

Formal Semantics of SQL (and Cypher)

Paolo Guagliardo



THE UNIVERSITY *of* EDINBURGH
informatics

SQL

- **Standard** query language for relational databases
- **\$30B/year** business
- Implemented in all major RDBMSs (free and commercial)
- First standardized in **1986** (ANSI) and **1987** (ISO)
- Several revision afterwards
(SQL-89, SQL-92, SQL:1999, SQL:2003, SQL:2006, SQL:2008, SQL:2011, SQL:2016)

“The nice thing about standards is that you have so many to choose from” — Andrew S. Tanenbaum

How standard is SQL?

```
SELECT * FROM ( SELECT R.A, R.A FROM R ) S
```

PostgreSQL outputs a table with two columns named "A"

Oracle throws an **ERROR**: reference to column "A" is ambiguous

```
SELECT * FROM R WHERE EXISTS (  
    SELECT * FROM (  
        SELECT R.A, R.A FROM R ) S )
```

Both **PostgreSQL** and **Oracle** output R

Who is right?

Let's have a look at the standard!

- A. If the `<select list> *` is *simply contained* in a `<subquery>` that is *immediately contained* in an `<exists predicate>`, then the `<select list>` is **equivalent** to a `<value expression>` that is an arbitrary `<literal>`.
- B. Otherwise, the `<select list> *` is **equivalent** to a `<value expression>` sequence in which each `<value expression>` is a column reference that references a column of *T* and each column of *T* is referenced **exactly once**. The columns are referenced in the ascending sequence of their ordinal position within *T*.

... which means

```
SELECT *  
FROM ( SELECT R.A, R.A  
        FROM R ) S
```

≡

```
SELECT S.A, S.A  
FROM ( SELECT R.A, R.A  
        FROM R ) S
```

```
SELECT * FROM R  
WHERE EXISTS (  
    SELECT *  
    FROM ( SELECT R.A, R.A  
            FROM R ) S  
)
```

≡

```
SELECT R.A FROM R  
WHERE EXISTS (  
    SELECT 1  
    FROM ( SELECT R.A, R.A  
            FROM R ) S  
)
```

The Need for a Formal Semantics

- Avoid ambiguity of natural language
- Clearly defined and not subject to interpretation
- Easy to understand and implement

Previous attempts

- Many simplifying assumptions: no **bags**, no **nulls**
- No justification of correctness

| | | | |
|---|------|---|------|
| R | A | S | A |
| | 1 | | NULL |
| | NULL | | |

Answer

```
SELECT R.A FROM R
EXCEPT
SELECT S.A FROM S
```

| |
|---|
| A |
| 1 |

```
SELECT R.A FROM R
WHERE R.A NOT IN (
  SELECT S.A FROM S)
```

| |
|---|
| A |
| |

```
SELECT R.A FROM R
WHERE NOT EXISTS (
  SELECT S.A FROM S
  WHERE S.A=R.A )
```

| |
|------|
| A |
| 1 |
| NULL |

Core SQL fragment

$\tau := (T_1, \dots, T_k), \quad \beta := (N_1, \dots, N_k), \quad k > 0$

$\alpha := (A_1, \dots, A_m), \quad \beta' := (N'_1, \dots, N'_m), \quad m > 0$

QUERIES:

$Q := \text{SELECT } [\text{DISTINCT}] (\alpha:\beta' \mid *) \text{ FROM } \tau:\beta \text{ WHERE } \theta$
 $\mid Q (\text{UNION} \mid \text{INTERSECT} \mid \text{EXCEPT}) [\text{ALL}] Q$

CONDITIONS:

$\theta := \text{TRUE} \mid \bar{t} (= \mid \neq) \bar{t} \mid t \text{ IS } [\text{NOT}] \text{ NULL}$
 $\mid \bar{t} [\text{NOT}] \text{ IN } Q \mid \text{ EXISTS } Q$
 $\mid \theta \text{ AND } \theta \mid \theta \text{ OR } \theta \mid \text{ NOT } \theta$

Essentially SQL without arithmetic, grouping and aggregation

Formal Semantics: Challenges

Data model

- Base relations / query outputs / intermediate results
- Primitive data manipulation operations

Attribute references

- Binding rules in subqueries
- *Environment* collects and propagates bindings

Proposed Semantics

$$\begin{aligned}
 \llbracket R \rrbracket_{D,\eta} &= R^D \\
 \llbracket \tau : \beta \rrbracket_{D,\eta} &= \llbracket (T_1, \dots, T_k) : (N_1, \dots, N_k) \rrbracket_{D,\eta} = N_1 \cdot \llbracket T_1 \rrbracket_{D,\eta} \times \dots \times N_k \cdot \llbracket T_k \rrbracket_{D,\eta} \\
 \llbracket \text{FROM } \tau : \beta \rrbracket_{D,\eta} &= \{ \bar{a} \in \llbracket \tau : \beta \rrbracket_{D,\eta} \mid \llbracket \theta \rrbracket_{D,\eta;\eta^{\bar{a}}} = \mathbf{t} \} \\
 \llbracket \text{SELECT } * \rrbracket_{D,\eta} &= \llbracket \text{FROM } \tau : \beta \rrbracket_{D,\eta} : \beta^{-1} \\
 \llbracket \text{SELECT } \alpha : \beta' \rrbracket_{D,\eta} &= \pi_\alpha \left(\llbracket \text{FROM } \tau : \beta \rrbracket_{D,\eta} \right) : \beta' \\
 \llbracket \text{SELECT DISTINCT } \alpha : \beta' \rrbracket_{D,\eta} &= \varepsilon \left(\llbracket \text{SELECT } \alpha : \beta' \rrbracket_{D,\eta} \right) \\
 \llbracket \text{TRUE} \rrbracket_{D,\eta} &= \mathbf{t} \\
 \llbracket t \rrbracket_{D,\eta} &= \begin{cases} \eta(A) & \text{if } t = A \\ t & \text{if } t \in \mathbf{C} \text{ or } t = \text{NULL} \end{cases} \\
 \llbracket t_1 = t_2 \rrbracket_{D,\eta} &= \begin{cases} \mathbf{t} & \text{if } \llbracket t_1 \rrbracket_{D,\eta} = \llbracket t_2 \rrbracket_{D,\eta} \text{ and } \llbracket t_1 \rrbracket_{D,\eta} \neq \text{NULL} \text{ and } \llbracket t_2 \rrbracket_{D,\eta} \neq \text{NULL} \\ \mathbf{f} & \text{if } \llbracket t_1 \rrbracket_{D,\eta} \neq \llbracket t_2 \rrbracket_{D,\eta} \text{ and } \llbracket t_1 \rrbracket_{D,\eta} \neq \text{NULL} \text{ and } \llbracket t_2 \rrbracket_{D,\eta} \neq \text{NULL} \\ \mathbf{u} & \text{if } \llbracket t_1 \rrbracket_{D,\eta} = \text{NULL} \text{ or } \llbracket t_2 \rrbracket_{D,\eta} = \text{NULL} \end{cases} \\
 \llbracket t \text{ IS NULL} \rrbracket_{D,\eta} &= \begin{cases} \mathbf{t} & \text{if } \llbracket t \rrbracket_{D,\eta} = \text{NULL} \\ \mathbf{f} & \text{if } \llbracket t \rrbracket_{D,\eta} \neq \text{NULL} \end{cases} \\
 \llbracket t \text{ IS NOT NULL} \rrbracket_{D,\eta} &= \neg \llbracket t \text{ IS NULL} \rrbracket_{D,\eta} \\
 \llbracket (t_1, \dots, t_n) = (t'_1, \dots, t'_n) \rrbracket_{D,\eta} &= \bigwedge_{i=1}^n \llbracket t_i = t'_i \rrbracket_{D,\eta} \quad \llbracket (t_1, \dots, t_n) \neq (t'_1, \dots, t'_n) \rrbracket_{D,\eta} = \bigvee_{i=1}^n \llbracket t_i \neq t'_i \rrbracket_{D,\eta} \\
 \llbracket \bar{t} \text{ IN } Q \rrbracket_{D,\eta} &= \begin{cases} \mathbf{t} & \text{if } \exists \bar{r} \in [Q]_{D,\eta} : \llbracket \bar{t} = \nu(\bar{r}) \rrbracket_{D,\eta} = \mathbf{t} \\ \mathbf{f} & \text{if } \forall \bar{r} \in [Q]_{D,\eta} : \llbracket \bar{t} = \nu(\bar{r}) \rrbracket_{D,\eta} = \mathbf{f} \\ \mathbf{u} & \text{if } \nexists \bar{r} \in [Q]_{D,\eta} : \llbracket \bar{t} = \nu(\bar{r}) \rrbracket_{D,\eta} = \mathbf{t} \text{ and } \exists \bar{r} \in [Q]_{D,\eta} : \llbracket \bar{t} = \nu(\bar{r}) \rrbracket_{D,\eta} \neq \mathbf{f} \end{cases} \\
 \llbracket \bar{t} \text{ NOT IN } Q \rrbracket_{D,\eta} &= \neg \llbracket \bar{t} \text{ IN } Q \rrbracket_{D,\eta} \\
 \llbracket \text{EXISTS } Q \rrbracket_{D,\eta} &= \begin{cases} \mathbf{t} & \text{if } [Q]_{D,\eta} \neq \emptyset \\ \mathbf{f} & \text{if } [Q]_{D,\eta} = \emptyset \end{cases} \\
 \llbracket \theta_1 \text{ AND } \theta_2 \rrbracket_{D,\eta} &= \llbracket \theta_1 \rrbracket_{D,\eta} \wedge \llbracket \theta_2 \rrbracket_{D,\eta} \quad \llbracket \theta_1 \text{ OR } \theta_2 \rrbracket_{D,\eta} = \llbracket \theta_1 \rrbracket_{D,\eta} \vee \llbracket \theta_2 \rrbracket_{D,\eta} \quad \llbracket \text{NOT } \theta \rrbracket_{D,\eta} = \neg \llbracket \theta \rrbracket_{D,\eta} \\
 \llbracket Q_1 \text{ UNION ALL } Q_2 \rrbracket_{D,\eta} &= [Q_1]_{D,\eta} \cup [Q_2]_{D,\eta} : \ell([Q_1]) \\
 \llbracket Q_1 \text{ INTERSECT ALL } Q_2 \rrbracket_{D,\eta} &= [Q_1]_{D,\eta} \cap [Q_2]_{D,\eta} : \ell([Q_1]) \\
 \llbracket Q_1 \text{ EXCEPT ALL } Q_2 \rrbracket_{D,\eta} &= [Q_1]_{D,\eta} - [Q_2]_{D,\eta} : \ell([Q_1]) \\
 \llbracket Q_1 \star Q_2 \rrbracket_{D,\eta} &= \varepsilon([Q_1 \star \text{ALL } Q_2]_{D,\eta}), \quad \star \in \{\text{UNION, INTERSECT}\} \\
 \llbracket Q_1 \text{ EXCEPT } Q_2 \rrbracket_{D,\eta} &= \varepsilon([Q_1]_{D,\eta}) - [Q_2]_{D,\eta} : \ell([Q_1])
 \end{aligned}$$

- Fits in one page
- Non-ambiguous
- Easy to understand
- Easy to implement
- Easy to modify

Formal Semantics: Validation

- Cannot **prove** that semantics is correct
- Provide sufficient **experimental evidence**
- Implemented in Python
- Validated on 100000+ random SQL queries

Formal Semantics of Cypher

- Collaboration between **Neo Technology** and the **University of Edinburgh**
- Preliminary meeting in December
- Legal agreements finalized recently
- Neo Technology sponsors a researcher (**Nadime Francis**)

Challenges

- Getting the (abstract) data model right
- Intermediate representation (QUIL?)
- Identify core fragment
- Language constantly evolving
- Follow the footsteps of SQL? (nulls)