

図 2 : Suffix DAWG(gttagtaaac)

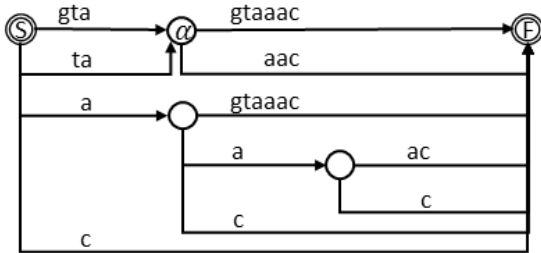


図 3 : Suffix CDAWG(gttagtaaac)

いる。ここで、Suffix Tree を CDAWG で表現することが可能と考え、このデータ構造を Suffix CDAWG と定義する。有限文字列 $w = \text{gttagtaaac}$ に対する Suffix CDAWG を図 3 に示す。ノード α から遷移し始める "gtaaac" のような、分岐なしに遷移が終端まで続く箇所が、CDAWG では 1 つの遷移にまとめられていることがわかる。この箇所は Separate String と見なすことができる。よって、ルートノードからノード α までは従来通りダブル配列に登録し、 α 以降の Separate String を TAIL として登録する。TAIL には接尾辞が格納されるので、集約することが可能となり、効率的である。本提案手法では、グラフ構築に際しての省空間化、計算時間の短縮を図る。

5. 評価

全文検索において提案手法の有効性を実証するために検証実験をおこなった。Intel Xeon 2.53GHz Quad-Core \times 2, L2 キャッシュ 8M を実験環境として用いた。また、検索対象となる全文は Pizza&Chili Corpus¹において公開されている DNA データ 50MB を使用し、そのうち先頭から文字列長分を抽出した。DNA データであるため、文字種は 4 種類のみである。今回は出現箇所を考慮せず、文書中にキーが含まれるか否かを検索する。全文のランダムな場所から抽出した 10000 個のキーに対する検索時間の合計を、検索時間として計測し、10 回の平均を取った。検索キー長は 81 文字とした。比較対象は簡潔データ構造によって実現されている圧縮接尾辞配列ライブラリ csalib100810²を用いた。圧縮する際のオプションは "-P3:512 -I:128:256" とした。検索

表 1 : キー長 81 文字での検索時間[ms]

文字列長	100,000	500,000	900,000
csalib	110.4	108.6	114.5
従来のダブル配列	25.9	—	—
提案手法	10.6	12.4	12.4

表 2 : 辞書容量

文字列長	100,000	500,000	900,000
csalib	30KiB	147KiB	262KiB
従来のダブル配列	1.0MiB	—	—
提案手法	0.8MiB	4.2MiB	7.6MiB

時間と辞書容量の比較結果を表 1, 2 に示す。また、Suffix DAWG をダブル配列で構築した場合を、従来のダブル配列として比較する。

結果として、csalib は文字列長が増加するにつれて、検索時間も増加するのに対し、提案手法は一定の検索速度を保つことが確認された。また、従来のダブル配列では長い文字列長において辞書構築できなかった。短い文字列から構築した辞書においては、従来のダブル配列と比べて、2 倍以上の検索速度を達成した。

6. おわりに

本論文では、キー検索法の 1 つであるダブル配列について、CDAWG を用いてグラフ構造を表現することにより、検索速度が高速な全文検索手法を提案した。今後は、より長い文字列長における実験と複数の文書を用いた全文検索実験を課題として設定する。複数文書における検索でも、部分共通文字列を統合できる本提案手法が活かせると考えている。

参考文献

- [1] Aoe, J.: An efficient digital search algorithm by using a double-array structure, Software Engineering, IEEE Transactions on, Vol. 15, No. 9, pp. 1066--1077 (1989).
- [2] Yata, S, Morita, K, Fuketa, M, and Aoe, J, "Fast string matching with space-efficient word graphs", In Proc. Innovations in Information Technology, pp. 79--83(2008).
- [3] Crochemore, M. and Rytter, W.: Jewels of Stringology, World Scientific (2002).
- [4] Crochemore, Maxime, and Renaud V erin. "On compact directed acyclic word graphs." Structures in Logic and Computer Science. Springer Berlin Heidelberg, pp. 192--211(1997).

¹ Pizza&Chili Corpus, <http://pizzachili.dcc.uchile.cl/>

² 圧縮接尾辞配列 csalib, <http://researchmap.jp/sada/csalib/>