

Abstract

To do : A lot :) (SAT/SMT SOLVER ???) to evaluate what we loose ?

1 Preliminaries

Definition 1. We denote the Set of well formed select query without agregation, full join and null test by $\llbracket SQL \rrbracket$

We denote the Set of well formed select query without agregation and full join by $\llbracket SQL \rrbracket_{\perp}$

Definition 2. Let's a Select query $Q \in \llbracket SQL \rrbracket$ a tuple (Σ, R, H, P) such that

$$\Sigma \subseteq \{r_i.a_i | \exists r_i \in R \wedge \exists a_i \in r_i\}$$

$$P \subseteq \{p_i = r_j.a_j | r_j \notin R\} \text{ a set of external parameter.}$$

R a set of relation.

H belongs to the following grammar

$$\begin{aligned} H ::= & r_i.a_i = c_i \mid r_i.a_i \neq c_j \mid r_i.a_i = r_j.a_j \mid r_i.a_i \neq r_j.a_j \mid r_i.a_i = p_i \mid \\ & exists(Q) \mid notexists(Q) \mid in(r_i.a_i, Q) \mid notin(r_i.a_i, Q) \mid \\ & H \wedge H \mid H \vee H \end{aligned}$$

Definition 3. Let's a Select query $Q \in \llbracket SQL \rrbracket_{\perp}$ a tuple (Σ, R, H, P) such that

$$\Sigma \subseteq \{r_i.a_i | \exists r_i \in R \wedge \exists a_i \in r_i\}$$

$$P \subseteq \{p_i = r_j.a_j | r_j \notin R\} \text{ a set of external parameter.}$$

R a set of relation.

H_{\perp} belongs to the following grammar

$$\begin{aligned}
H_{\perp} ::= & r_i.a_i = c_i \mid r_i.a_i \neq c_j \mid r_i.a_i = r_j.a_j \mid r_i.a_i \neq r_j.a_j \mid \\
& r_i.a_i = p_i \mid \text{null}(r_i.a_i) \mid \text{const}(r_i.a_i) \\
& \text{exists}(Q_{\perp}) \mid \text{notexists}(Q_{\perp}) \mid \text{in}(r_i.a_i, Q_{\perp}) \mid \text{notin}(r_i.a_i, Q_{\perp}) \mid \\
& H_{\perp} \wedge H_{\perp} \mid H_{\perp} \vee H_{\perp}
\end{aligned}$$

We denote $(\Sigma, R, H, P)[x]$ the query $(\Sigma, R, H, P \cup x)$
 We denote $(\Sigma, R, H, P)_*$ the query $(*, R, H, P)$

Proposition 1.

$$\llbracket SQL \rrbracket \subset \llbracket SQL \rrbracket_{\perp}$$

Definition 4. We call a bag B a function $D \rightarrow \mathbb{N}$ such that $B(x)$ represents the multiplicity of x in the bag B .

Definition 5.

$$\begin{aligned}
& \forall x, \emptyset(x) = 0 \\
& \forall x, (B_1 \cap B_2)(x) = \min(B_1(x), B_2(x)) \\
& \forall x, (B_1 \cup B_2)(x) = \max(B_1(x), B_2(x)) \\
& \forall x, (B_1 \uplus B_2)(x) = B_1(x) + B_2(x) \\
& \forall x, (B_1 \setminus B_2)(x) = \max(0, B_1(x) - B_2(x)) \\
& \forall x, \llbracket a \rrbracket(x) = \begin{cases} 1 & \text{if } x = a \\ 0 & \text{otherwise} \end{cases} \\
& \forall x, \llbracket a^n \rrbracket(x) = \begin{cases} n & \text{if } x = a \\ 0 & \text{otherwise} \end{cases} \\
& \forall x, \llbracket y^n \mid P(y, n) \rrbracket(x) = \max(\{i \mid P(x, i)\}) \\
& x \in B \iff B(x) \geq 1 \\
& x \in^n B \iff B(x) \geq n \\
& x \notin B \iff B(x) = 0 \\
& B_1 = B_2 \iff \forall x, B_1(x) = B_2(x) \\
& B_1 \subseteq B_2 \iff \forall x, B_1(x) \leq B_2(x) \\
& \{B\} = \{x \mid B(x) \geq 1\}
\end{aligned}$$

2 Semantics

Definition 6. Let Q be a query (Σ, R, H, P) then R^{\times} denotethe bag coming from the cartesian product of R .

$$R^{\times} = \prod_{r \in R} r$$

Definition 7.

$$\begin{aligned}\sigma_\Sigma(x) &= (x[r_i.a_i] | r_i.a_i \in \Sigma) \\ \sigma_*(x) &= x\end{aligned}$$

Definition 8.

$$\sigma_\Sigma(B) = \llbracket y^n | n = \sum_{x \in \{z | z \in \{B\} \wedge \sigma_\Sigma(z) = y\}} B(x) \rrbracket$$

Definition 9.

$$\begin{aligned}Eval_{SQL}((\Sigma, R, H_1 \wedge H_2, P), D) &= \sigma_\Sigma(Eval_{SQL}((*, R, H_1, P), D) \cap Eval_{SQL}((*, R, H_2, P), D)) \\ Eval_{SQL}((\Sigma, R, H_1 \vee H_2, P), D) &= \sigma_\Sigma(Eval_{SQL}((*, R, H_1, P), D) \cup Eval_{SQL}((*, R, H_2, P), D)) \\ Eval_{SQL}((\Sigma, R, null(r_i.a_i), P), D) &= \sigma_\Sigma(\llbracket x^n | R^\times(x) = n \wedge x[r_i.a_i] = \perp \rrbracket) \\ Eval_{SQL}((\Sigma, R, const(r_i.a_i), P), D) &= \sigma_\Sigma(\llbracket x^n | R^\times(x) = n \wedge x[r_i.a_i] \neq \perp \rrbracket) \\ Eval_{SQL}((\Sigma, R, r_i.a_i = c_i, P), D) &= \sigma_\Sigma(\llbracket x^n | R^\times(x) = n \wedge x[r_i.a_i] = c_i \rrbracket) \\ Eval_{SQL}((\Sigma, R, r_i.a_i = r_j.a_j, P), D) &= \sigma_\Sigma(\llbracket x^n | R^\times(x) = n \wedge x[r_i.a_i] = x[r_j.a_j] \wedge x[r_i.a_i] \neq \perp \wedge x[r_j.a_j] \neq \perp \rrbracket) \\ Eval_{SQL}((\Sigma, R, r_i.a_i = p_i, P), D) &= \sigma_\Sigma(\llbracket x^n | R^\times(x) = n \wedge x[r_i.a_i] = P[p_i] \wedge x[r_i.a_i] \neq \perp \wedge P[p_i] \neq \perp \rrbracket) \\ Eval_{SQL}((\Sigma, R, r_i.a_i \neq c_i, P), D) &= \sigma_\Sigma(\llbracket x^n | R^\times(x) = n \wedge x[r_i.a_i] \neq c_i \wedge x[r_i.a_i] \neq \perp \rrbracket) \\ Eval_{SQL}((\Sigma, R, r_i.a_i \neq r_j.a_j, P), D) &= \sigma_\Sigma(\llbracket x^n | R^\times(x) = n \wedge x[r_i.a_i] \neq x[r_j.a_j] \wedge x[r_i.a_i] \neq \perp \wedge x[r_j.a_j] \neq \perp \rrbracket) \\ Eval_{SQL}((\Sigma, R, r_i.a_i \neq p_i, P), D) &= \sigma_\Sigma(\llbracket x^n | R^\times(x) = n \wedge x[r_i.a_i] \neq P[p_i] \wedge x[r_i.a_i] \neq \perp \wedge P[p_i] \neq \perp \rrbracket) \\ Eval_{SQL}((\Sigma, R, exists(Q), P), D) &= \sigma_\Sigma(\llbracket x^n | R^\times(x) = n \wedge Eval_{SQL}(Q[x], D) \neq \emptyset \rrbracket) \\ Eval_{SQL}((\Sigma, R, notexists(Q), P), D) &= \sigma_\Sigma(\llbracket x^n | R^\times(x) = n \wedge Eval_{SQL}(Q[x], D) = \emptyset \rrbracket) \\ Eval_{SQL}((\Sigma, R, in(\Sigma_1, Q), P), D) &= \sigma_\Sigma(\llbracket x^n | R^\times(x) = n \\ &\quad \wedge \forall y \in Eval_{SQL}(Q, D), \forall i \in \llbracket 1; |\Sigma_1| \rrbracket, \sigma_{\Sigma_1}(x)[i] = y[i] \\ &\quad \wedge \sigma_{\Sigma_1}(x)[i] \neq \perp \wedge y[i] \neq \perp \rrbracket) \\ Eval_{SQL}((\Sigma, R, notin(\Sigma_1, Q), P), D) &= \sigma_\Sigma(\llbracket x^n | R(x) = n \\ &\quad \wedge \forall y \in Eval_{SQL}(Q, D), \forall i \in \llbracket 1; |\Sigma_1| \rrbracket, \sigma_{\Sigma_1}(x)[i] \neq y[i] \\ &\quad \vee \sigma_{\Sigma_1}(x)[i] = \perp \vee y[i] = \perp \rrbracket)\end{aligned}$$

Definition 10. *With marked nulls we have:*

$$cert_\perp(Q, D) = \llbracket x^n | \forall h, h(x) \in {}^n Eval_{SQL}(Q, h(D)) \rrbracket$$

Proposition 2. *With SQL nulls we have:*

$$\forall Q \in \llbracket SQL \rrbracket, cert_\perp((\Sigma, R, H, P), D) = \sigma_\Sigma(\llbracket x^n | R(x) = n \wedge \forall h, h(x) \in \{Eval_{SQL}((*, R, H, h(P)), h(D))\} \rrbracket)$$

Proof. Dont even know where to begin. \square

Definition 11.

$$\forall Q \in \llbracket SQL \rrbracket, posi_\perp((\Sigma, R, H, P), D) = \sigma_\Sigma(\llbracket x^n | R(x) = n \wedge \exists h, h(x) \in \{Eval_{SQL}((*, R, H, h(P)), h(D))\} \rrbracket)$$

3 Translation

$$\begin{aligned}(\Sigma, R, H, P)^+ &\rightarrow (\Sigma, R, H^*, P) \\ (\Sigma, R, H, P)^? &\rightarrow (\Sigma, R, H^{**}, P)\end{aligned}$$

$$\begin{aligned}(H_1 \wedge H_2)^* &\rightarrow H_1^* \wedge H_2^* \\ (H_1 \vee H_2)^* &\rightarrow H_1^* \vee H_2^* \\ (r_i.a_i = c_i)^* &\rightarrow r_i.a_i = c_i \\ (r_i.a_i \neq c_i)^* &\rightarrow r_i.a_i \neq c_i \\ (r_i.a_i = r_j.a_j)^* &\rightarrow r_i.a_i = r_j.a_j \\ (r_i.a_i \neq r_j.a_j)^* &\rightarrow r_i.a_i \neq r_j.a_j \\ \text{null}(r_i.a_i)^* &\rightarrow \text{null}(r_i.a_i) \\ \text{const}(r_i.a_i)^* &\rightarrow \text{const}(r_i.a_i) \\ \text{exists}(Q)^* &\rightarrow \text{exists}(Q^+) \\ \text{notexists}(Q)^* &\rightarrow \text{notexists}(Q^?) \\ \text{in}(\Sigma_1, Q)^* &\rightarrow \text{in}(\Sigma, Q^+) \\ \text{notin}(\Sigma_1, (\Sigma, R, H, P))^* &\rightarrow \text{notexists}(\Sigma, R, H \wedge (\Sigma = \Sigma_1), P \cup \Sigma_1)^*\end{aligned}$$

$$\begin{aligned}(H_1 \wedge H_2)^{**} &\rightarrow H_1^{**} \wedge H_2^{**} \\ (H_1 \vee H_2)^{**} &\rightarrow H_1^{**} \vee H_2^{**} \\ (r_i.a_i = c_i)^{**} &\rightarrow r_i.a_i = c_i \vee \text{null}(r_i.a_i) \\ (r_i.a_i \neq c_i)^{**} &\rightarrow r_i.a_i \neq c_i \vee \text{null}(r_i.a_i) \\ (r_i.a_i = r_j.a_j)^{**} &\rightarrow r_i.a_i = r_j.a_j \vee \text{null}(r_i.a_i) \vee \text{null}(r_j.a_j) \\ (r_i.a_i \neq r_j.a_j)^{**} &\rightarrow r_i.a_i \neq r_j.a_j \vee \text{null}(r_i.a_i) \vee \text{null}(r_j.a_j) \\ \text{null}(r_i.a_i)^{**} &\rightarrow \text{null}(r_i.a_i) \\ \text{const}(r_i.a_i)^{**} &\rightarrow \text{const}(r_i.a_i) \\ \text{exists}(Q)^{**} &\rightarrow \text{exists}(Q^?) \\ \text{notexists}(Q)^{**} &\rightarrow \text{notexists}(Q^+) \\ \text{in}(r_i.a_i, Q)^{**} &\rightarrow \text{in}(r_i.a_i, Q^?) \\ \text{notin}(\Sigma_1, (\Sigma, R, H, P))^{**} &\rightarrow \text{notexists}(\Sigma, R, H \wedge (\Sigma = \Sigma_1), P \cup \Sigma_1)^{**}\end{aligned}$$

Proposition 3.

$$\forall Q \in \llbracket SQL \rrbracket, \text{Eval}_{SQL}(Q^+, D) \subseteq \text{cert}_\perp(Q, D)$$

Proposition 4.

$$\forall Q \in \llbracket SQL \rrbracket, \text{posi}_\perp(Q, D) \subseteq \text{Eval}_{SQL}(Q^?, D)$$

Proof. Assume (5).

By induction :

$$Q = (\Sigma, R, H_1 \wedge H_2, P)$$

$$x \in^n Eval_{SQL}(Q^+, D) \Rightarrow \sum_{z \in \{y | y \in Eval_{SQL}(Q^+, D) \wedge \sigma_\Sigma(y) = x\}} Eval_{SQL}(Q^+, D)(z) \geq n$$

Moreover

$$\begin{aligned} Eval_{SQL}(Q^+, D)(z) = k &\Rightarrow Eval_{SQL}((*, R, H_1^* \wedge H_2^*, P), D)(z) = k \\ &\Rightarrow (Eval_{SQL}((*, R, H_1^*, P), D) \cap Eval_{SQL}((*, R, H_2^*, P), D))(z) = k \\ &\Rightarrow (Eval_{SQL}((*, R, H_1, P)^+, D) \cap Eval_{SQL}((*, R, H_2, P)^+, D))(z) = k \\ &\Rightarrow (Eval_{SQL}((*, R, H_1, P)^+, D))(x) \geq k \wedge (Eval_{SQL}((*, R, H_2, P)^+, D))(z) \geq k \\ &\Rightarrow (cert_\perp((*, R, H_1, P), D))(x) \geq k \wedge (cert_\perp((*, R, H_2, P), D))(z) \geq k \\ &\Rightarrow (cert_\perp((*, R, H_1, P), D) \cap cert_\perp((*, R, H_2, P), D))(z) \geq k \\ &\Rightarrow (cert_\perp((*, R, H_1 \wedge H_2, P), D))(z) \geq k \\ &\Rightarrow (cert_\perp(Q_*, D))(z) \geq k \end{aligned}$$

Then

$$\begin{aligned} x \in^n Eval_{SQL}(Q^+, D) &\Rightarrow \sum_{z \in \{y | y \in Eval_{SQL}(Q^+, D) \wedge \sigma_\Sigma(y) = x\}} cert_\perp(Q_*, D)(z) \geq n \\ &\Rightarrow x \in^n cert_\perp(Q, D) \end{aligned}$$

$$Q = (\Sigma, R, H_1 \vee H_2, P)$$

$$x \in^n Eval_{SQL}(Q^+, D) \Rightarrow \sum_{z \in \{y | y \in Eval_{SQL}(Q_*^+, D) \wedge \sigma_\Sigma(y) = x\}} Eval_{SQL}(Q_*^+, D)(z) \geq n$$

Moreover

$$\begin{aligned} Eval_{SQL}(Q_*^+, D)(z) = k &\Rightarrow Eval_{SQL}((*, R, H_1^* \vee H_2^*, P), D)(z) = k \\ &\Rightarrow (Eval_{SQL}((*, R, H_1^*, P), D) \cup Eval_{SQL}((*, R, H_2^*, P), D))(z) = k \\ &\Rightarrow (Eval_{SQL}((*, R, H_1, P)^+, D) \cup Eval_{SQL}((*, R, H_2, P)^+, D))(z) = k \\ &\Rightarrow (Eval_{SQL}((*, R, H_1, P)^+, D))(x) \geq k \vee (Eval_{SQL}((*, R, H_2, P)^+, D))(z) \geq k \\ &\Rightarrow (cert_\perp((*, R, H_1, P), D))(x) \geq k \vee (cert_\perp((*, R, H_2, P), D))(z) \geq k \\ &\Rightarrow (cert_\perp((*, R, H_1, P), D) \cup cert_\perp((*, R, H_2, P), D))(z) \geq k \\ &\Rightarrow (cert_\perp((*, R, H_1 \vee H_2, P), D))(z) \geq k \\ &\Rightarrow (cert_\perp(Q_*, D))(z) \geq k \end{aligned}$$

Then

$$\begin{aligned} x \in^n Eval_{SQL}(Q^+, D) &\Rightarrow \sum_{z \in \{y | y \in Eval_{SQL}(Q_*^+, D) \wedge \sigma_\Sigma(y) = x\}} cert_\perp(Q_*, D)(z) \geq n \\ &\Rightarrow x \in^n cert_\perp(Q, D) \end{aligned}$$

$$Q = (\Sigma, R, notexists(Q'), P)$$

$$x \in^n Eval_{SQL}(Q^+, D) \Rightarrow \sum_{z \in \{y | y \in Eval_{SQL}(Q_*^+, D) \wedge \sigma_\Sigma(y) = x\}} Eval_{SQL}(Q_*^+, D)(z) \geq n$$

Moreover

$$\begin{aligned} Eval_{SQL}(Q_*^+, D)(z) = k &\Rightarrow Eval_{SQL}((*, R, notexists(Q')^*, P), D)(z) = k \\ &\Rightarrow Eval_{SQL}((*, R, notexists(Q'^?), P), D)(z) = k \\ &\Rightarrow R(z) = k \wedge Eval_{SQL}(Q'[z]^?, D) = \emptyset \\ &\Rightarrow R(z) = k \wedge posi_\perp(Q'[z], D) = \emptyset \\ &\Rightarrow R(z) = k \wedge \forall h, \forall w, h(w) \notin Eval_{SQL}(Q'[h(z)], h(D)) \\ &\Rightarrow R(z) = k \wedge \forall h, Eval_{SQL}(Q'[h(z)], h(D)) = \emptyset \\ &\Rightarrow R(z) = k \wedge \forall h, h(z) \in Eval_{SQL}((*, R, notexists(Q'), h(P)), h(D)) \\ &\Rightarrow cert_\perp(Q_*, D)(z) = k \end{aligned}$$

Then

$$\begin{aligned} x \in^n Eval_{SQL}(Q^+, D) &\Rightarrow \sum_{z \in \{y | y \in Eval_{SQL}(Q_*^+, D) \wedge \sigma_\Sigma(y) = x\}} cert_\perp(Q_*, D)(z) \geq n \\ &\Rightarrow x \in^n cert_\perp(Q, D) \end{aligned}$$

$$Q = (\Sigma, R, r_i.a_i \neq p_i, P)$$

$$x \in^n Eval_{SQL}(Q^+, D) \Rightarrow \sum_{z \in \{y | y \in Eval_{SQL}(Q_*^+, D) \wedge \sigma_\Sigma(y) = x\}} Eval_{SQL}(Q_*^+, D)(z) \geq n$$

Moreover

$$\begin{aligned} Eval_{SQL}(Q_*^+, D)(z) = k &\Rightarrow Eval_{SQL}((*, R, (r_i.a_i \neq p_i)^*, P), D)(z) = k \\ &\Rightarrow Eval_{SQL}((*, R, r_i.a_i \neq p_i, P), D)(z) = k \\ &\Rightarrow R(x) = k \wedge x[r_i.a_i] \neq P[p_i] \wedge x[r_i.a_i] \neq \perp \wedge P[p_i] \neq \perp \end{aligned}$$

Moreover

$$x[r_i.a_i] \neq \perp \wedge P[p_i] \neq \perp \Rightarrow \forall h, h(x)[r_i.a_i] = x[r_i.a_i] \wedge h(P)[p_i] = P[p_i]$$

Then

$$\begin{aligned} Eval_{SQL}(Q_*^+, D)(z) = k &\Rightarrow R(x) = k \wedge \forall h, h(x)[r_i.a_i] \neq h(P)[p_i] \wedge h(x)[r_i.a_i] \neq \perp \wedge h(P)[p_i] \neq \perp \\ &\Rightarrow cert_\perp(Q_*, D)(z) = k \end{aligned}$$

Then

$$\begin{aligned} x \in^n Eval_{SQL}(Q^+, D) &\Rightarrow \sum_{z \in \{y | y \in Eval_{SQL}(Q_*^+, D) \wedge \sigma_\Sigma(y) = x\}} cert_\perp(Q_*, D)(z) \geq n \\ &\Rightarrow x \in^n cert_\perp(Q, D) \end{aligned}$$

□

Proof. Assume (4).

By induction ...

$$Q = (\Sigma, R, H_1 \wedge H_2, P)$$

$$x \in^n \text{posi}_\perp(Q, D) \Rightarrow \sum_{z \in \{y \mid y \in \text{posi}_\perp(Q_*, D) \wedge \sigma_\Sigma(y) = x\}} \text{posi}_\perp(Q_*, D)(z) \geq n$$

Moreover

$$\begin{aligned} \text{posi}_\perp(Q_*, D)(z) = k &\Rightarrow R(z) = k \wedge \exists h, h(z) \in \text{Eval}_{SQL}((*, R, H_1 \wedge H_2, h(P)), h(D)) \\ &\Rightarrow R(z) = k \wedge \exists h, h(z) \in \text{Eval}_{SQL}((*, R, H_1, h(P)), h(D)) \wedge h(z) \in \text{Eval}_{SQL}((*, R, H_2, h(P)), h(D)) \\ &\Rightarrow \text{posi}_\perp((*, R, H_1, P), D)(z) \geq k \wedge \text{posi}_\perp((*, R, H_2, P), D)(z) \geq k \\ &\Rightarrow \text{Eval}_{SQL}((*, R, H_1, P)^\intercal, D)(z) \geq k \wedge \text{Eval}_{SQL}((*, R, H_2, P)^\intercal, D)(z) \geq k \\ &\Rightarrow \text{Eval}_{SQL}((*, R, H_1^{**}, P), D)(z) \geq k \wedge \text{Eval}_{SQL}((*, R, H_2^{**}, P), D)(z) \geq k \\ &\Rightarrow \text{Eval}_{SQL}((*, R, H_1^{**} \wedge H_2^{**}, P), D)(z) \geq k \\ &\Rightarrow \text{Eval}_{SQL}((*, R, H_1 \wedge H_2, P)^\intercal, D)(z) \geq k \\ &\Rightarrow \text{Eval}_{SQL}(Q_*^\intercal, D)(z) \geq k \end{aligned}$$

Then

$$\begin{aligned} x \in^n \text{posi}_\perp(Q, D) &\Rightarrow \sum_{z \in \{y \mid y \in \text{posi}_\perp(Q_*, D) \wedge \sigma_\Sigma(y) = x\}} \text{Eval}_{SQL}(Q_*^\intercal, D)(z) \geq n \\ &\Rightarrow x \in^n \text{Eval}_{SQL}(Q^\intercal, D) \end{aligned}$$

$$Q = (\Sigma, R, \text{notexists}(Q'), P)$$

$$x \in^n \text{posi}_\perp(Q, D) \Rightarrow \sum_{z \in \{y \mid y \in \text{posi}_\perp(Q_*, D) \wedge \sigma_\Sigma(y) = x\}} \text{posi}_\perp(Q_*, D)(z) \geq n$$

Moreover

$$\begin{aligned} \text{posi}_\perp(Q_*, D)(z) = k &\Rightarrow R(z) = k \wedge \exists h, h(z) \in \text{Eval}_{SQL}((*, R, \text{notexists}(Q'), h(P)), h(D)) \\ &\Rightarrow R(z) = k \wedge \exists h, \text{Eval}_{SQL}(Q'[h(z)], h(D)) = \emptyset \\ &\Rightarrow R(z) = k \wedge \llbracket w^n \mid \forall g, g(w) \in^n \text{Eval}_{SQL}((Q'[g(z)]), g(D)) \rrbracket = \emptyset \\ &\Rightarrow R(z) = k \wedge \text{cert}_\perp(Q'[z], D) = \emptyset \\ &\Rightarrow R(z) = k \wedge \text{Eval}(Q'[z]^+, D) = \emptyset \\ &\Rightarrow R(z) = k \wedge z \in \text{Eval}_{SQL}((*, R, \text{notexists}(Q'^+), P), D) \\ &\Rightarrow R(z) = k \wedge z \in \text{Eval}_{SQL}((*, R, \text{notexists}(Q')^{**}, P), D) \\ &\Rightarrow \text{Eval}_{SQL}((*, R, \text{notexists}(Q'), P)^?, D)(z) = k \end{aligned}$$

Then

$$\begin{aligned} x \in^n \text{posi}_\perp(Q, D) &\Rightarrow \sum_{z \in \{y \mid y \in \text{posi}_\perp(Q_*, D) \wedge \sigma_\Sigma(y) = x\}} \text{Eval}_{SQL}(Q_*^?, D)(z) \geq n \\ &\Rightarrow x \in^n \text{Eval}_{SQL}(Q_*^?, D) \end{aligned}$$

$$Q = (\Sigma, R, r_i.a_i = r_j.a_j, P)$$

$$x \in^n \text{posi}_\perp(Q, D) \Rightarrow \sum_{z \in \{y \mid y \in \text{posi}_\perp(Q_*, D) \wedge \sigma_\Sigma(y) = x\}} \text{posi}_\perp(Q_*, D)(z) \geq n$$

Moreover

$$\begin{aligned} \text{posi}_\perp(Q_*, D)(z) = k &\Rightarrow R(z) = k \wedge \exists h, h(z) \in \text{Eval}_{SQL}((*, R, r_i.a_i = r_j.a_j, h(P)), h(D)) \\ &\Rightarrow R(z) = k \wedge \exists h, h(z)[r_i.a_i] = h(z)[r_j.a_j] \wedge h(z)[r_i.a_i] \neq \perp \wedge h(z)[r_j.a_j] \neq \perp \end{aligned}$$

Moreover

$$\begin{aligned} \exists h, h(z)[r_i.a_i] = h(z)[r_j.a_j] &\Rightarrow z[r_i.a_i] = z[r_j.a_j] \vee z[r_i.a_i] = \perp \vee z[r_j.a_j] = \perp \\ \forall g, g(z)[r_i.a_i] \neq \perp &\Rightarrow \text{TRUE} \end{aligned}$$

Then

$$\begin{aligned} \text{posi}_\perp(Q_*, D)(z) = k &\Rightarrow R(z) = k \wedge z[r_i.a_i] = z[r_j.a_j] \vee z[r_i.a_i] = \perp \vee z[r_j.a_j] = \perp \\ &\Rightarrow R(z) = k \wedge z \in \text{Eval}_{SQL}((*, R, (r_i.a_i = r_j.a_j \vee \text{null}(r_i.a_i) \vee \text{null}(r_j.a_j)), P), D) \\ &\Rightarrow R(z) = k \wedge z \in \text{Eval}_{SQL}((*, R, (r_i.a_i = r_j.a_j)^{**}, P), D) \\ &\Rightarrow \text{Eval}_{SQL}((*, R, r_i.a_i = r_j.a_j, P)^?, D)(z) = k \end{aligned}$$

Then

$$\begin{aligned} x \in^n \text{posi}_\perp(Q, D) &\Rightarrow \sum_{z \in \{y \mid y \in \text{posi}_\perp(Q_*, D) \wedge \sigma_\Sigma(y) = x\}} \text{Eval}_{SQL}(Q_*^?, D)(z) \geq n \\ &\Rightarrow x \in^n \text{Eval}_{SQL}(Q^?, D) \end{aligned}$$

□

4 Removing useless null check

The translation $Q \rightarrow Q^+$ has an heavy cost has explained in the paper ..., in order to remove some useless null test we offer an optimization.

Definition 12. For each Query and nested Query we maintain a set of attributes that have to be not null in order for the query to be true resp. false we denote this set \perp_Q^T resp. \perp_Q^F .

Definition 13.

$$\begin{aligned}
\perp_{(\Sigma, R, H_1 \wedge H_2, P)}^T &= \perp_{(\Sigma, R, H_1, P)}^T \cup \perp_{(\Sigma, R, H_2, P)}^T \\
\perp_{(\Sigma, R, H_1 \vee H_2, P)}^T &= \perp_{(\Sigma, R, H_1, P)}^T \cap \perp_{(\Sigma, R, H_2, P)}^T \\
\perp_{(\Sigma, R, r_i.a_i = c_i, P)}^T &= \{r_i.a_i\} \\
\perp_{(\Sigma, R, r_i.a_i \neq c_i, P)}^T &= \{r_i.a_i\} \\
\perp_{(\Sigma, R, r_i.a_i = r_j.a_j, P)}^T &= \{r_i.a_i, r_j.a_j\} \\
\perp_{(\Sigma, R, r_i.a_i \neq r_j.a_j, P)}^T &= \{r_i.a_i, r_j.a_j\} \\
\perp_{(\Sigma, R, r_i.a_i = p_i, P)}^T &= \{r_i.a_i, p_i\} \\
\perp_{(\Sigma, R, r_i.a_i \neq p_i, P)}^T &= \{r_i.a_i, p_i\} \\
\perp_{(\Sigma, R, \text{null}(r_i.a_i), P)}^T &= \emptyset \\
\perp_{(\Sigma, R, \text{const}(r_i.a_i), P)}^T &= \{r_i.a_i\} \\
\perp_{(\Sigma, R, \text{exists}(Q), P)}^T &= \perp_Q^T[?] \\
\perp_{(\Sigma, R, \text{notexists}(Q), P)}^T &= \perp_Q^F[?]
\end{aligned}$$

$$\begin{aligned}
\perp_{(\Sigma, R, H_1 \wedge H_2, P)}^F &= \perp_{(\Sigma, R, H_1, P)}^F \cap \perp_{(\Sigma, R, H_2, P)}^F \\
\perp_{(\Sigma, R, H_1 \vee H_2, P)}^F &= \perp_{(\Sigma, R, H_1, P)}^F \cup \perp_{(\Sigma, R, H_2, P)}^F \\
\perp_{(\Sigma, R, r_i.a_i = c_i, P)}^F &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i \neq c_i, P)}^F &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i = r_j.a_j, P)}^F &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i \neq r_j.a_j, P)}^F &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i = p_i, P)}^F &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i \neq p_i, P)}^F &= \emptyset \\
\perp_{(\Sigma, R, \text{null}(r_i.a_i), P)}^F &= \emptyset \\
\perp_{(\Sigma, R, \text{const}(r_i.a_i), P)}^F &= \emptyset \\
\perp_{(\Sigma, R, \text{exists}(Q), P)}^F &= \perp_Q^F[?] \\
\perp_{(\Sigma, R, \text{notexists}(Q), P)}^F &= \perp_Q^T[?]
\end{aligned}$$

Proposition 5.

$$\begin{aligned}
x \in \text{Eval}_{SQL}(Q_*, D) &\Rightarrow \forall r_i.a_i \in \perp_Q^T, x[r_i.a_i] \neq \perp \\
x \notin \text{Eval}_{SQL}(Q_*, D) &\Rightarrow \forall r_i.a_i \in \perp_Q^F, x[r_i.a_i] \neq \perp
\end{aligned}$$

Proof.

□

Definition 14.

$$\begin{aligned}
nested^+(H_1 \wedge H_2) &= \{H_1 \wedge H_2\} \cup nested^+(H_1) \cup nested^+(H_2) \\
nested^+(H_1 \vee H_2) &= \{H_1 \vee H_2\} \cup nested^+(H_1) \cup nested^+(H_2) \\
nested^+(r_i.a_i = c_i) &= \{r_i.a_i = c_i\} \\
nested^+(r_i.a_i \neq c_i) &= \{r_i.a_i \neq c_i\} \\
nested^+(r_i.a_i = r_j.a_j) &= \{r_i.a_i = r_j.a_j\} \\
nested^+(null(r_i.a_i)) &= \{null(r_i.a_i)\} \\
nested^+(const(r_i.a_i)) &= \{const(r_i.a_i)\} \\
nested^+(exists(Q)) &= \{exists(Q)\} \cup nested^+(Q) \\
nested^+(notexists(Q)) &= \{notexists(Q)\} \cup nested^-(Q)
\end{aligned}$$

$$\begin{aligned}
nested^-(H_1 \wedge H_2) &= nested^-(H_1) \cup nested^-(H_2) \\
nested^-(H_1 \vee H_2) &= nested^-(H_1) \cup nested^-(H_2) \\
nested^-(r_i.a_i = c_i) &= \emptyset \\
nested^-(r_i.a_i \neq c_i) &= \emptyset \\
nested^-(r_i.a_i = r_j.a_j) &= \emptyset \\
nested^-(null(r_i.a_i)) &= \emptyset \\
nested^-(const(r_i.a_i)) &= \emptyset \\
nested^-(exists(Q)) &= nested^-(Q) \\
nested^-(notexists(Q)) &= nested^+(Q)
\end{aligned}$$

$$nested(Q) = nested^-(Q) \cup nested^+(Q)$$

Definition 15.

$$\begin{aligned}
notexists(Q')_{\bar{Q}}^{\perp} &\rightarrow notexists(Q'_{\bar{Q}}^{\perp}) \\
exists(Q')_{\bar{Q}}^{\perp} &\rightarrow exists(Q'_{\bar{Q}}^{\perp})
\end{aligned}$$

Definition 16.

$$\begin{aligned}
(\Sigma, R, H \vee null(r_i.a_i), P)_{\bar{Q}}^{\perp} &\rightarrow (\Sigma, R, H_{\bar{Q}}^{\perp}, P) \\
&\text{if } \exists Q' \in nested^+(Q), (H \vee null(r_i.a_i)) \in nested(Q'), r_i.a_i \in \perp_{Q'}^T, \\
&\text{if } \exists Q' \in nested^-(Q), (H \vee null(r_i.a_i)) \in nested(Q'), r_i.a_i \in \perp_{Q'}^F,
\end{aligned}$$

Proposition 6.

$$Eval_{SQL}(Q, D) = Eval_{SQL}(Q_Q^\perp, D)$$

Proof.

Definition 17.

$$\begin{aligned} (\Sigma, R, H \vee null(r_i.a_i), P)_{Q_Q^\perp}^{\perp^F} &\rightarrow (\Sigma, R, H_Q^{\perp^F}, P) \\ &\text{if } \exists Q' \in nested^-(Q), (H \vee null(r_i.a_i)) \in nested(Q'), r_i.a_i \in \perp_{Q'}^T, \\ &\text{if } \exists Q' \in nested^+(Q), (H \vee null(r_i.a_i)) \in nested(Q'), r_i.a_i \in \perp_{Q'}^F, \end{aligned}$$

Proposition 7.

$$\forall Q' \in nested^+(Q), Eval(Q, D) = Eval(Q_{Q'}^\perp, D)$$

Proposition 8.

$$\forall Q' \in nested^-(Q), Eval(Q, D) = Eval(Q_{Q'}^{\perp^F}, D)$$

Proof. Assume (10).

By Induction ...

$$H_1 \wedge H_2 \in nested^+(Q)$$

$$\exists r_i.a_i \in (\perp_{H_1}^T \setminus \perp_{H_2}^T) \wedge (H \vee null(r_i.a_i)) \in nested(H_2)$$

Then if $x[r_i.a_i] = \perp$, $Eval(H_{2,H_2}^T, D, x)$ might be different from $Eval(H_{2,H_1 \wedge H_2}^T, D, x)$ □

□

Proposition 9.

5 Moving around null check

Definition 18. For each Query and nested Query we maintain a set of attribute that have to be not null in order for the query to be true resp. false we denote this set \perp_Q^T resp. \perp_Q^F .

Definition 19.

$$\begin{aligned}
\perp_{(\Sigma, R, H_1 \wedge H_2, P)}^T &= \perp_{(\Sigma, R, H_1, P)}^T \cup \perp_{(\Sigma, R, H_2, P)}^T \\
\perp_{(\Sigma, R, H_1 \vee H_2, P)}^T &= \perp_{(\Sigma, R, H_1, P)}^T \cap \perp_{(\Sigma, R, H_2, P)}^T \\
\perp_{(\Sigma, R, r_i.a_i = c_i, P)}^T &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i \neq c_i, P)}^T &= \{r_i.a_i\} \\
\perp_{(\Sigma, R, r_i.a_i > c_i, P)}^T &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i < c_i, P)}^T &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i = r_j.a_j, P)}^T &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i \neq r_j.a_j, P)}^T &= \{r_i.a_i, r_j.a_j\} \\
\perp_{(\Sigma, R, r_i.a_i > r_j.a_j, P)}^T &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i < r_j.a_j, P)}^T &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i = p_i, P)}^T &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i \neq p_i, P)}^T &= \{r_i.a_i, p_i\} \\
\perp_{(\Sigma, R, r_i.a_i > p_i, P)}^T &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i < p_i, P)}^T &= \emptyset \\
\perp_{(\Sigma, R, \text{const}(r_i.a_i), P)}^T &= \{r_i.a_i\} \\
\perp_{(\Sigma, R, \text{const}(p_i), P)}^T &= \{p_i\} \\
\perp_{(\Sigma, R, \text{null}(r_i.a_i), P)}^T &= \emptyset \\
\perp_{(\Sigma, R, \text{null}(p_i), P)}^T &= \emptyset \\
\perp_{(\Sigma, R, \text{exists}(Q), P)}^T &= \perp_Q^T[?] \\
\perp_{(\Sigma, R, \text{notexists}(Q), P)}^T &= \perp_Q^F[?]
\end{aligned}$$

$$\begin{aligned}
\perp_{(\Sigma, R, H_1 \wedge H_2, P)}^F &= \perp_{(\Sigma, R, H_1, P)}^F \cap \perp_{(\Sigma, R, H_2, P)}^F \\
\perp_{(\Sigma, R, H_1 \vee H_2, P)}^F &= \perp_{(\Sigma, R, H_1, P)}^F \cup \perp_{(\Sigma, R, H_2, P)}^F \\
\perp_{(\Sigma, R, r_i.a_i = c_i, P)}^F &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i \neq c_i, P)}^F &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i > c_i, P)}^F &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i < c_i, P)}^F &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i = r_j.a_j, P)}^F &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i \neq r_j.a_j, P)}^F &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i > r_j.a_j, P)}^F &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i < r_j.a_j, P)}^F &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i = p_i, P)}^F &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i \neq p_i, P)}^F &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i > p_i, P)}^F &= \emptyset \\
\perp_{(\Sigma, R, r_i.a_i < p_i, P)}^F &= \emptyset \\
\perp_{(\Sigma, R, \text{const}(r_i.a_i), P)}^F &= \emptyset \\
\perp_{(\Sigma, R, \text{const}(p_i), P)}^F &= \emptyset \\
\perp_{(\Sigma, R, \text{null}(r_i.a_i), P)}^F &= \{r_i.a_i\} \\
\perp_{(\Sigma, R, \text{null}(p_i), P)}^F &= \{p_i\} \\
\perp_{(\Sigma, R, \text{exists}(Q), P)}^F &= \perp_{Q[?]}^F \\
\perp_{(\Sigma, R, \text{notexists}(Q), P)}^F &= \perp_{Q[?]}^T
\end{aligned}$$

Definition 20.

$$\begin{aligned}
(\Sigma, R, H, P)^{up} &\rightarrow (\Sigma, R, (H \setminus \{\text{null}(p_i)\}) \wedge \text{const}(p_i), P) \\
&\quad \text{if } p_i \in \perp_{(\Sigma, R, H, P)}^T \\
(\Sigma, R, H, P)^{up} &\rightarrow (\Sigma, R, (H \setminus \{\text{null}(r_i.a_i)\}) \wedge \text{notexists}((*, R, \text{null}(r_i.a_i), P)), P) \\
&\quad \text{if } r_i.a_i \in \perp_{(\Sigma, R, H, P)}^T
\end{aligned}$$

Not well written, actually false (notexists Query is different than that) but intuitive.

Proposition 10.

$$\forall Q \in \llbracket SQL \rrbracket_{\perp} \text{Eval}_{SQL}(Q, D) = \text{Eval}_{SQL}(Q^{up}, D)$$

6 Optimizer

We can read in postgres source code : "We stop as soon as we hit a non-AND item." Indeed whenever postgres try to optimize join thanks to index or hash function it stop when a disjunction appears. Moreover as seen in previous section our translation involve adding disjunction in sub-queries, meaning the optimizer will choose to do a nested-loop in order to compute the answer, resulting in a blow up in time. In order to keep the computation time reasonable we offer an rewritting of the query which might seem counter-intuitive but allow the optimizer to use hash and index in order to perform join.

First of all please notice that the only disjunction that might appear in our rewritting are of the form of $r_i.a = r_j.b \vee null(r_i.a) \vee null(r_j.b)$ then we will focus on the case even if most of the optimization might be applicable in more general case (ie. if the null check doesn't concern the same attribute then we have to care when performing *UNIONALL*, if the first literal of the disjunction isn't an equality optimization is useless because the optimizer won't be able to do a hash or index join.) Moreover if the primary query already contain disjunction we won't optimize it then we can consider that our query after translation is in CNF.

6.1 Compute possible answer

Proposition 11.

$$Eval_{SQL}(\Sigma, R, (r_1.a = c \vee null(r_1.a)) \wedge H, P) = Eval_{SQL}(\Sigma, (R \setminus \{r_1\}) \cup \{r_1^\vee\}, H, P)$$

$$r_1^\vee = Eval_{SQL}(*, \{r_1\}, r_1.a = c \vee null(r_1.a), \emptyset)$$

Proof. $Q = (\Sigma, R, (r_1.a = c \vee null(r_1.a)) \wedge H, P)$

$$\begin{aligned} x \in^n Eval_{SQL}(Q*, D) &\Leftrightarrow x \in^n Eval_{SQL}(*, R, r_1.a = c \vee null(r_1.a), P) \cap Eval_{SQL}(*, R, H, P) \\ &\Leftrightarrow x \in^n Eval_{SQL}(*, R, r_1.a = c \vee null(r_1.a), P) \wedge x \in^n Eval_{SQL}(*, R, H, P) \\ &\Leftrightarrow R^\times(x) \geq n \wedge (x[r_1.a] = c \vee x[r_1.a] = \perp) \wedge x \in^n Eval_{SQL}(*, R, H, P) \\ &\dots \end{aligned}$$

□

Proposition 12.

$$Eval_{SQL}(\Sigma, R, ((r_1.a = r_2.b \vee null(r_1.a) \vee null(r_2.b)) \wedge H, P) = Eval_{SQL}(\Sigma, (R \setminus \{r_1, r_2\}) \cup \{r_{1 \vee 2}\}, H, P)$$

$$r_{1 \vee 2} = Eval_{SQL}(*, \{r_1, r_2\}, r_1.a = r_2.b, P) \uplus Eval_{SQL}(*, \{r_1, r_2^+\}, null(r_1.a), P) \uplus Eval_{SQL}(*, \{r_1, r_2\}, null(r_2.b), P)$$

$$r_2^+ = Eval_{SQL}(*, \{r_2\}, const(r_2.b), H, P)$$

6.2 Compute certain answer

In order to compute certain answer we do not necessarily have to compute possible answer, indeed we only have to compute a set representing the possible answer which interest us.

Proposition 13.

$$(\Sigma, R, \text{notexists}((*, R', H \wedge (r_1.a = r_2.b \vee \text{null}(r_1.a)), P')), P) = (\Sigma, R, \text{notexists}((*, R', H \wedge r_1.a = r_2.b, P')) \wedge \text{not}$$