

Algebra for full-text queries II

Oleg Bartunov, Teodor Sigaev, Mikhail Prokhorov
Moscow University, Sternberg Astronomical Institute
2008



Операция \$: *BEFORE (AFTER)*

- Поиск фразы из 2 слов

$$a \$[n] b \iff \exists i, j : \text{pos}(b)_i - \text{pos}(a)_j = n$$

- Оператор *BEFORE* (\$n) ищет фразу
 - с заданным порядком слов (операндов) — *a BEFORE b* (или *b AFTER a*)
 - с заданным расстоянием *n* (default is 1)



Обобщения: операции $\$[m,n]$ и т.п.

- $a \ \$[m,n] \ b \ == \ \exists \ i,j: \ m \leq \text{pos}(b)_i - \text{pos}(a)_j \leq n$
- $a \ \$[<n] \ b \ == \ \exists \ i,j: \ 0 < \text{pos}(b)_i - \text{pos}(a)_j \leq n$
- $a \ \$[>n] \ b \ == \ \exists \ i,j: \ \text{pos}(b)_i - \text{pos}(a)_j > n$
- И т.п.



Действия с \$

- $a \$[n] b = b \$[-n] a$
- $!(a \$[n] b) = \forall i, j : \text{pos}(b)_i - \text{pos}(a)_j \neq n$
- $!!(a \$[n] b) = a \$[n] b$
- $a \$!b = a \& (\exists i, j : \text{pos}(!b)_i - \text{pos}(a)_j = 1)$
- $!a \$ b = b \& (\exists i, j : \text{pos}(b)_i - \text{pos}(!a)_j = 1)$
- $!a \$!b = (\exists i, j : \text{pos}(!b)_i - \text{pos}(!a)_j = 1)$
- $a \$ (b | c) = a \$ b | a \$ c$
- $(b | c) \$ a = b \$ a | c \$ a$
- $a \$ (b \& c) = b \& c \& (a \$ b | a \$ c);$
- $(b \& c) \$ a = b \& c \& (b \$ a | c \$ a)$



Действия общего вида

- $a \$[n] b$ – в тексте присутствует b через n позиций после a (или a за n позиций перед b)
- $a \$[n] b = b \$[-n] a$ – n позиций **перед** то же самое, что на $-n$ позиций **за (после)**
- $!(a \$[n] b)$ – в тексте нет b через n позиций после a
- $!!(a \$[n] b) = a \$[n] b$ – двойное отрицание $\$$ дает исходный результат



Операции с отрицанием

- $a \$!b$ – в тексте есть **a** за которым идет **не b**
- $!a \$ b$ – в тексте есть **b** перед которым стоит **не a**
- $!a \$!b$ – в тексте есть последовательная пара слов первое из которых **не a**, а второе **не b**



Сочетание \$ с AND и OR

- $a \$ (b | c)$ – в тексте есть **a** за которым идет **b или c**
- $(b | c) \$ a$ – в тексте есть **b или c** за которым идет **a**
- $a \$ (b \& c)$ – в тексте есть **a** за которым идет **b И c**
(см. пояснение ниже)
- $(b \& c) \$ a$ – в тексте есть **b И c** за которым идет **a**
(см. пояснение ниже)



Ускорение поиска

- Рассмотрим определение основной операции

$$a \$[n] b \iff \exists i, j : \text{pos}(b)_i - \text{pos}(a)_j = n$$

- Для ускорения выполнения запроса усложним выражение

$$a \$[n] b \iff a \& b \& (\exists i, j : \text{pos}(b)_i - \text{pos}(a)_j = n)$$

- Смысл усложнения: если в тексте нет одного из аргументов (a или b), то фразы быть не может и проверка позиций не нужна



Отрицание \$

- Рассмотрим отрицание основной операции

$$!(a \$[n] b) == \forall i,j : \text{pos}(b)_i - \text{pos}(a)_j \neq n$$

- Для ускорения запроса усложним выражение

$$a \$[n] b == !a \mid !b \mid (\forall i,j : \text{pos}(b)_i - \text{pos}(a)_j \neq n)$$

- Смысл усложнения: если в тексте нет одного из аргументов (a или b), то запрос верен и проверка списков позиций не требуется



Сочетание \$ и &

- Рассмотрим сочетание операций \$ и &: **$a \$ (b \& c)$**
- Возможны две интерпретации:
 - $a \$ (b \& c) = (a \$ b) \& (a \$ c)$ – в тексте должны присутствовать обе пары $a \$ b$ и $a \$ c$
 - $a \$ (b \& c) = b \& c \& (a \$ b \mid a \$ c)$ – в тексте должны присутствовать оба аргумента b и c и один из них должен следовать за a : $(a \$ b \mid a \$ c)$
- Пример с заменой синонимов
 $'some' \$ 'telefonsvarer' \Rightarrow 'some' \$ ('telefon' \& 'svar')$
служит аргументом **в пользу 2-го варианта**



Запросы с отрицаниями

Вырождение запросов с отрицаниями

- $a \$!b - \dots a \times \dots a \times \dots aab \dots a \times \dots$
- $a \$ [2] !b - \dots a \times b \dots a \times b \dots aab \dots a \times b \dots$
- $a \$ [<3] !b - \dots a \times b \dots acb \dots aab \dots azb \dots$
- $!a \$ b - \dots a \times \dots a \times \dots aab \dots a \times \dots$
- $!a \$!b - \dots a \times \dots a \times \dots a \times \dots a \times \dots xy \dots$

НО данному выражению не удовлетворяет только текст

$n \geq 0$ раз \rightarrow **aaa...aaabbb...bbb** \leftarrow $m \geq 0$ раз



Регулярные выражения и \$

- $a [xyz] b \Leftrightarrow a \$ (x | y | z) \$ b$
- $a . b \Leftrightarrow a \$[2] b$
- $a \dots b \Leftrightarrow a \$[4] b$
- $a .^* b \Leftrightarrow a \$[1, \infty] b$
- $a [xy]\{1,2\} b \Leftrightarrow a \$ (x | y | x\$x | x\$y | y\$x | y\$y) \$ b$



Фразы (>2 слов)

Фраза задается:

- 1) a ... b ... c ... b ... d x - порядком слов
- 2) \$ \$[4] \$[<3] \$[3,7] \$[n] ... \$[m] - интервалами

В тексте положение фразы задается позициями:

posL – самое левое
слово

posR – самое правое
слово

*) Задание posL и posR может определить несколько фраз

***) для \$ и \$[n] $\text{posR}(A) - \text{posL}(A) = \sum n_i$

****) Для фразы из одного слова "a": $\text{posL}(a) = \text{posR}(a) = \text{pos}(a)$



Обобщение \$ на фразы

Пусть даны фразы

$$A = (a \$ b \$ c \dots y \$ z),$$

$$B = (a1 \$ \dots \$ z1),$$

тогда

$$A \$[n] B \iff \exists i, j : \text{posL}(B)_i - \text{posR}(A)_j = n$$



Поиск фраз в тексте

Определение через поиск слов

$a \text{ } \$[n_a] \text{ } b \text{ } \$[n_b] \text{ } c \text{ } \dots \text{ } x \text{ } \$[n_x] \text{ } y \text{ } \$[n_y] \text{ } z \text{ } ==$

$$\begin{aligned} \exists i_a, i_b, \dots, i_y, i_z: & \quad \text{pos}(b)_{i_b} - \text{pos}(a)_{i_a} = n_a \ \& \\ & \quad \text{pos}(c)_{i_c} - \text{pos}(b)_{i_b} = n_b \ \& \\ & \quad \dots \\ & \quad \text{pos}(y)_{i_y} - \text{pos}(x)_{i_x} = n_x \ \& \\ & \quad \text{pos}(z)_{i_z} - \text{pos}(y)_{i_y} = n_y \end{aligned}$$



Поиск фраз в тексте

Определение через поиск фраз

$A == (a1 \text{ } [n_a] \text{ } b1 \text{ } \dots \text{ } y1 \text{ } [n_y] \text{ } z1) ==$

$\exists i_a, i_b, \dots, i_y, i_z: \text{pos}(b1)_{i_b} - \text{pos}(a1)_{i_a} = n_a \ \& \ \dots \ \& \ \text{pos}(z1)_{i_z} - \text{pos}(y1)_{i_y} = n_y$

$B == (a2 \text{ } [m_a] \text{ } b2 \text{ } \dots \text{ } y2 \text{ } [m_y] \text{ } z2) ==$

$\exists j_a, j_b, \dots, j_y, j_z: \text{pos}(b2)_{j_b} - \text{pos}(a2)_{j_a} = m_a \ \& \ \dots \ \& \ \text{pos}(z2)_{j_z} - \text{pos}(y2)_{j_y} = m_y$

$A \text{ } [k] \text{ } B == \exists i, j: \text{posL}(B)_i - \text{posR}(A)_j = k \Rightarrow$

$\text{posR}(A) = \text{pos}(z1) \ \& \ \text{posL}(B) = \text{pos}(a2) \Rightarrow$

$\exists i_a, \dots, i_z, j_a, \dots, j_z:$

$\text{pos}(b1)_{i_b} - \text{pos}(a1)_{i_a} = n_a \ \& \ \dots \ \& \ \text{pos}(z1)_{i_z} - \text{pos}(y1)_{i_y} = n_y$

$\ \& \ \text{pos}(a2)_{j_a} - \text{pos}(z1)_{i_z} = k \ \&$

$\text{pos}(b2)_{j_b} - \text{pos}(a2)_{j_a} = m_a \ \& \ \dots \ \& \ \text{pos}(z2)_{j_z} - \text{pos}(y2)_{j_y} = m_y$



Поиск фраз в тексте

Вывод:

Разбиение фразы на произвольные части не меняет результата.

Это позволяет организовывать поиск рекурсивно, наиболее удобным образом.



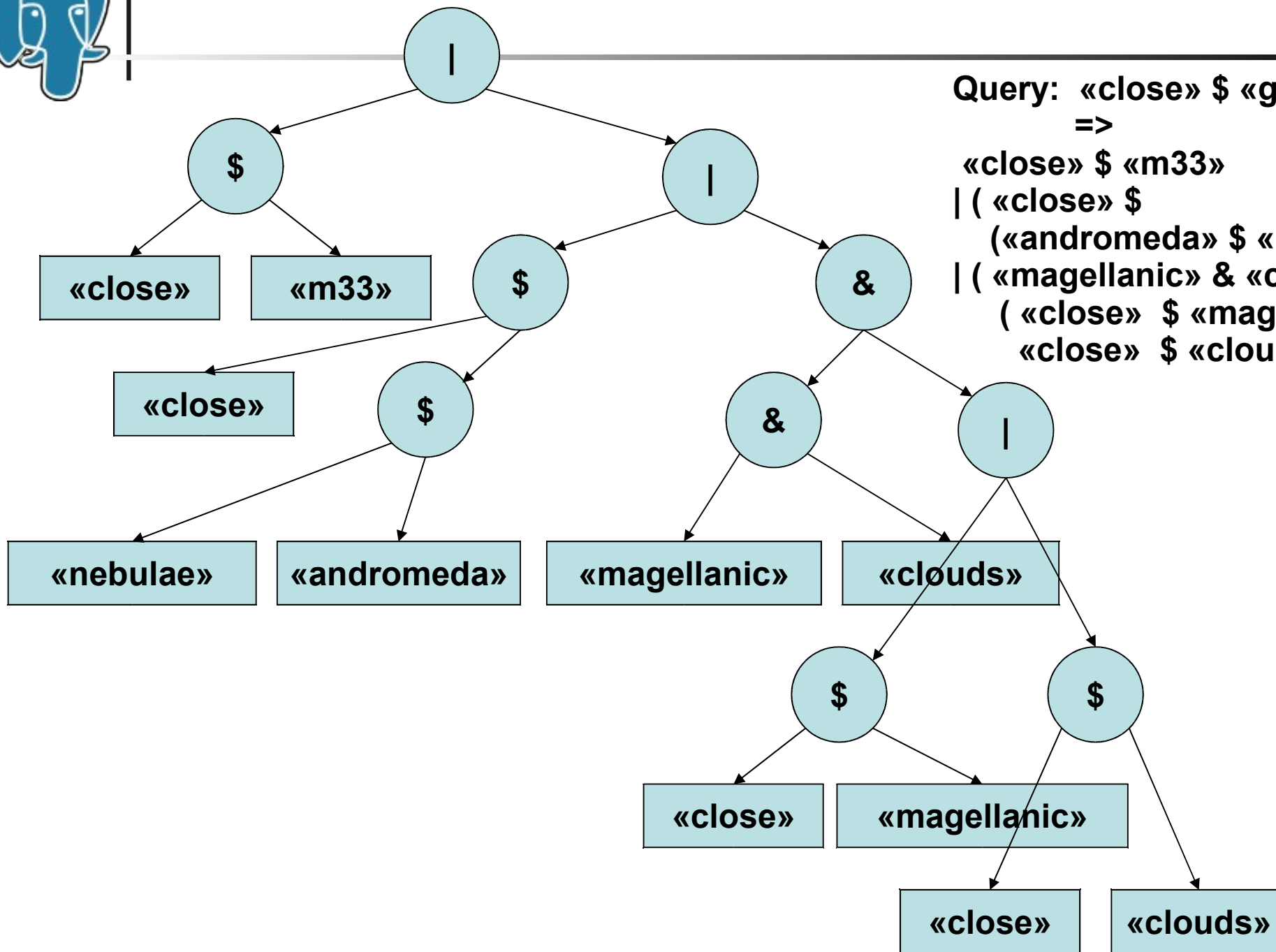
Частный случай (2 слова+1)

$$\begin{aligned} (a \ \$[n] \ b) \ \$[m] \ c &= (a \ \$[n] \ b) \ \& \ c \ \& \\ & \quad (\exists i, j: \text{posL}(c)_j - \text{posR}(\text{"ab"})_i = m) \Rightarrow \\ (a \ \$[n] \ b) \ \& \ c \ \& \ (\exists i, j: \text{pos}(c)_j - \text{posR}(\text{"ab"})_i = m) & \Rightarrow \\ (a \ \& \ b \ \& \ (\exists k, l: \text{pos}(b)_k - \text{pos}(a)_l = n)) \ \& \\ & \quad c \ \& \ (\exists i, j: \text{pos}(c)_j - \text{posR}(\text{"ab"})_i = m) = \\ = a \ \& \ b \ \& \ c \ \& \\ & \quad (\exists k, l: \text{pos}(b)_k - \text{pos}(a)_l = n) \ \& \ (\exists j: \text{pos}(c)_j - \text{pos}(b)_k = m) = \\ = a \ \& \ b \ \& \ c \ \& \ (\exists j, k, l: \text{pos}(b)_k - \text{pos}(a)_l = n \ \& \\ & \quad \text{pos}(c)_j - \text{pos}(b)_k = m) = \\ = a \ \$[n] \ b \ \$[m] \ c \end{aligned}$$

Аналогично: $a \ \$[n] \ (b \ \$[m] \ c) = a \ \$[n] \ b \ \$[m] \ c$



Оптимизация поиска



Query: «close» \$ «galaxies»
=>
«close» \$ «m33»
| («close» \$
(«andromeda» \$ «nebulae»))
| («magellanic» & «clouds» &
(«close» \$ «magellanic» |
«close» \$ «clouds»))



Поглощение выражений

- $a \$ b \subset a \$[n] b \quad (n \geq 1)$
- $a \$[m] b \subset a \$[n] b \quad (n \geq m)$
- $a \$[m_1, n_1] b \subset a \$[m, n] b \quad (m \leq m_1 \ \& \ n_1 \leq n)$
- $b \$ c \subset a \$ b \$ c \$ d$

Если правое выражение TRUE => левое тоже TRUE

Если левое выражение известно – ускоряется вычисление правого

