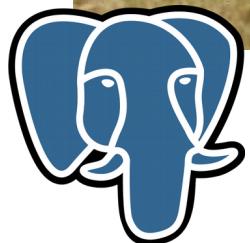
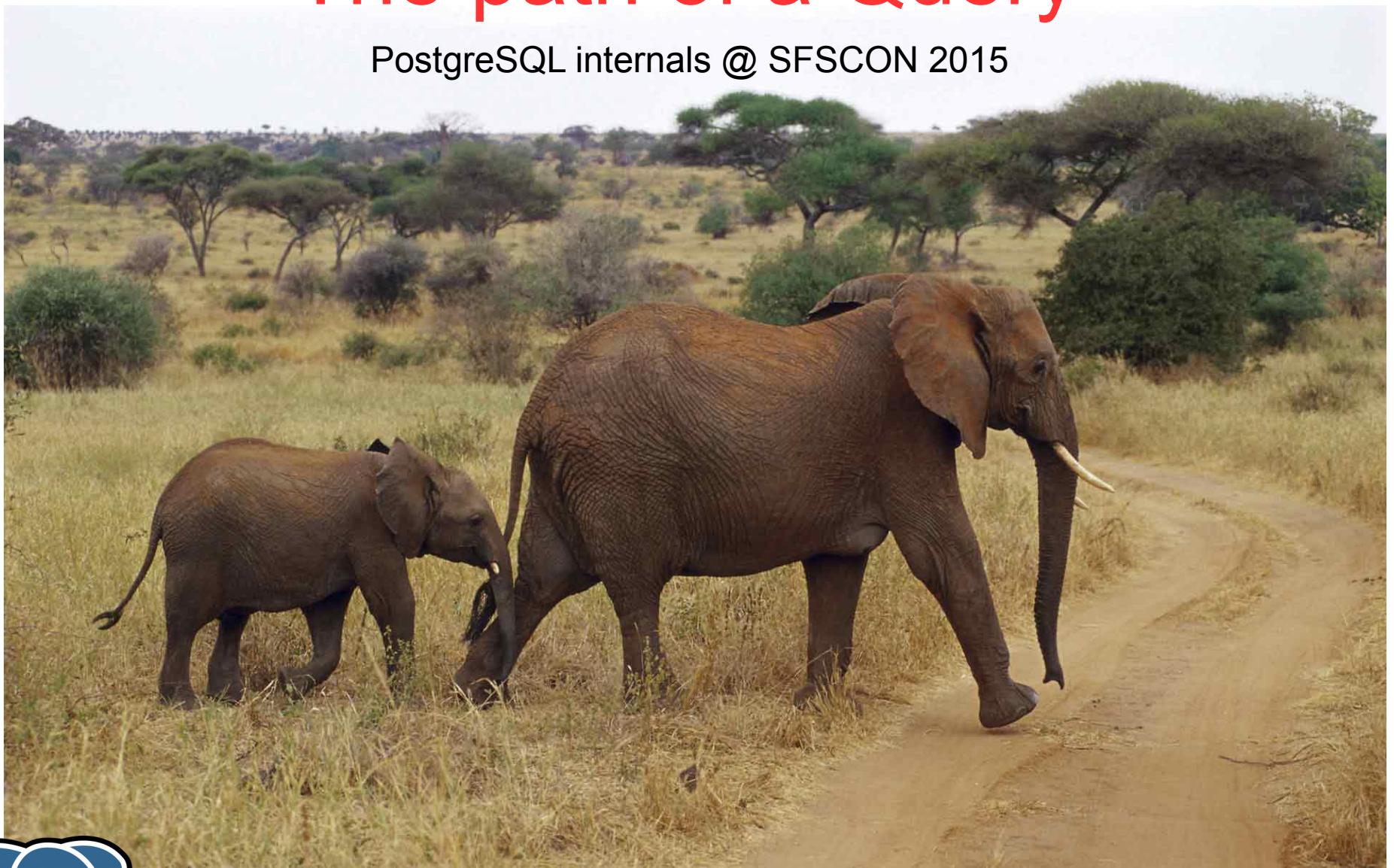


The path of a Query

PostgreSQL internals @ SFSCON 2015



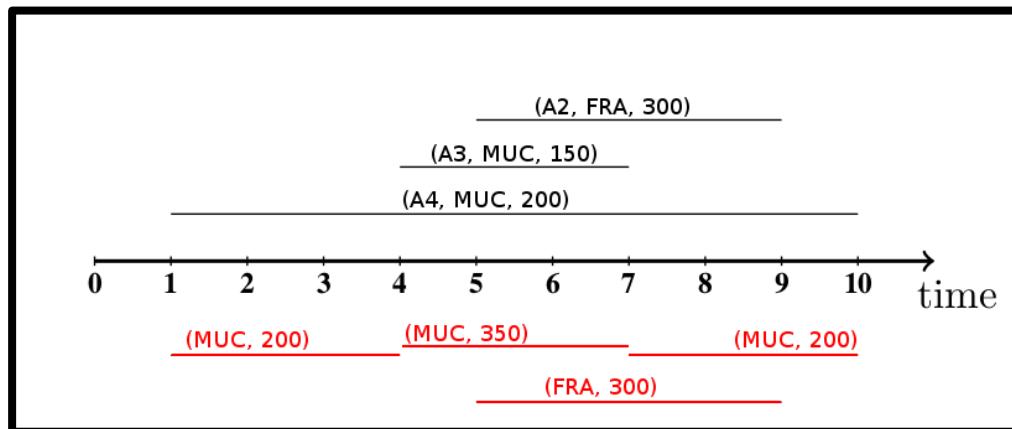
A glance into the
PostgreSQL kernel

by Peter Moser,
master student @ **unibz**

as an
intern @ **TiS**
innovation park

Our Project: Temporal PostgreSQL

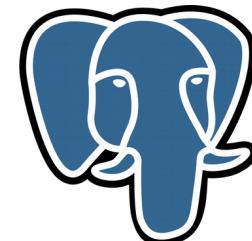
- Temporal queries **missing**



Temporal Aggregation Example

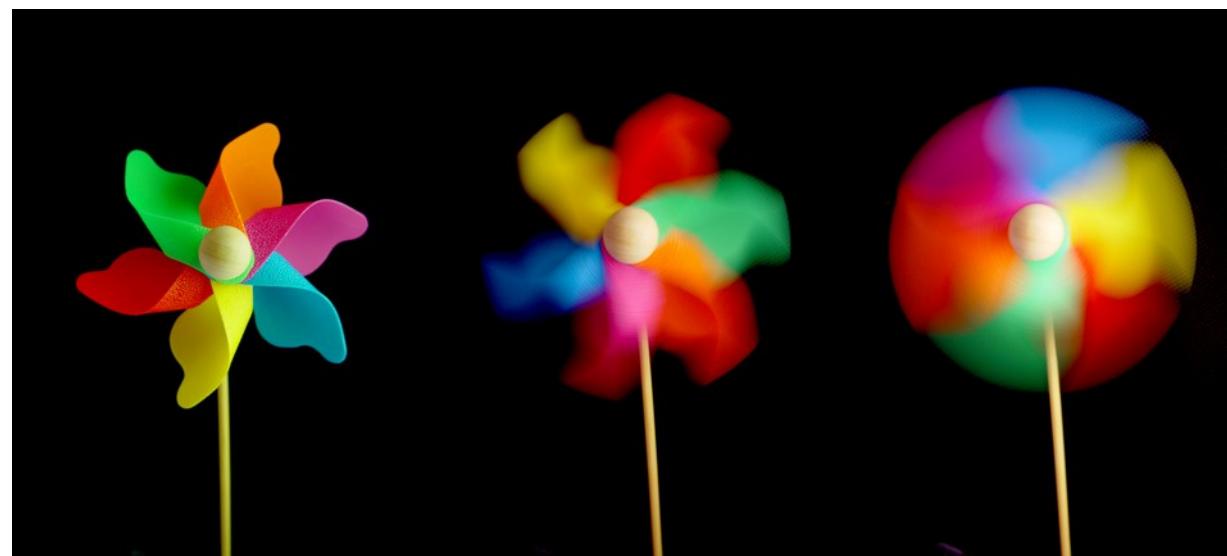


- Ongoing **research** at the Free University of BZ
- Extending the PostgreSQL kernel
with **temporal operators**



Why creating patches for PostgreSQL?

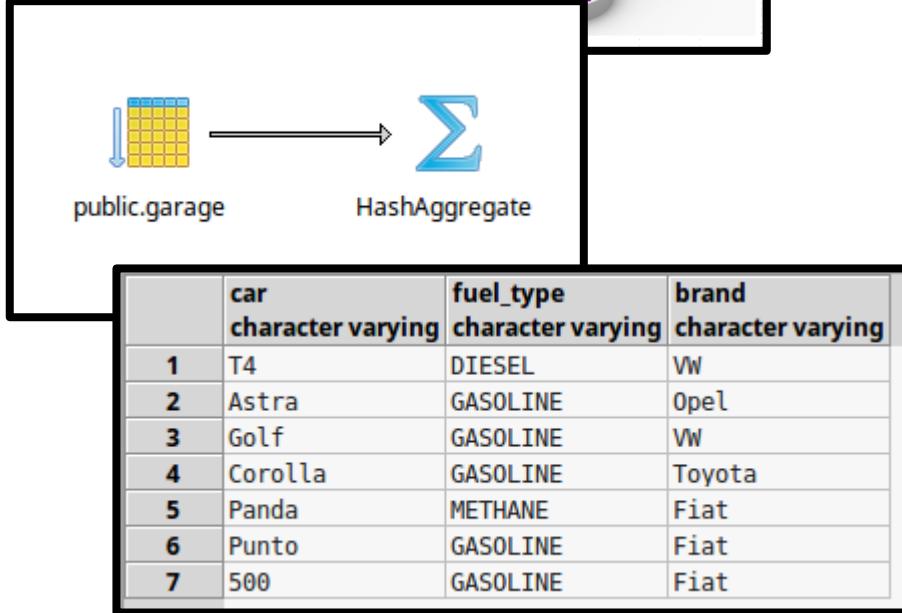
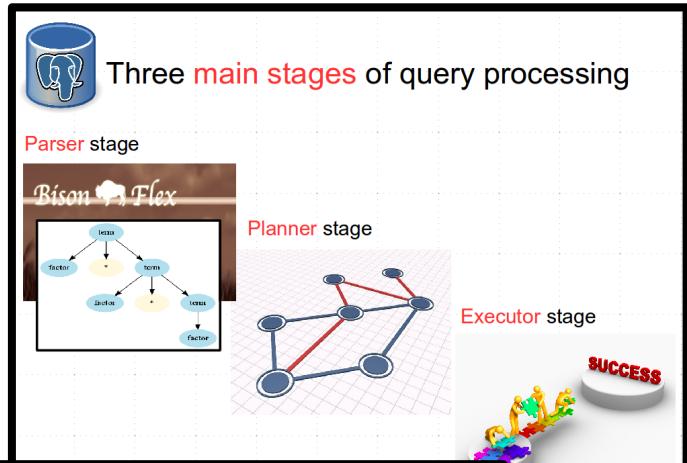
- Extending SQL **operators** (not only functions)
- **SQL standard** compliance (missing features)
- **Performance / Memory** usage



Talk (overview)

vs.

Workshop (code)



The screenshot shows a developer's environment with three main windows:

- Terminal:** Displays PostgreSQL logs from a 'postgretemporal' database, showing queries related to table creation and analysis.
- IDE (C/C++ project):** Shows a C++ project named 'parser_rewritten'. The code is a rewritten parser for SQL-like statements, specifically handling 'SELECT' clauses. It includes comments explaining the logic for handling various parts of the query, such as 'WITH' clauses, table references, and sorting.
- Code Editor (File browser):** Shows a file named 'relation.c'. This file contains a large switch statement that handles different types of relations based on their names. It includes many annotations and comments explaining the logic for each case, such as handling 'WITH' clauses, table references, and sorting.

How does Postgres process a Query?



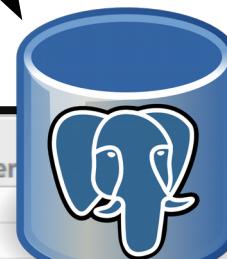
How many cars with a specific fuel type are in our garage?

WHAT? (Query)

```
SELECT  
    fuel_type, COUNT(*)  
FROM garage  
GROUP BY fuel_type;
```



HOW?



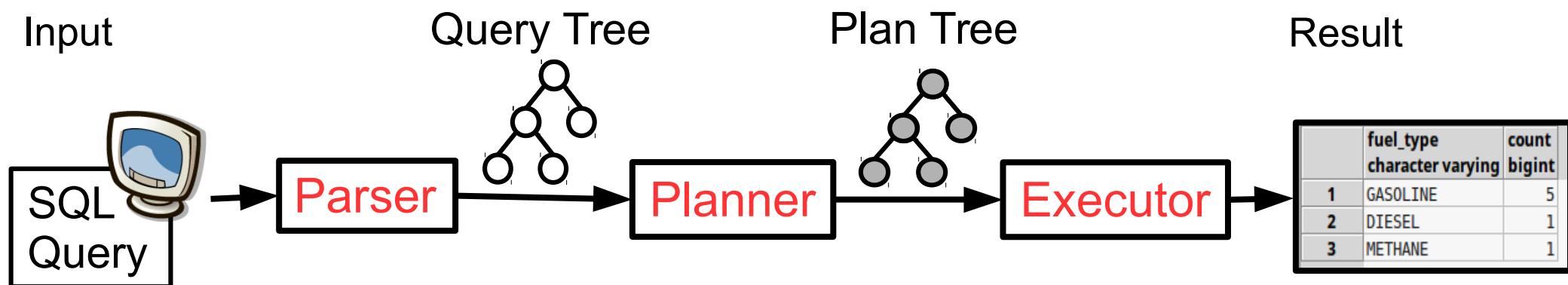
	car	fuel_type	brand
	character varying	character varying	character
1	T4	DIESEL	VW
2	Astra	GASOLINE	Opel
3	Golf	GASOLINE	VW
4	Corolla	GASOLINE	Toyota
5	Panda	METHANE	Fiat
6	Punto	GASOLINE	Fiat
7	500	GASOLINE	Fiat

DATA

Result

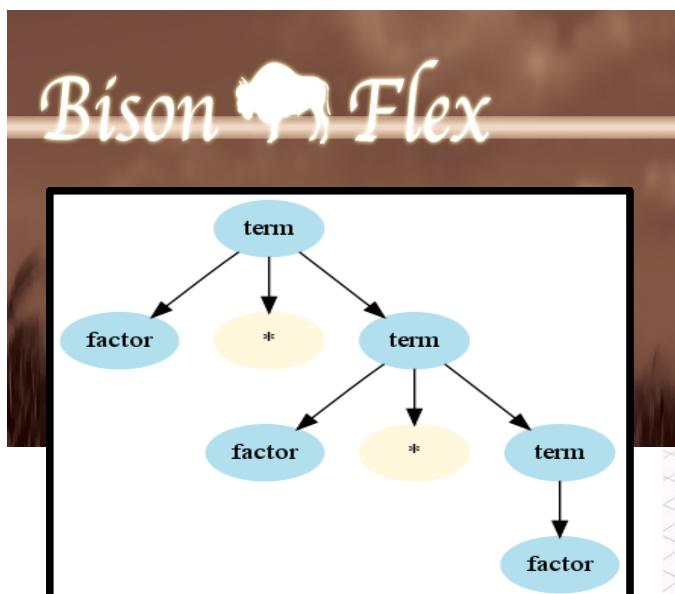
	fuel_type	count
	character varying	bigint
1	GASOLINE	5
2	DIESEL	1
3	METHANE	1

The path of a query from bird's-eye view

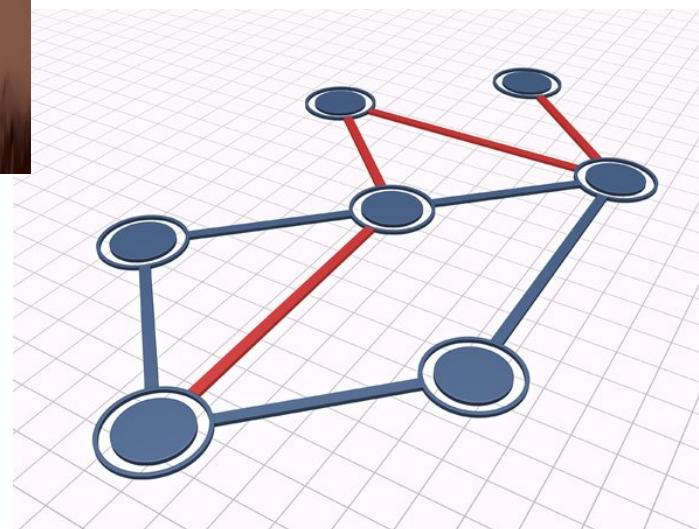


Three main stages of query processing

Parser stage

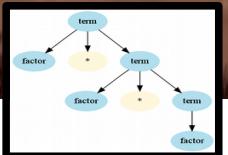


Planner stage



Executor stage





Parser stage

- Checks the SQL syntax

- Syntax error:

```
SELECT * FRAM garage;
```

“FRAM” is not
a SQL key word

- Makes lookups in a catalog

- Error in semantics:

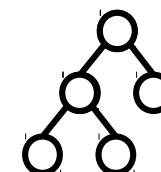
```
SELECT * FROM garauge;
```

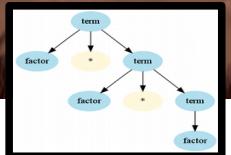
- Adds information
(eg. column names, aliases)

Table “garauge”
does not exist

- Builds the Query Tree

- Internal representation of a Query





Parser stage: Identify tokens

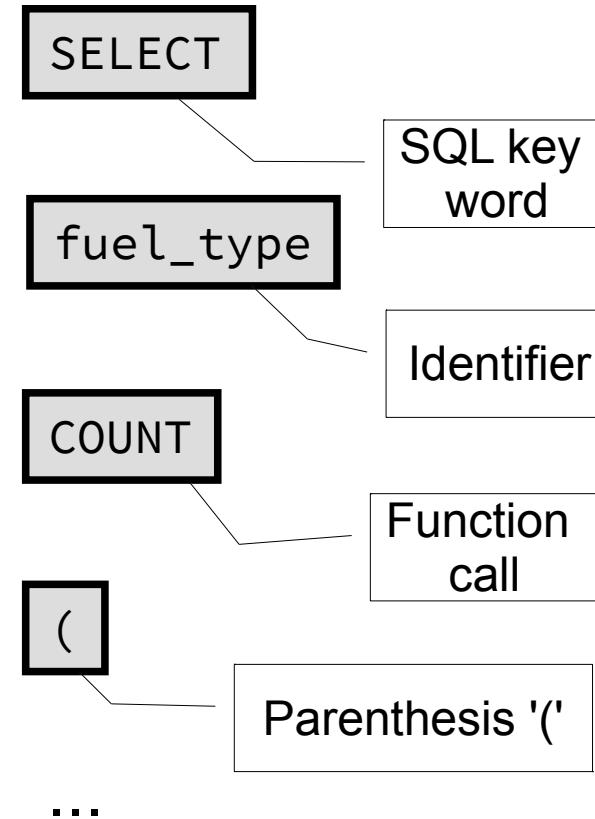
Query String
arrives as plain text

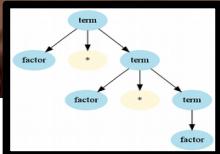
```
SELECT
  fuel_type, COUNT(*)
FROM garage
GROUP BY fuel_type;
```



(C) wiseGeek

The parser
generates tokens





Parser stage: Build query tree

Grammar rules

simple_select:

```
SELECT target_list
FROM from_clause
GROUP BY group_clause
```

target_list:

```
target_el
| target_list ',' target_el
```

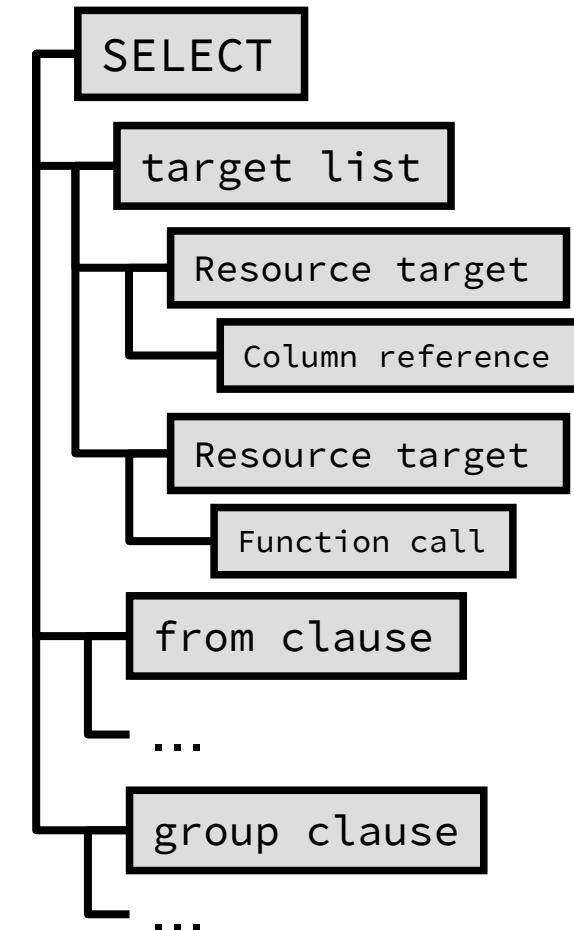
target_el:

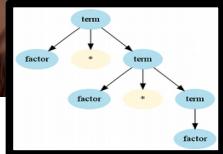
```
a_expr AS ColLabel
| a_expr
| '*'
```

a_expr:

```
...
```

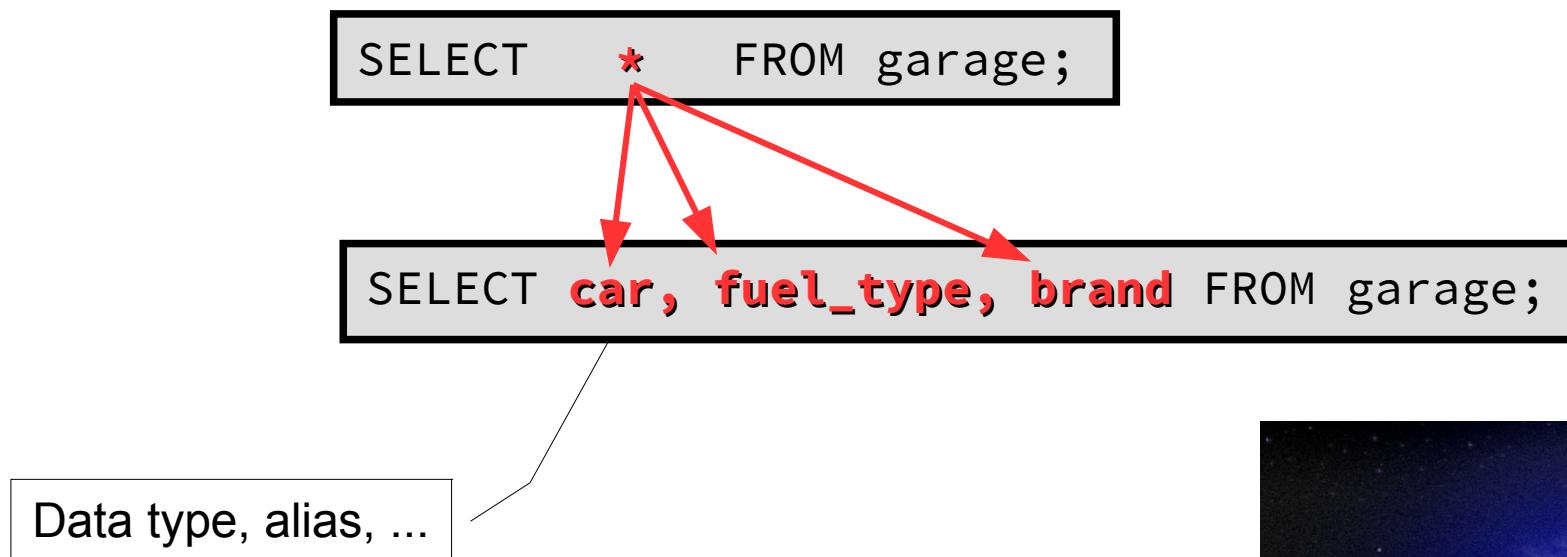
Query tree





Query Tree Annotation

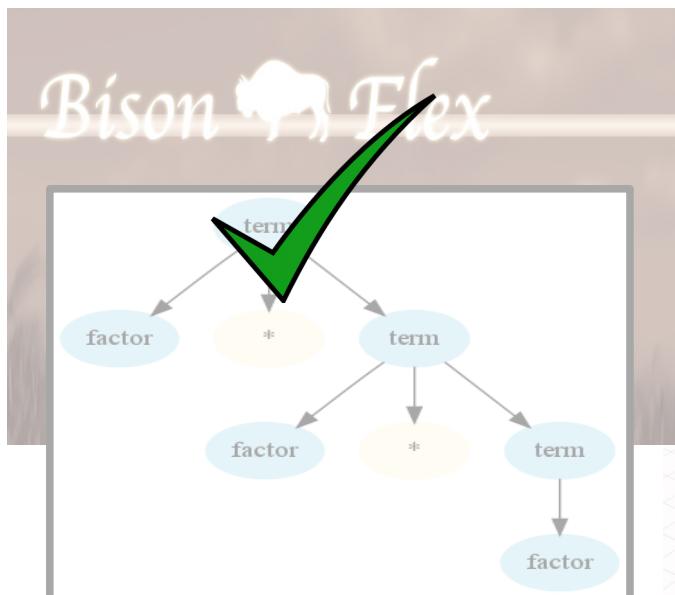
- ◆ Example: A_STAR expansion



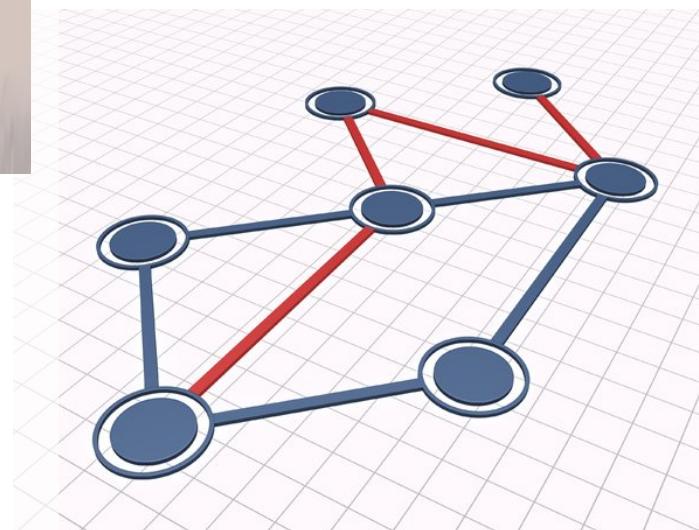
(C) wallpaperhi

Three main stages of query processing

Parser stage

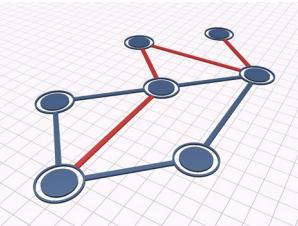


Planner stage



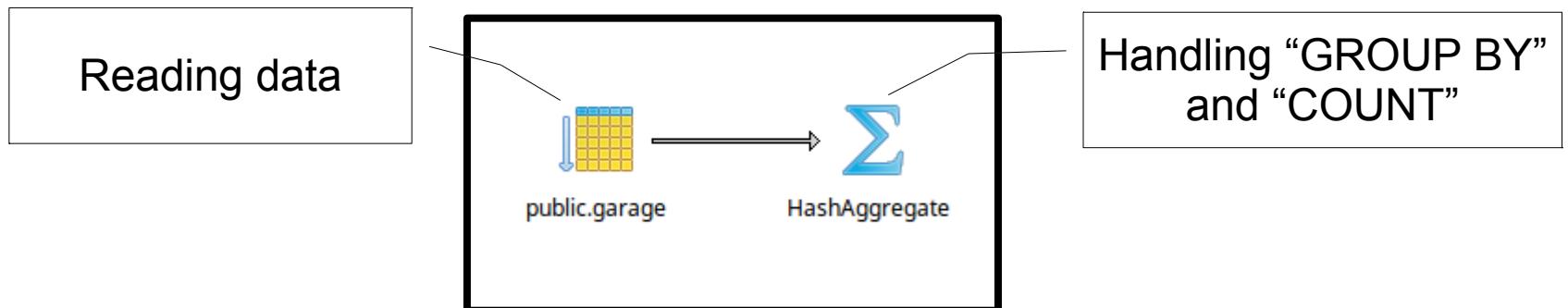
Executor stage

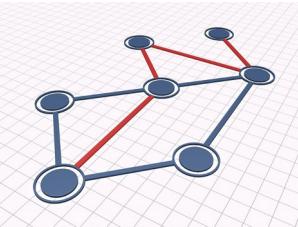




Planner stage

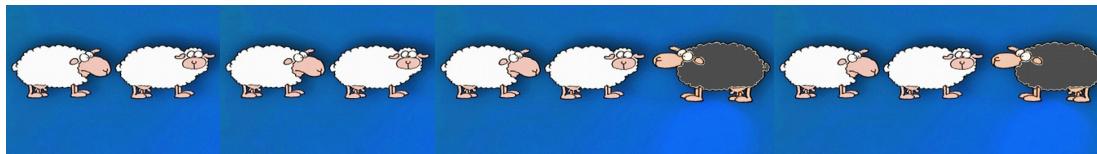
- ◆ There is a standard way to execute SQL
- ◆ The planner tries to find the best **plan tree**
- ◆ Creates the **Plan Tree**





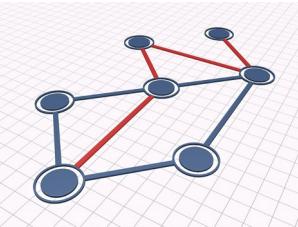
Planner decisions

- Access Scan Types
 - Sequential Scan, (Bitmap) Index Scan, Index-Only Scan



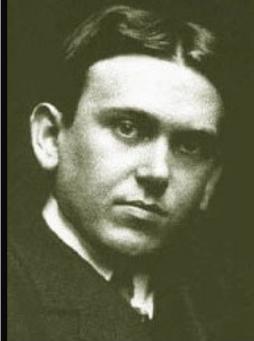
- Join
 - Join order
 - Join strategy: nested loop, merge join, hash join
 - Inner vs. Outer join
- Aggregation
 - Plain, sorted, hashed

MIN, MAX, COUNT, AVG, SUM, ...



Planner stage

- ◆ Query planning is a **complex problem**
- ◆ **Finding all** plans could be **too expensive**
- ◆ If it is not feasible, it uses **heuristics**

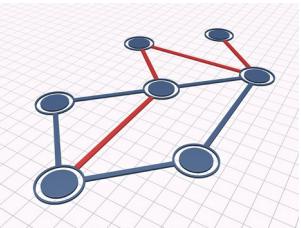


For every complex problem there is an answer that is clear, simple, and wrong.

(H. L. Mencken)

izquotes.com

A black and white portrait of H. L. Mencken is on the left side of the quote box. The quote itself is in white text on a dark background. The source, izquotes.com, is at the bottom right of the box.

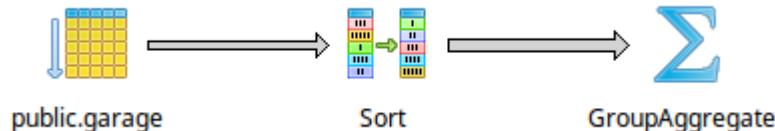


Planner example

```
SELECT fuel_type, COUNT(*)  
FROM garage  
GROUP BY fuel_type;
```

50

Plan #1



Cost estimates

- Scan: 16
- Sort: 48
- Group: 50

20

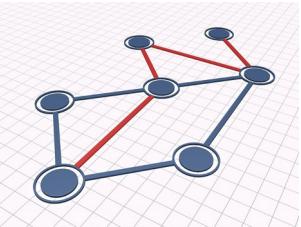
Plan #2



Cost estimates

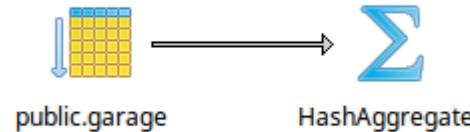
- Scan: 16
- HashAgg: 20

Smaller cost means faster!



is Plan #2

```
SELECT fuel_type, COUNT(*)
FROM garage
GROUP BY fuel_type;
```



FINAL CHOICE:
Sequential Scan
on “public.garage”
and
aggregation with
a hash

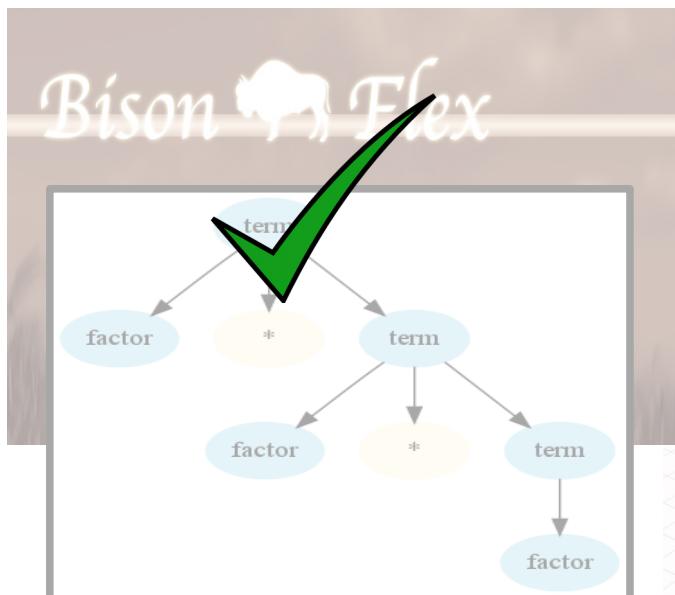
Explaining the Query Plan Tree in pgAdmin

QUERY PLAN	
	text
1	HashAggregate (cost=19.75..21.75 rows=200 width=32) (actual time=0.024..0.025 rows=3 loops=1)
2	Group Key: fuel type
3	-> Seq Scan on garage (cost=0.00..16.50 rows=650 width=32) (actual time=0.010..0.012 rows=7 loops=1)
4	Planning time: 0.071 ms
5	Execution time: 0.075 ms

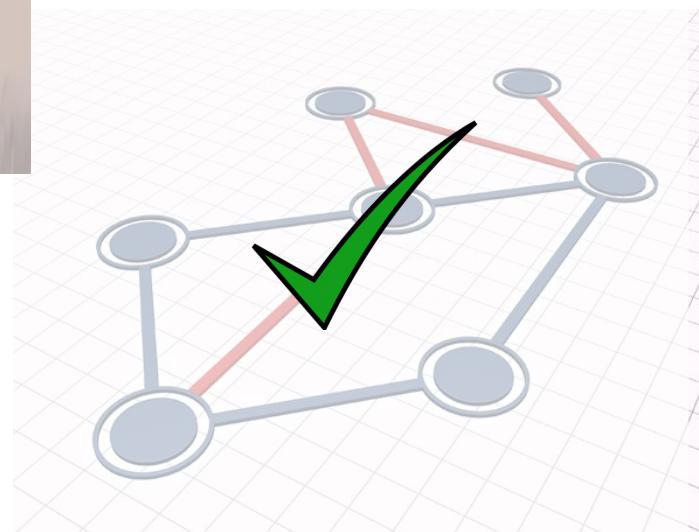


Three main stages of query processing

Parser stage



Planner stage



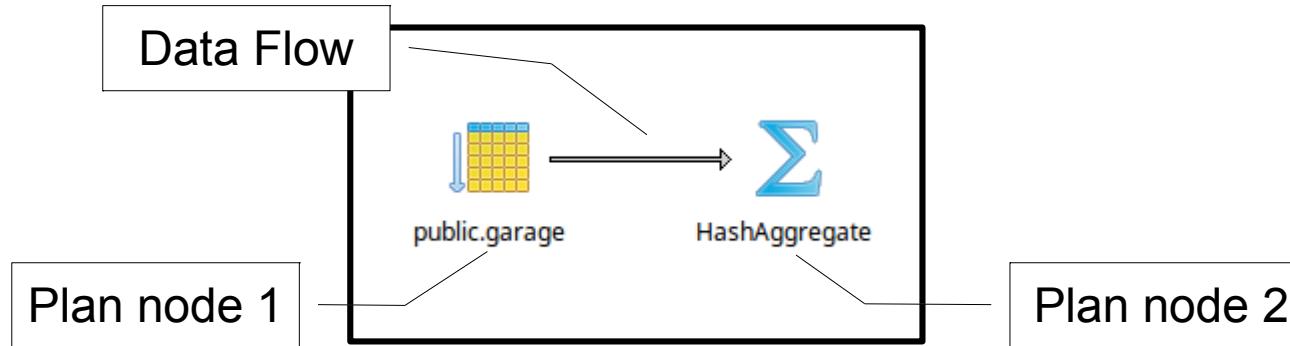
Executor stage





Executor stage

- ◆ Processes the (fastest) Plan Tree



- ◆ Each plan node returns a single row **on demand**
- ◆ Creates the **Result**

	fuel_type character varying	count bigint
1	GASOLINE	5
2	DIESEL	1
3	METHANE	1



Executor implementation

- ◆ Each executor node has a strict **interface**

ExecInitNode()

ExecProcNode()

...

ExecEndNode()

- ◆ Interface routines:

ExecInitNode():

initialize a plan node and its
subplans

ExecProcNode():

get a tuple by executing
the plan node
return the tuple, or NULL

ExecEndNode():

shut down a plan node and its
subplans

Calling
“ExecProcNode”
twice

	car character varying	fuel_type character varying
1	T4	DIESEL
2	Astra	GASOLINE
3	Golf	GASOLINE
4	Corolla	GASOLINE
5	Panda	METHANE
6	Punto	GASOLINE
7	500	GASOLINE



Executor stage by example

- ◆ Our query (+ ordering):

```
SELECT fuel_type, COUNT(*)  
FROM garage  
GROUP BY fuel_type  
ORDER BY count DESC;
```

Data:

	car character varying	fuel_type character varying	brand character varying
1	T4	DIESEL	VW
2	Astra	GASOLINE	Opel
3	Golf	GASOLINE	VW
4	Corolla	GASOLINE	Toyota
5	Panda	METHANE	Fiat
6	Punto	GASOLINE	Fiat
7	500	GASOLINE	Fiat

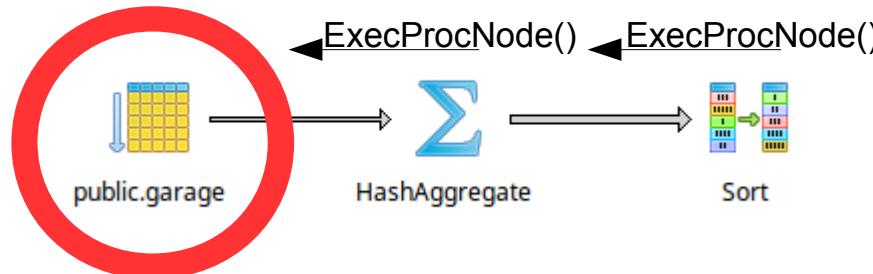
- ◆ Execution Plan:



- ◆ The executor **recursively calls itself** to process subplans (starting from the top node Sort)



Executor stage



```
SELECT fuel_type, COUNT(*)  
FROM garage  
GROUP BY fuel_type  
ORDER BY count DESC;
```

- ◆ Run the **Sequential Scan** node
- ◆ Return the result to the node **HashAggregate**
- ◆ Input:

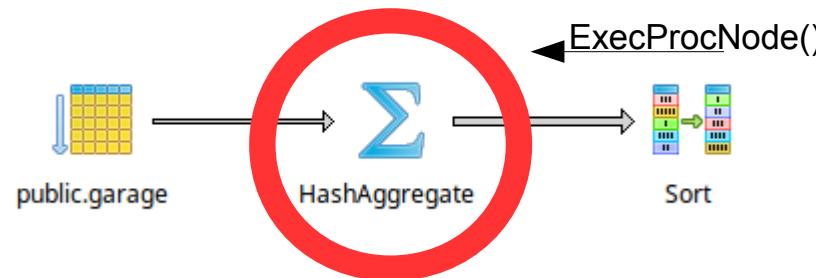


Result:

	car character varying	fuel_type character varying	brand character varying
1	T4	DIESEL	VW
2	Astra	GASOLINE	Opel
3	Golf	GASOLINE	VW
4	Corolla	GASOLINE	Toyota
5	Panda	METHANE	Fiat
6	Punto	GASOLINE	Fiat
7	500	GASOLINE	Fiat



Executor stage



```
SELECT fuel_type, COUNT(*)  
FROM garage  
GROUP BY fuel_type  
ORDER BY count DESC;
```

- ◆ Build a hash upon “fuel_type”
- ◆ Use the hash to find candidates for fuel type groups
- ◆ Grouping by hash

	car character varying	fuel_type character varying	brand character varying	hash text
Group #1	1 T4	DIESEL	VW	167
Group #2	2 Panda	METHANE	Fiat	b01
Group #3	3 Golf	GASOLINE	VW	e6c
	4 Punto	GASOLINE	Fiat	e6c
	5 500	GASOLINE	Fiat	e6c
	6 Corolla	GASOLINE	Toyota	e6c
	7 Astra	GASOLINE	Opel	e6c

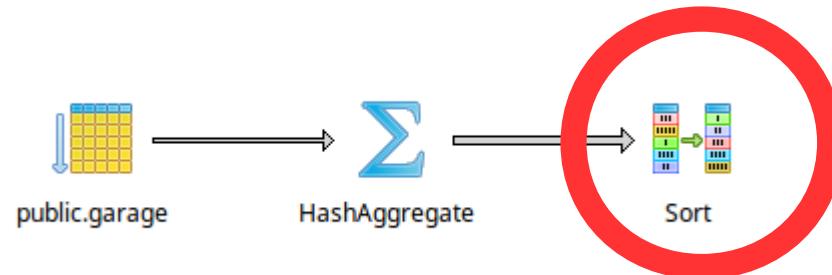


Result

	fuel_type character varying	count bigint
1	DIESEL	1
2	GASOLINE	5
3	METHANE	1



Executor stage



```
SELECT fuel_type, COUNT(*)  
FROM garage  
GROUP BY fuel_type  
ORDER BY count DESC;
```

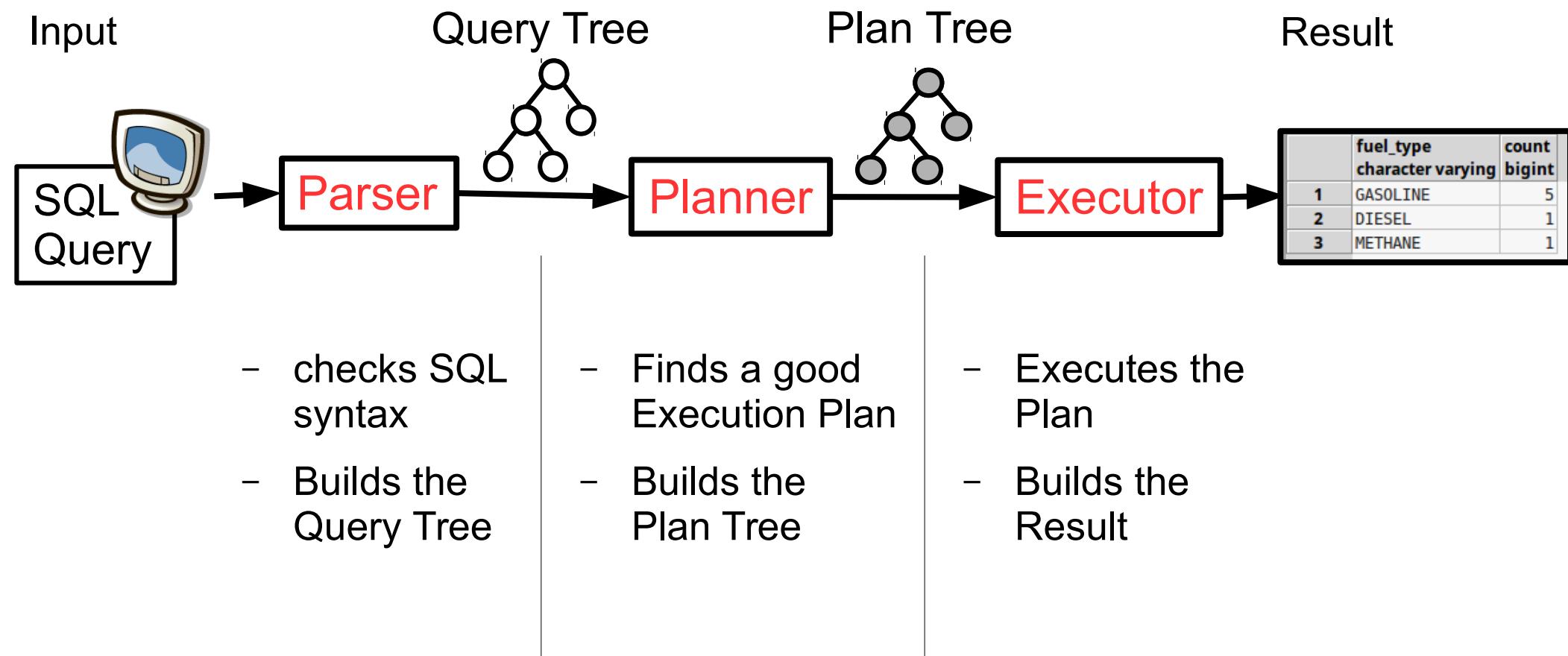
- ◆ Sort the rows from the subplan HashAggregate with sort key “count”
- ◆ Return the sorted groups as result

	fuel_type character varying	count bigint
1	DIESEL	1
2	GASOLINE	5
3	METHANE	1



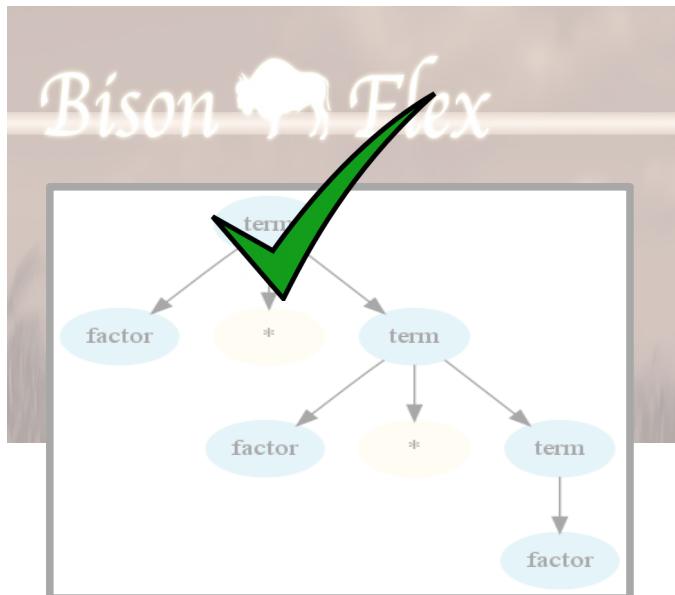
	fuel_type character varying	count bigint
1	GASOLINE	5
2	DIESEL	1
3	METHANE	1

In summary, the path of a query...

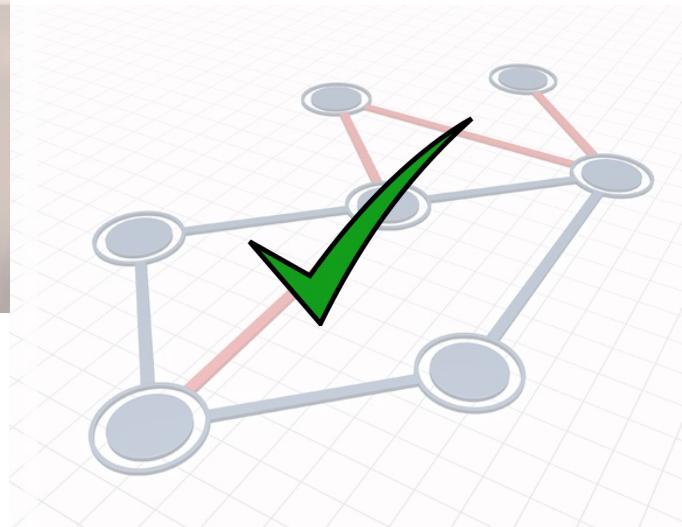


Three main stages of query processing

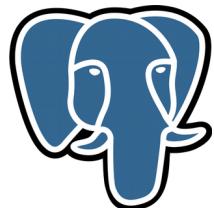
Parser stage



Planner stage



Executor stage



Thank you!



Additional information

Workshop Tutorial @ github

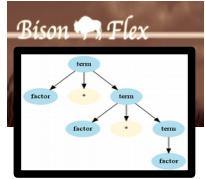
[https://github.com/Piit/postgres-twice-patch
/blob/master/SFSCON15-README.md](https://github.com/Piit/postgres-twice-patch/blob/master/SFSCON15-README.md)

...or...

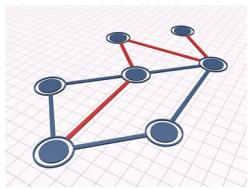
<http://bit.ly/10yU0tD>



In summary, the path of a query...



- **Parser**
 - checks SQL syntax
 - builds a query tree
- **Planner / Optimizer**
 - Finds an (not always) optimal plan
 - Tests all possibilities or uses heuristics
- **Executor**
 - Executes the plan
 - Is implemented as an iterator



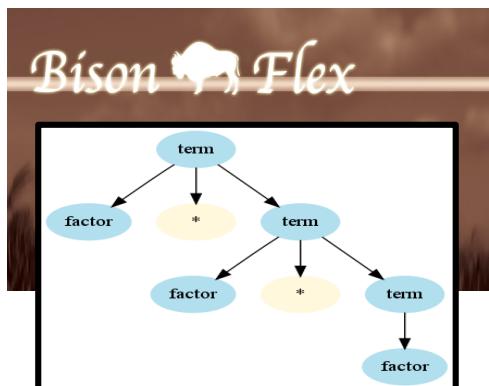
Three main stages of query processing

SQL Query

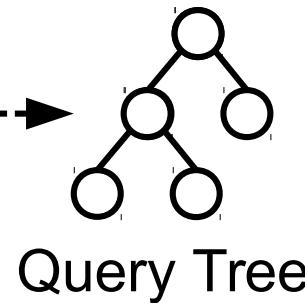
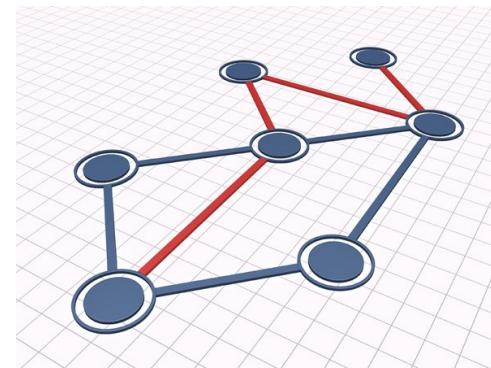
```
SELECT fuel_type, COUNT(*)  
FROM garage  
GROUP BY fuel_type;
```



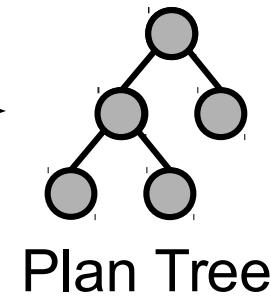
Parser stage



Planner stage



Query Tree



Plan Tree

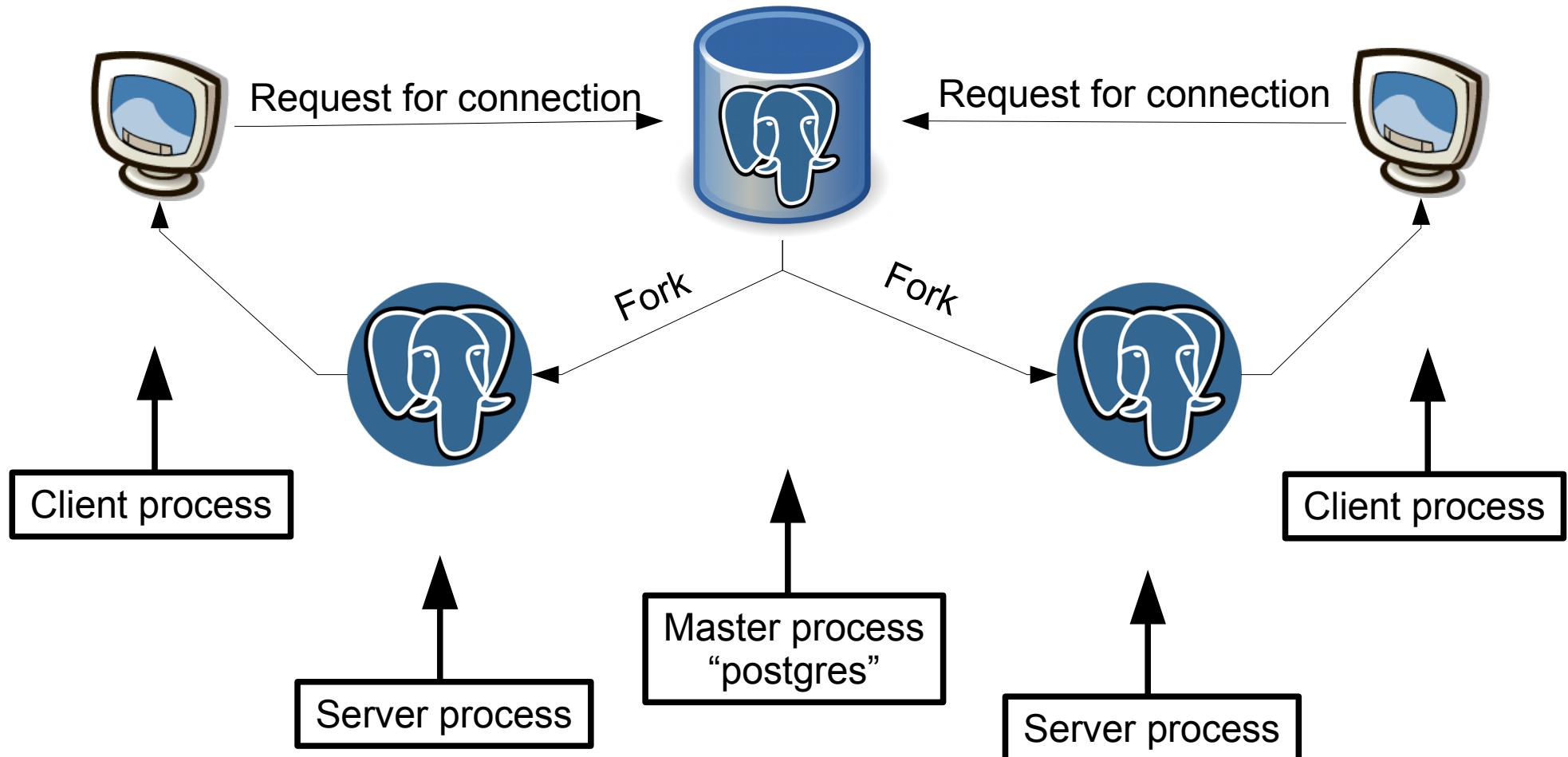
Result

	fuel_type	count
	character varying	bigint
1	DIESEL	1
2	GASOLINE	5
3	METHANE	1

Executor stage



PostgreSQL Server / Client interaction



```
{SELECT
:distinctClause <>
:intoClause <>
:targetList (
{RESTARTER
:name <>
:indirection <>
:val
{COLUMNREF
:fields ("car")
:location 7
}
)
{RESTARTER
:name <>
:indirection <>
:val
{FUNCCALL
:funcname ("count")
:args <>
:agg_star true
:over <>
:location 12
}
)
}
```

Simplified Parse Tree

```
:fromClause (
{RANGEVAR
:relname garage
:alias <>
:location 26
})
:whereClause <>
:groupClause (
{COLUMNREF
:fields ("car")
:location 42
})
:havingClause <>
>windowClause <>
:valuesLists <>
:sortClause <>
:limitCount <>
:lockingClause <>
:withClause <>
:op 0
:all false
:larg <>
:rarg <>
)
```

pretty-print
output of
pprint()