

【VLDB 2013勉強会】

Session 12: Spatial and Text

担当：石川佳治

Some figures are copied from VLDB 2013 proceedings.

Efficient Error-tolerant Query Autocompletion

- ▶ C. Xiao, J. Qin, W. Wang, Y. Ishikawa, K. Tsuda, K. Sadakane (Nagoya U./UNSW/NII/AIST)

- ▶ Chuan Xiao (肖川): 石川研ポスドク

- ▶ サーチ処理におけるキーワード補完

- ▶ 入力文字列の誤りも許す

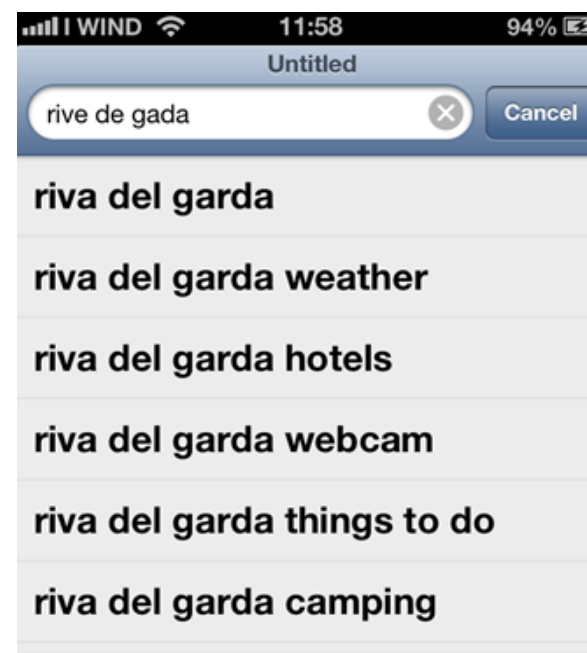
- ▶ 編集距離に基づく接頭辞探索 (Edit Prefix Search)

- ▶ 例: 編集距離の閾値 $\tau = 1$

- ▶ 検索文字列 $q = \text{"abc"}$

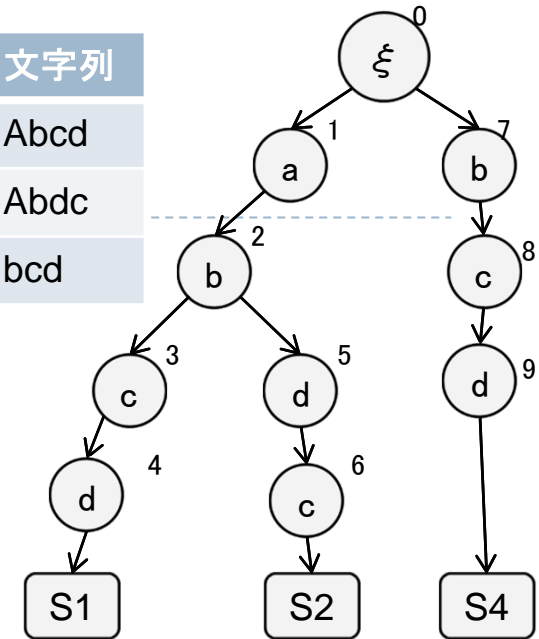
- ▶ 文字列集合 $S = \{\text{"acdefg"}, \text{"cda"}, \dots\}$

- ▶ “acdefg”は条件にマッチ: 接頭辞の“ac”について $ED(\text{"abc"}, \text{"ac"}) = 1$ が成立



アプローチ

ID	文字列
S1	Abcd
S2	Abdc
S4	bcd



▶ 既存手法

- ▶ 文字列集合 S について**トライ(trie)**を構築
- ▶ 検索時には編集距離を考慮しながら探索
- ▶ 問題点: **多数のノードへのアクセス** $O(|\Sigma|^\tau)$

▶ 本論文の貢献

- ▶ 新たなデータ構造 (**Deletion Variants Trie**) の提案
- ▶ Edit Prefix Search問題を**Exact Prefix Search問題**に変換
- ▶ **1,000倍**の高速化(ただし20倍の索引サイズ)

▶ **Deletion Variant**の概念を利用: $s = \text{“abcd”}$ の場合

- ▶ 0-variants: abcd {}
- ▶ 1-variants: bcd {1}, acd {2}, abd {3}, **abc {4}**
- ▶ 2-variants: cd {1,1}, bd {1,2}, bc {1,3}, ad {2,2}, ac {2,3}, ab {3,3}

4番目の文字を削除

アプローチ(続き)

▶ Deletion Variantを用いた判定

▶ 例: $s = \text{“abcd”}$, $q = \text{“abxd”}$ のとき, 両者とも $(abd, \{3\})$ という Deletion Variantを持ち, $\tau = 1$ の編集距離

▶ 性質: 二つの Deletion Variants (x, D_x) , (y, D_y) があるとき, $x = y$ かつ $|D_x \cup D_y| \leq \tau$ なら $ED(s, q) \leq \tau$

▶ **Deletion Variants Trie**: $s = \text{“abcd”}$ で $\tau = 2$ のとき, $\{abcd, \#bcd, a\#cd, ab\#d, abc\#, \#\#cd, \#b\#d, \dots\}$ をトライに挿入

▶ 検索処理: $q = \text{“abc”}$, $\tau = 2$ のとき

▶ $\tau = 2$ を考慮し, $abc, \#bc, a\#c, ab\#, \#\#c, \#b\#, a\#\#$ と展開

▶ 検索文字列の列挙(それぞれを実際に探す文字列に展開)

▶ $a\#c$ の場合, $ac, a\#c, \#a\#c, a\#\#c, a\#c\#$

▶ 実際には, 検索文字入力に応じて適応的に処理 $O(\tau \cdot (|q| + \tau)^\tau)$