.000 cars, 1.000.000 drivers, 1.000 accidents, 2 cars per accident, 10 accidents on 1/11/17

« Name and address of drivers in accidents on 1/11/2017? »

Three plans, same scan costs (neglected below); join costs dominant

The best plan reads only the accidents that have to be consulted

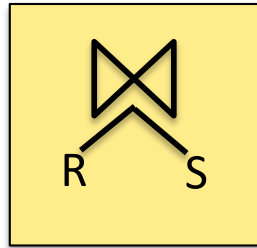
- **Selective data access**
- Typically supported by an **index**
 - Auxiliary data structure, built on top of the data collection
 - Allows to access directly objects satisfying a certain condition



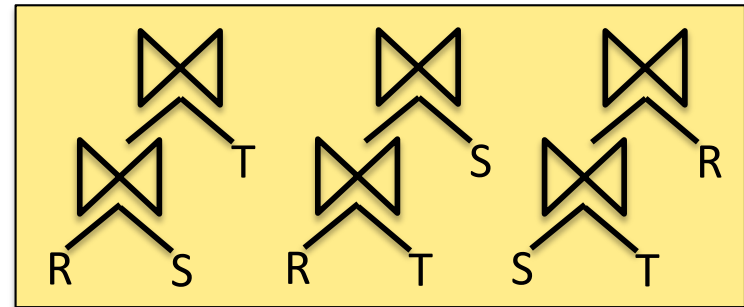
$$10^7 + 2 * 10^7 \sim 3 * 10^7$$

Join ordering is the main problem in logical query optimization

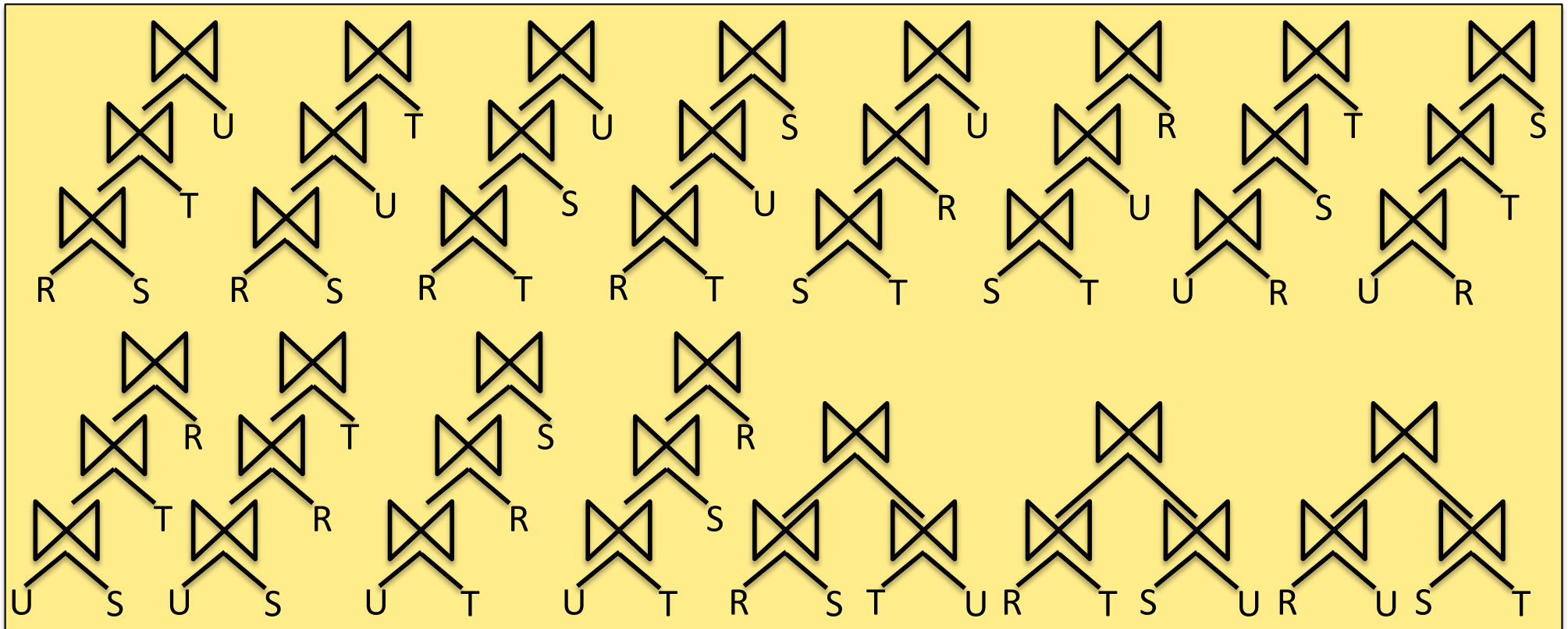
N=2:



N=3:



N=4:



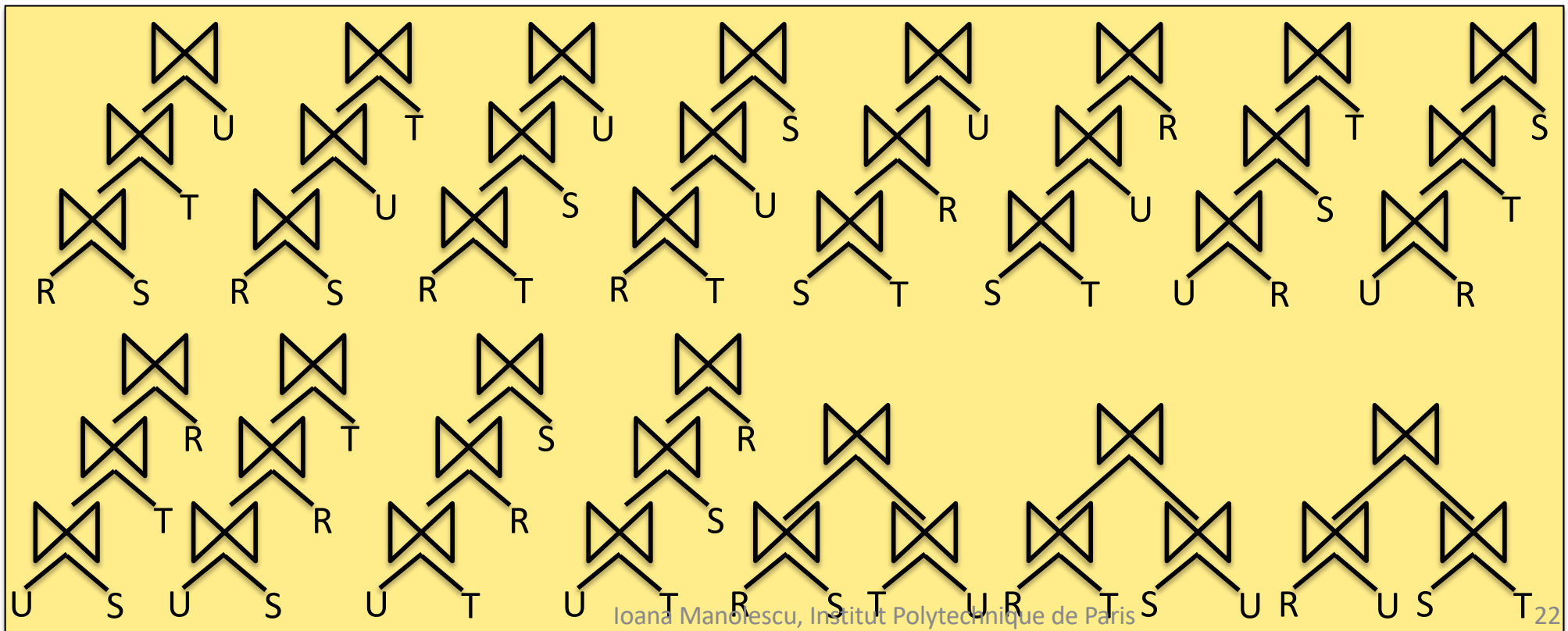
Join ordering is the main problem in logical query optimization

$$\text{Plans}(n+1) = (n+1) * \text{Plans}(n) + \frac{1}{2} * \sum_{i=1}^{(n/2)} \text{Plans}(i) * \text{Plans}(n+1-i)$$

High (exponential) complexity → many heuristics

- Exploring only left-linear plans etc.

N=4:



Logical query optimization needs statistics

Exact statistics (on base data):

- 1.000.000 cars, 1.000.000 drivers, 1.000 accidents

Approximate / estimated statistics (on intermediary results)

- "1.75 cars involved in every accident"

Statistics are gathered

- When **loading** the data: take advantage of the scan
- **Periodically** or upon **request** (e.g. analyze in the Postgres RDBMS)
- At **runtime**: modern systems may do this to change the data layout

Statistics on the **base data** vs. on **results of operations not evaluated** (yet):

- « On average 2 cars per accident »
- For each column R.a, store:
 $|R|$, $|R.a|$ (number of distinct values), $\min\{R.a\}$, $\max\{R.a\}$
- Assume **uniform distribution** in R.a
- Assume **independent distribution**
 - of values in R.a vs values in R.b; of values in R.a vs values in S.c
- + simple probability computations

More on statistics

- For each column R.a, store:
 - $|R|$, $|R.a|$ (number of distinct values), $\min\{R.a\}$, $\max\{R.a\}$
- Assume **uniform distribution** in R.a
- Assume **independent distribution**
 - of values in R.a vs values in R.b; of values in R.a vs values in S.c
- The **uniform distribution** assumption is **frequently wrong**
 - Real-world distribution are skewed (popular/frequent values)
- The **independent distribution** assumption is **sometimes wrong**
 - « Total » counter-example: *functional dependency*
 - Partial but strong enough to ruin optimizer decisions: *correlation*
- Actual optimizers use more sophisticated statistic informations
 - **Histograms**: equi-width, equi-depth
 - Trade-offs: size vs. maintenance cost vs. control over estimation error

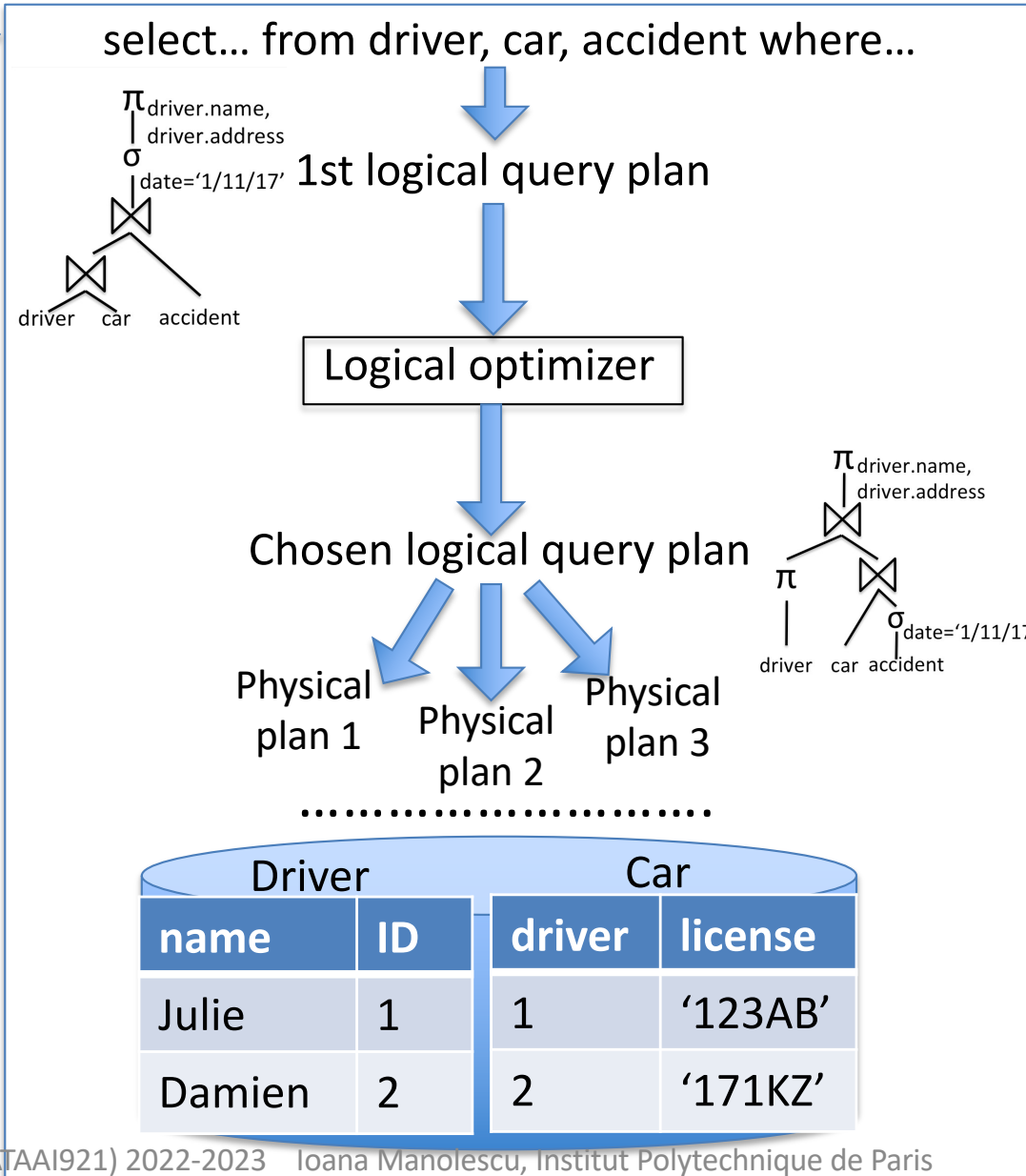
Database internal: query optimizer

SQL



```

select driver.name,
driver.address
from driver, car,
accident
where
driver.ID=car.driver
and
car.license=accident
.carLicense and
accident.date='1/11
/17'
    
```



Query language



Chosen logical plan

Results

Physical query plans

Made up of **physical operators** =
algorithms for implementing logical operators

Example: equi-join (R.a=S.b)

Nested loops join:

```
foreach t1 in R{
  foreach t2 in S {
    if t1.a = t2.b then output (t1 || t2)
  }
}
```

Merge join: // requires sorted inputs

```
repeat{
  while (!aligned) { advance R or S };
  while (aligned) { copy R into topR, S into topS };
  output topR x topS;
} until (endOf(R) or endOf(S));
```

Hash join: // builds a hash table in memory

```
While (!endOf(R)) {  $t_R \leftarrow R.next$ ; put(hash( $t_R.a$ ),  $t_R$ ); }
While (!endOf(S)) {  $t_S \leftarrow S.next$ ;
  matchingR = get(hash( $t_S.b$ ));
  output(matchingR x  $t_S$ );
}
```

Physical query plans

Made up of **physical operators** =
algorithms for implementing logical operators

Example: equi-join (R.a=S.b)

Nested loops join:

$O(|R| \times |S|)$

```
foreach t1 in R{
  foreach t2 in S {
    if t1.a = t2.b then output (t1 || t2)
  }
}
```

Merge join: // requires sorted inputs

```
repeat{
  while (!aligned) { advance R or S };
  while (aligned) { copy R into topR, S into topS };
  output topR x topS;
} until (endOf(R) or endOf(S));
```

$O(|R| + |S|)$

Hash join: // builds a hash table in memory

```
While (!endOf(R)) {  $t_R \leftarrow R.next$ ; put(hash( $t_R.a$ ),  $t_R$ ); }
```

```
While (!endOf(S)) {  $t_S \leftarrow S.next$ ;
```

```
    matchingR = get(hash( $t_S.b$ ));
```

```
    output(matchingR x  $t_S$ );
```

$O(|R| + |S|)$

```
}
```

Also:

Block nested loops join

Index nested loops join

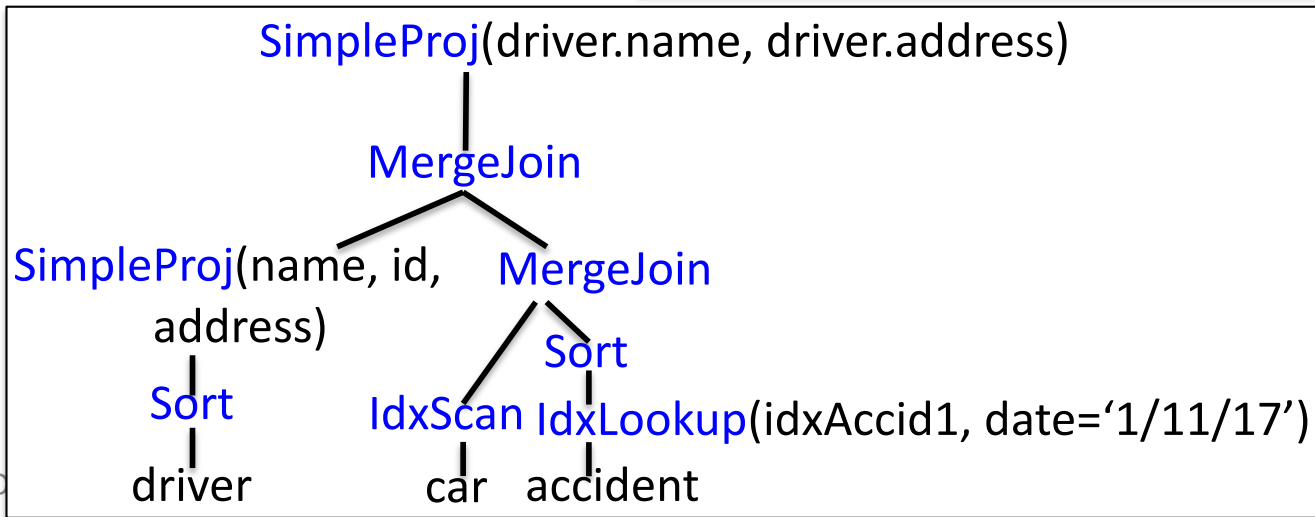
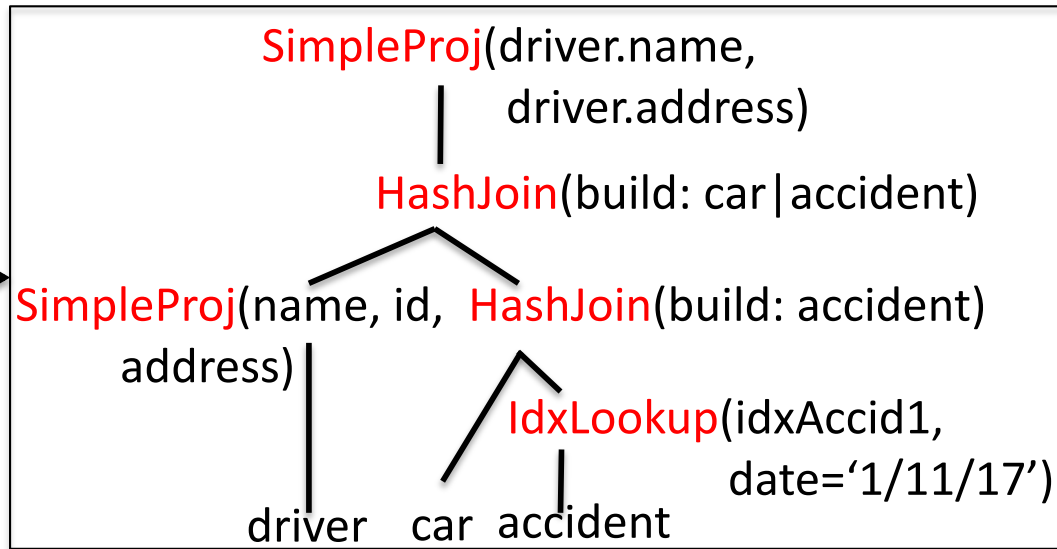
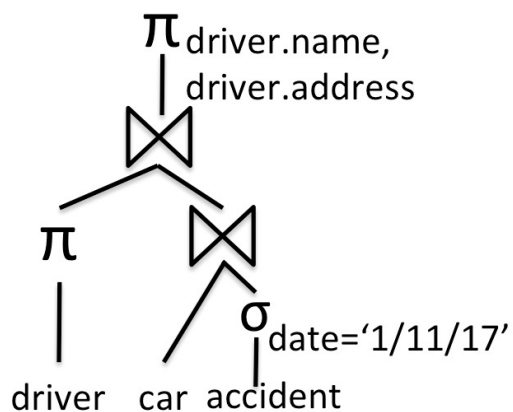
Hybrid hash join

Hash groups / teams

...

Physical optimization

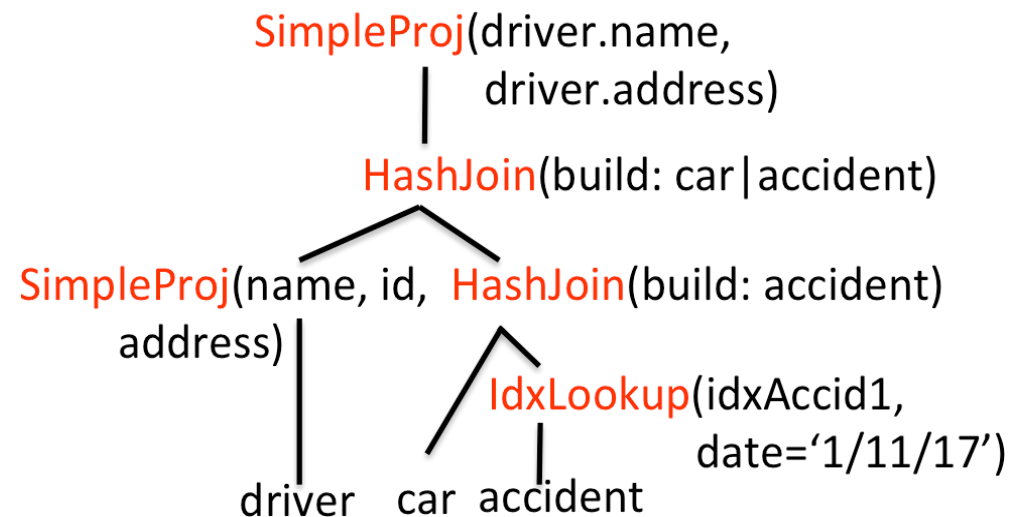
Possible physical plans produced by physical optimization for our sample logical plan:



Physical plan performance

Metrics characterizing a physical plan

- **Response time:** between the time the query starts running to the we know it's end of results
- **Work** (resource consumption)
 - How many **I/O** calls (blocks read)
 - Scan, IdxScan, IdxAccess; Sort; HybridHash (or spilling HashJoin)
 - How much **CPU**
 - All operators
 - Distributed plans: **network** traffic
- **Total work:** work made by all operators



Query optimizers in action

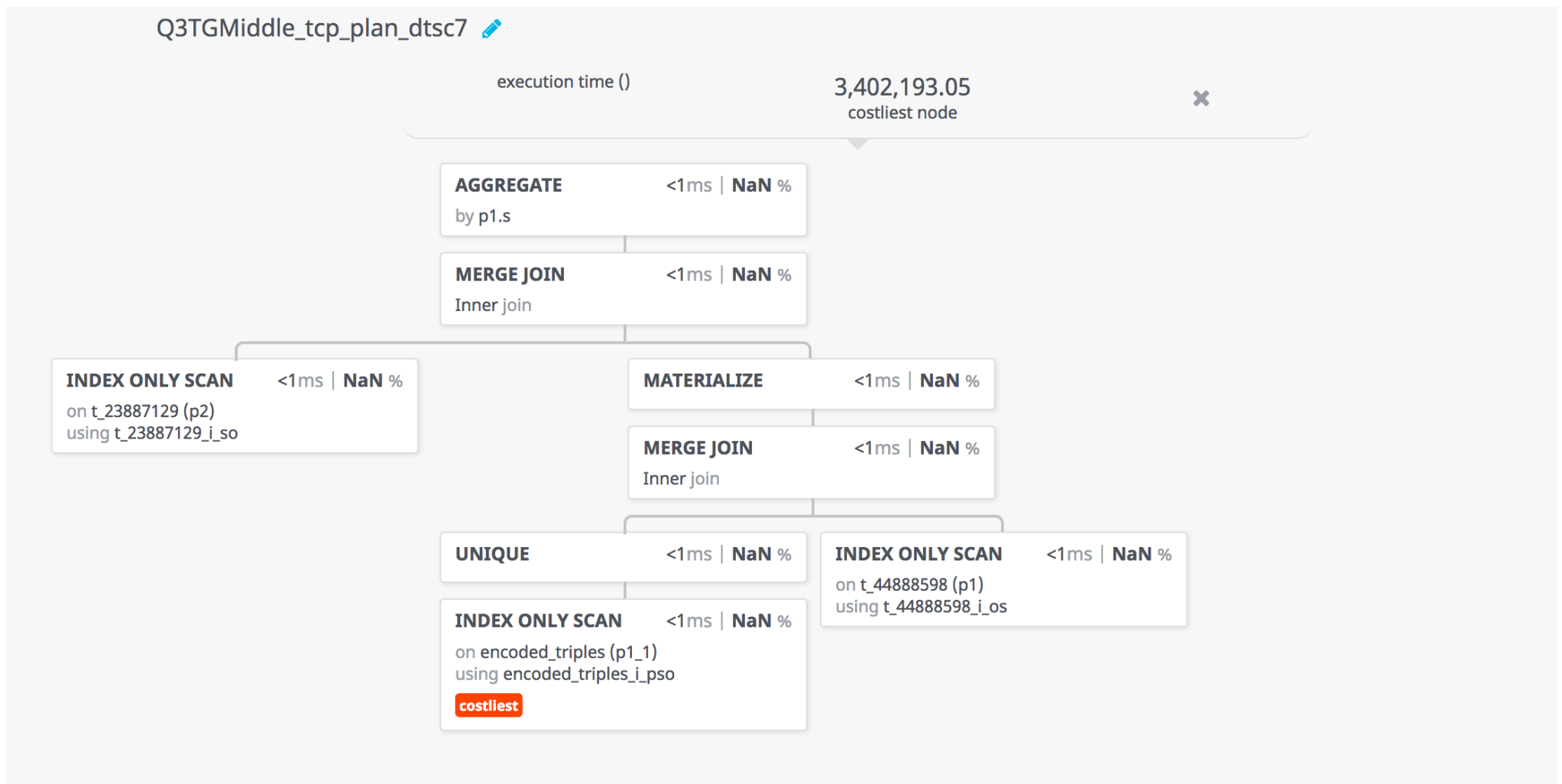
Most database management systems have an « explain » functionality → physical plans. Below sample Postgres output:

```
EXPLAIN SELECT * FROM tenk1;
      QUERY PLAN
-----
Seq Scan on tenk1 (cost=0.00..458.00 rows=10000 width=244)
```

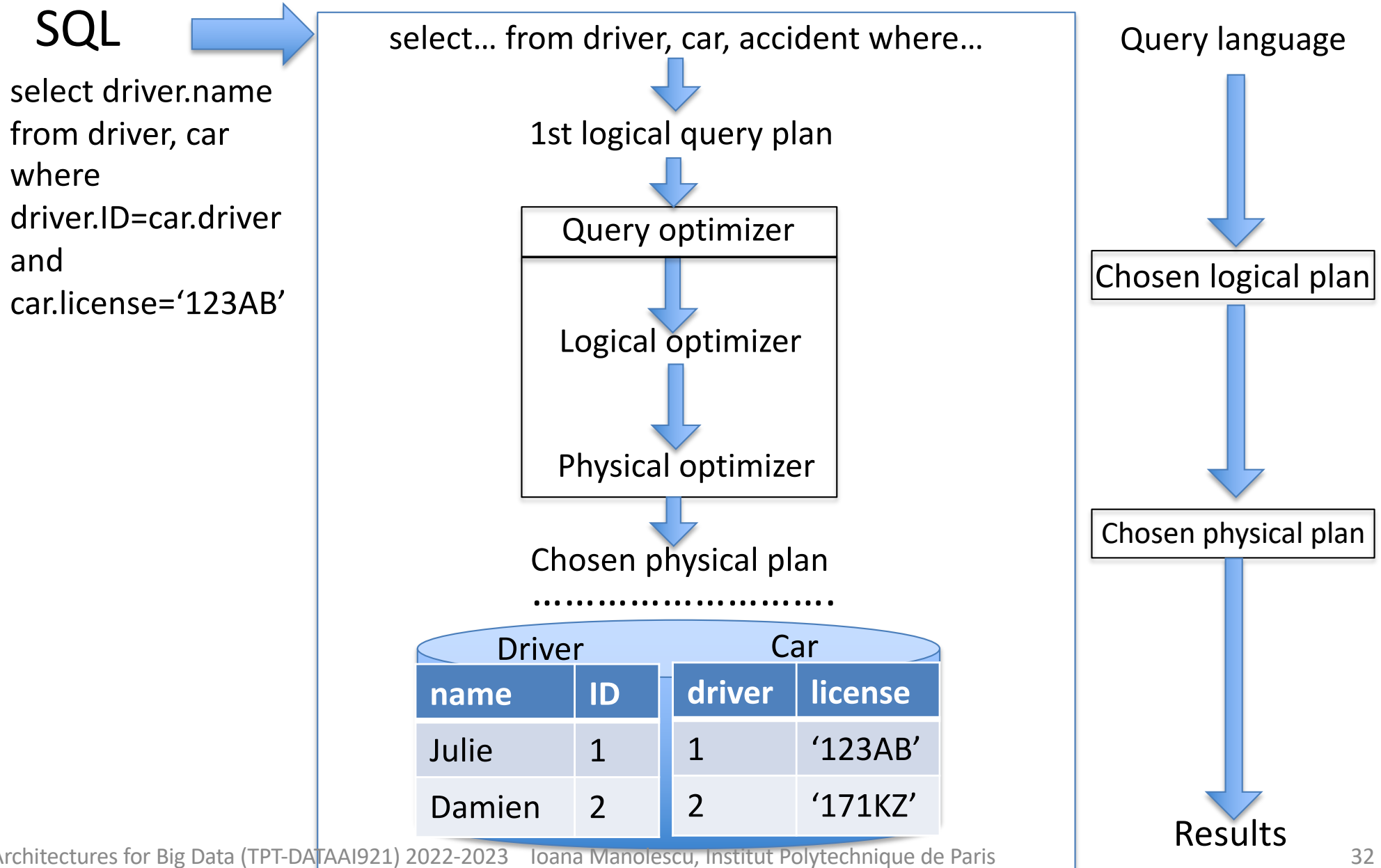
```
EXPLAIN SELECT * FROM tenk1 t1, tenk2 t2
WHERE t1.unique1 < 100 AND t1.unique2 = t2.unique2;
      QUERY PLAN
-----
Hash Join (cost=232.61..741.67 rows=106 width=488)
  Hash Cond: ("outer".unique2 = "inner".unique2)
    -> Seq Scan on tenk2 t2 (cost=0.00..458.00 rows=10000 width=244)
    -> Hash (cost=232.35..232.35 rows=106 width=244)
          -> Bitmap Heap Scan on tenk1 t1 (cost=2.37..232.35 rows=106 width=244)
                Recheck Cond: (unique1 < 100)
                -> Bitmap Index Scan on tenk1_unique1 (cost=0.00..2.37 rows=106 width=0)
                      Index Cond: (unique1 < 100)
```

Inspecting query plans

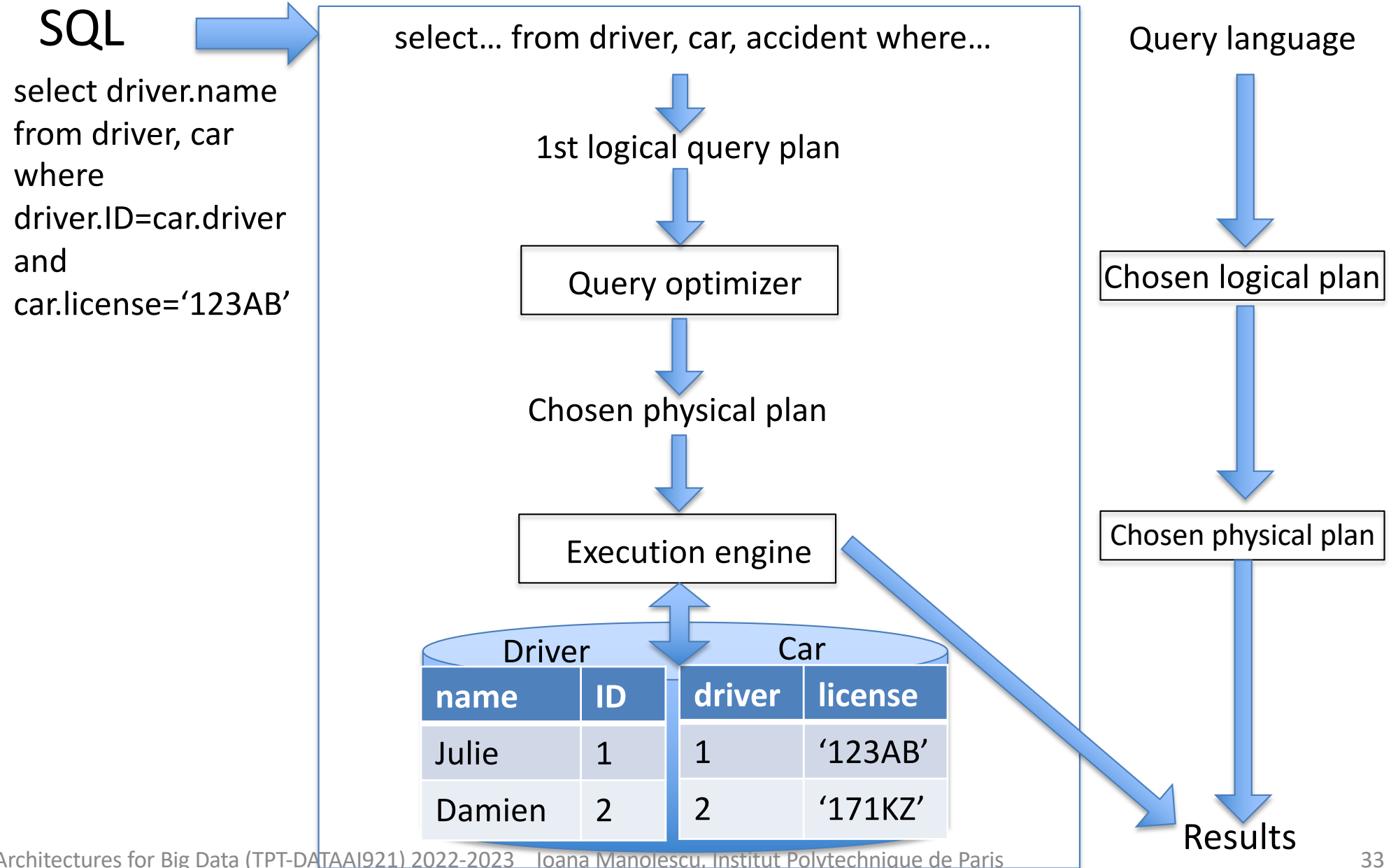
- Here using <https://tatiyants.com/pev/#/plans>



Database internal: physical plan



Database internals: query processing pipeline

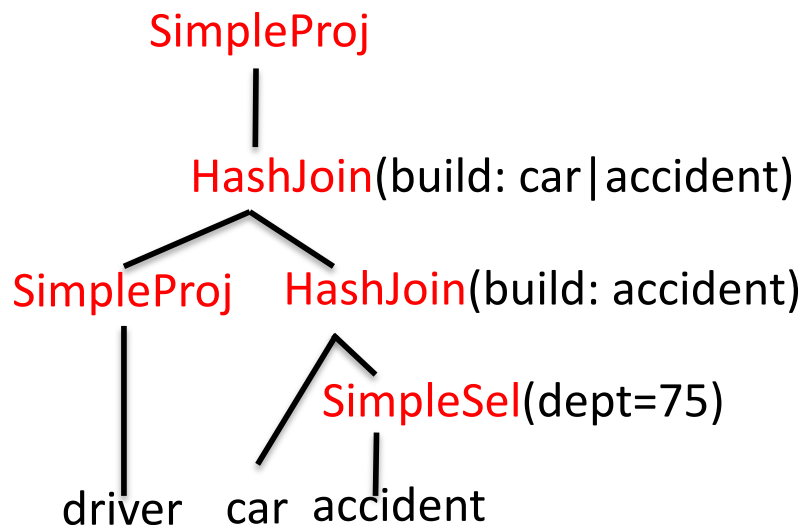


Advanced query optimization techniques:

Dynamic Query Optimization

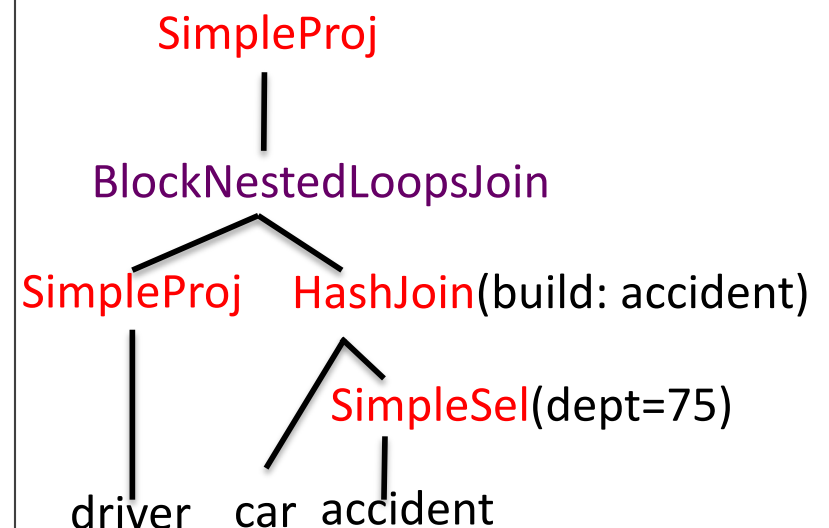
- Sizes (cardinalities) of intermediary results are estimated, which may lead to estimation errors
- A cardinality estimation error may lead to choosing a logical plan and a set of physical operators that perform significantly different from expectation (especially for the worse)

Initially chosen plan:



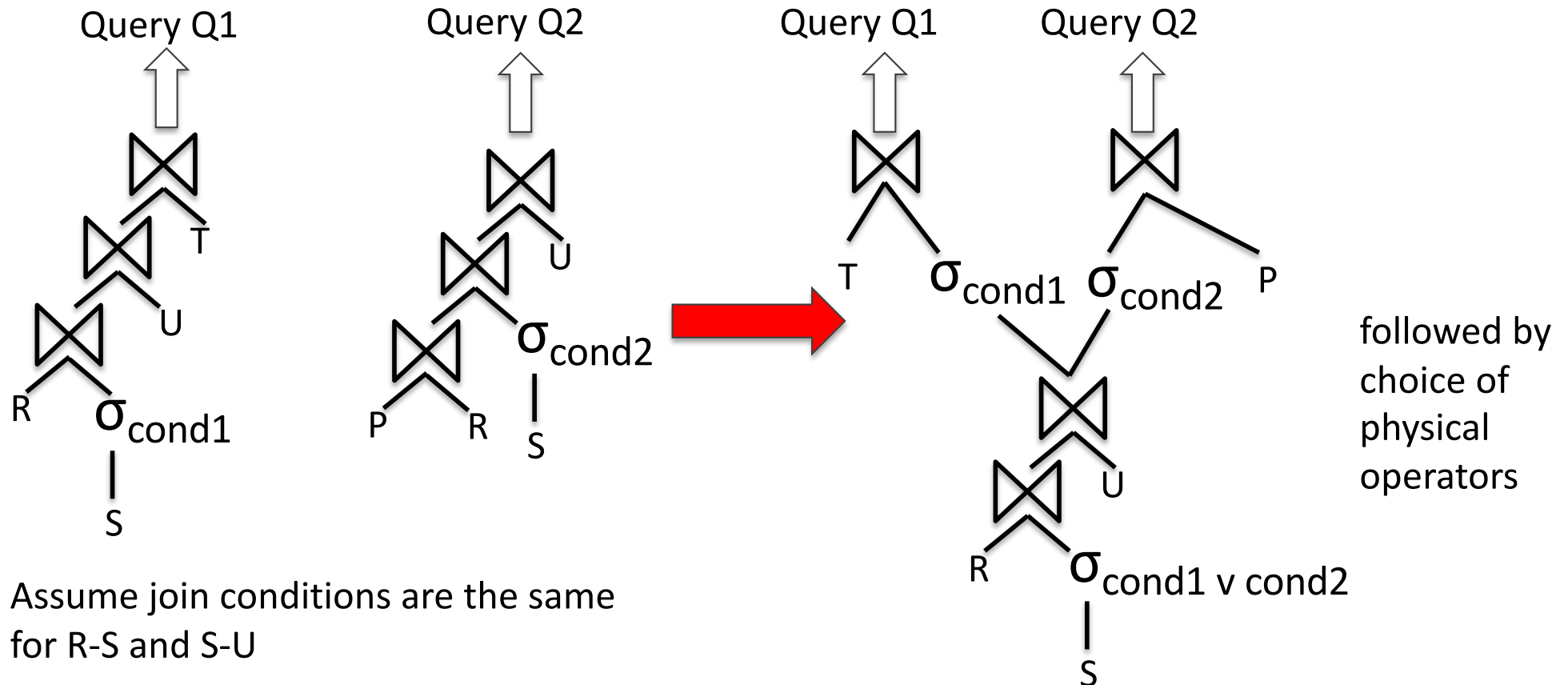
At execution time, we see that the lower HashJoin output is larger than expected: memory insufficient to build

Modified plan:



Advanced query optimization techniques: Multi-Query Optimization

Multiple queries sharing sub-expressions can be optimized together into a single plan with **shared subexpressions**




DATABASE FUNDAMENTALS (RECALL/CRASH COURSE)

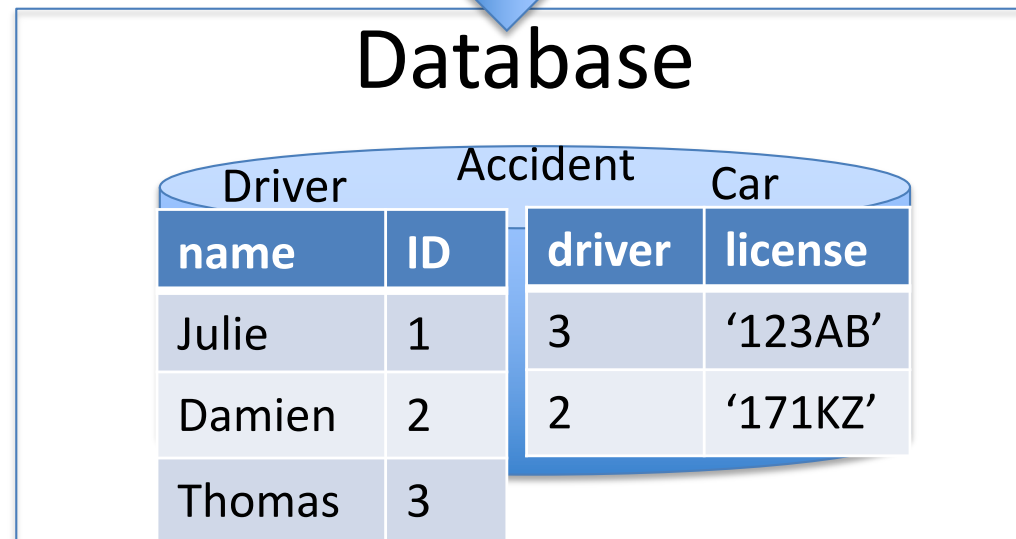
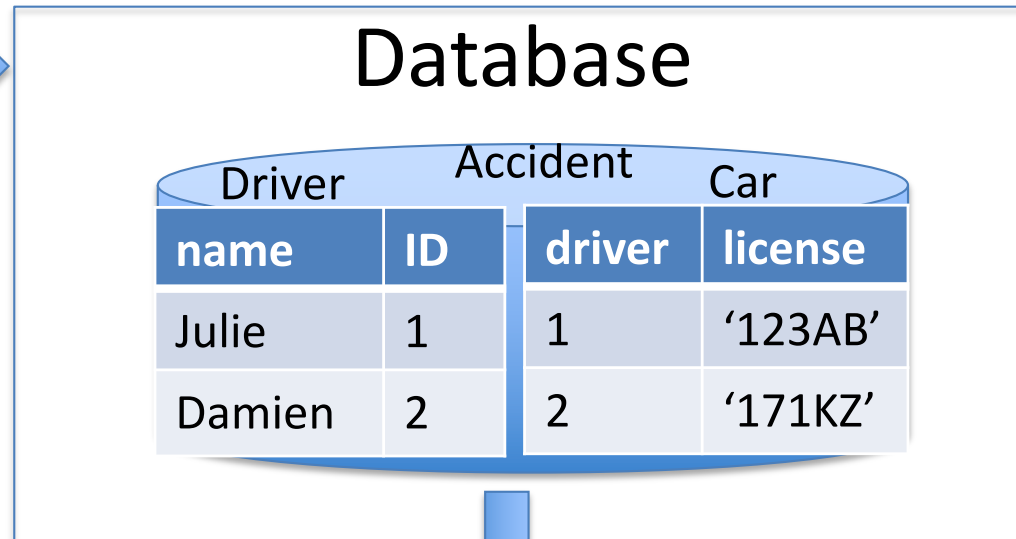
Updating the database

What's in a database?

SQL
update



insert into driver
values ('Thomas',
3);
update car set
driver=3 where
license='123AB';



Database updates

- A set of operations atomically executed (either all, or none) is called a **transaction**
- There may be some **dependencies** between the operations of a transaction
 - First read the bank account balance
 - Then write that value reduced by 50€
- A total order over the operations of several concurrent transaction is called a **scheduling**
- The DB component that receives all incoming transactions and decides what operation will be executed when (i.e., global order over the operations of all transactions) is the **scheduler**

Database updates

- The **scheduler** is in charge of ordering all operations so that they will appear executed one after the other (serially)

```
T1: BEGIN  A=A+100,  B=B-100  END
T2: BEGIN  A=1.06*A, B=1.06*B  END
```

```
T1: A=A+100, B=B-100,
T2:                A=1.06*A,  B=1.06*B
```

```
T1: A=A+100,                B=B-100
T2:                A=1.06*A,                B=1.06*B
```

Ensuring consistency of concurrent updates

- **Scheduling** is implemented through specific algorithms and with the help of protocols
- A **protocol** is a rule that holds on the order in which a transaction performs its operations
 - E.g., "once a transaction has released a lock, the transaction will never take another lock"
- If all transactions follow a given protocol, then, regardless on the order in which they are executed, certain **good properties** are guaranteed
 - E.g., "there is no deadlock" or "the result is the same as if the transactions had been executed one after the other"

Fundamental database features

1. Data storage

- Protection against unauthorized access, data loss

2. Ability to at least **add** to and **remove** data to the database

- Also: **updates; active behavior** upon update (triggers)

3. Support for **accessing** the data

- Declarative query languages: say what data you need, not how to find it

Fundamental properties of database stores: ACID

- **A**tomicity: either all operations involved in a transactions are done, or none of them is
 - E.g. bank payment
- **C**onsistency: application-dependent constraint
 - E.g. every client has a single birthdate
- **I**solation: concurrent operations on the database are executed as if each ran alone on the system
 - E.g. if a debit and a credit operation run concurrently, the final result is still correct
- **D**urability: data will not be lost nor corrupted even in the presence of system failure during operation execution

Jim Gray, ACM Turing Award 1998 for « fundamental contributions to databases and transaction management »

ACID properties

- **Atomicity:** per transaction (cf. boundaries)
- **Consistency:** difference in the expressive power of the constraints
- Illustrated below for relational databases, **create table** statement:

```
CREATE TABLE tbl_name (create_definition,...) [table_options] [partition_options]
```

```
create_definition: col_name column_definition |
```

```
  [CONSTRAINT [symbol]] PRIMARY KEY [index_type] (index_col_name,...) [index_option] ... |
```

```
  {INDEX|KEY} [index_name] [index_type] (index_col_name,...) [index_option] ... |
```

```
  [CONSTRAINT [symbol]] UNIQUE [INDEX|KEY] [index_name] [index_type]
```

```
                                (index_col_name,...) [index_option] (...) |
```

```
  CHECK (expr)
```

```
column_definition: data_type [NOT NULL | NULL] [DEFAULT default_value]
```

```
  [AUTO_INCREMENT] [UNIQUE [KEY] | [PRIMARY] KEY] (...)
```

ACID properties

Consistency (continued)

- SQL constraint syntax (within create table):

```
[CONSTRAINT [symbol]] FOREIGN KEY [index_name]  
(index_col_name, ...)
```

```
    REFERENCES tbl_name (index_col_name,...)
```

```
    [ON DELETE reference_option]
```

```
    [ON UPDATE reference_option]
```

```
reference_option: RESTRICT | CASCADE | SET NULL | NO ACTION
```

- Key-value store: REDIS
 - a data item can have only one value for a given property
- Key-value store: DynamoDB
 - The value of a data item can be constrained to be unique, *or* allowed to be a set
- Hadoop File System (HDFS): no constraints

ACID properties

- **Isolation**: concurrent operations on the database are executed as if each ran alone on the system
 - Watch out for: read-write (RW) or write-write (WW) conflicts
 - Conflict granularity depends on the data model
- **An example of advanced isolation support: SQL**
 - E.g. SQL

Isolation Level	Dirty Read	Non Repeatable Read	Phantom
Read uncommitted	Yes	Yes	Yes
Read committed	No	Yes	Yes
Repeatable read	No	No	Yes
Snapshot	No	No	No
Serializable	No	No	No

- High isolation conflicts with high transaction throughput
- E.g. HDFS: a file is never modified (written only once and integrally)

DATABASE FUNDAMENTALS (RECALL/CRASH COURSE)

Takeaway

Main principles behind correct and scalable data management...

... pioneered in database management systems:

1. **Declarative query language** allows users to just state what they want
2. For one query there are several **logical plans**; for each, several **physical plans**
 - Optimizer picks **best plan**
3. **ACID** properties crucial for "faith in the system" ("my salary, payments, and social security are within a reliable system")