

図3 提案手法

い遷移であることが確認できる。もし CHECK が一致しなかった場合、クエリが単語集合に存在しないとして、探索を打ち切る。接続節 t から次の HP の先頭節、つまり t の HC への遷移は、 t に対応する NEXT により決定する。例えば、接続節 4 において、NEXT[4]=6 より次の HP の先頭節 6 へ遷移する。

以上の探索により SP 節へ到達したとき、SP 節に対応する BASE、または NEXT の負の値により TAIL へのリンクを実現し、以降の探索は SP ストリングとクエリの未探索の文字列を比較することで達成する。例えば、接続節 10 において、NEXT[10]=-4 より TAIL[4]から“e#”を参照する。本手法では、ある節の HP に対応する文字列との比較に失敗したとき、節に対応する BASE を用いて接続節への遷移をおこなう。そのため、トライの単語を決定づける節を、SP 節として用いることができない節が存在する。図2の節9,12はSP節だが、図3の節4,7はSP節ではない。そこで、従来のSP節を利用できない節に対し、自身の子をSP節として実現した。図3の節5,8がこれにあたる。

本手法によって、図3の節1から5の遷移と6から8の遷移では式1を必要とせず、文字列の比較のみで辿ることができるため、高速である。また、接続節への遷移も式1により $O(1)$ でおこなうため、高速である。

4. 評価

ダブル配列へ CPD を適用することで、ダブル配列の検索が高速化されたことを、実験により検証する。単語集合は日本語、英語の Wikipedia のページタイトル¹から、ランダムに抽出した 10 万語から 50 万語とした。それぞれ単語集合に含まれる単語すべてを検索し、1 単語あたりの検索時間により、検証をおこなった。実験環境は CPU が Core i7 の 3.1GHz、L1, L2 キャッシュが 256KB, 1MB、メモリが 8GB でおこなった。

図4に単語数と検索時間の関係を示し、図5, 6に単語数とデータ容量の関係を示す。Wikipedia のページタイトルの日本語を JP とし、英語を EN とする。平均データ容量は JP, EN で、提案手法が従来の 13% 増となったものの、平均検索時間において、JP では従来の 91% に短縮し、EN では 85% に短縮することができた。データ容量増加の原因は、一部の SP 節をトライの後方へ移動したためと考えられる。これらの実験結果により、提案手法を適用することでデータ容量の増加はあるものの、高速な検索を実現できることがわかる。

¹ Wikipedia ページタイトル <http://dumps.wikimedia.org/>

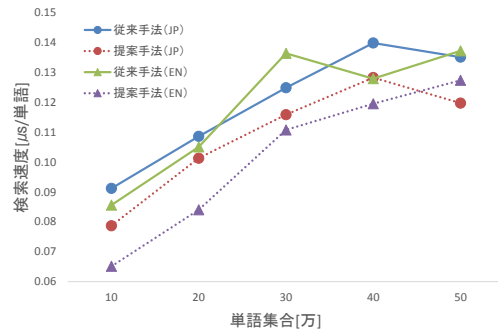


図4 Wikipedia 日英の検索時間

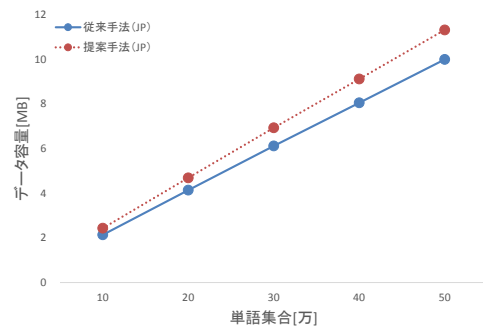


図5 Wikipedia 日本語のデータ容量

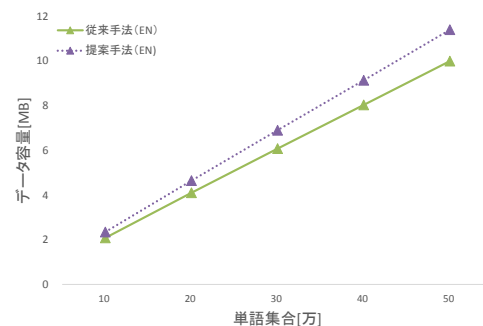


図6 Wikipedia 英語のデータ容量

5. おわりに

本稿では、トライの検索の高速化を実現する CPD をダブル配列に適用することにより、HP 上を文字列比較のみで辿ることで、ダブル配列の検索を高速化できることを示した。今後の課題として、増加したデータ容量を削減することが挙げられる。

参考文献

- [1] 青江 順一, “トライとその応用 (<連載講座> キー検索技法 4)”, 情報処理, Vol.34, No.2, pp.224-251 (1993)
- [2] 青江 順一, “ダブル配列による高速デジタル検索アルゴリズム”, 電子情報通信学会論文誌 D 情報・システム, Vol.71, No.9, pp.1592-1600 (1989)
- [3] 矢田 晋, 森田 和宏, 泓田 正雄, 平石 互, 青江 順一, “ダブル配列におけるキャッシュの効率化”, FIT2006, pp.1066-1077 (2006)
- [4] Paolo Ferragina, Roberto Grossi, Ankur Gupta, Rahul Shah, Jeffrey Scott Vitter, “On Searching Compressed String Collections Cache-Obliviously”, PODS, pp.181-190 (2008)