

動的なコード生成を用いた正規表現マッチャの実装

新屋良磨[†] 河野真治[†]

与えられた正規表現から、等価な有限状態オートマトンに変換し、オートマトンにおける状態遷移を C や Continuation based C 等のコンパイル型言語での関数遷移に変換する正規表現コンパイラを Python で開発した¹⁾。本論文では、その正規表現コンパイラを利用した、テキストファイルに対しパターンマッチを行いマッチする行を出力するツールである grep に特化した実装方法を説明し、実装した grep の性能を評価する。

Implementation of Regular Expression Engine with Dynamic Code Generation.

RYOMA SHINYA[†] and SHINJI KONO[†]

We implemented regular expression compiler using Python. It convertes regular expression to equivalence automaton(DFA). Then, translate automaton transition to functional/code-segmental transition which is wirtten in compiled language like C, Continuation based C. In this paper, we introduce efficient implementation of grep that uses this compiler, and evaluate its performance.

1. はじめに

コンパイラ理論の発展と共に、コンパイルにかかる時間はより短く、また得られるプログラムはアセンブラレベルで最適化が施され、より高速になってきている。

完全に静的なコンパイルが可能な対象として、正規表現マッチャ(エンジン)に着目した。現在、正規表現のエンジンは、プログラミング言語の組み込み機能やライブラリ等、さまざまな実装が存在するが、それらの殆どは仮想マシン方式を採用している³⁾。仮想マシンを採用した実装でも、正規表現を内部表現に変換する処理を行っており、それらを“コンパイル”と呼ぶことが多い。本研究で実装したエンジンの“コンパイル”とは、正規表現を内部形式に変換することではなく、正規表現から実行バイナリを生成することを指す(3.3節)。本研究では、実行バイナリの生成にはコンパイラ基盤である LLVM, GCC を用いており、エンジン生成系は Python で実装した。

本論文では、まず正規表現のコンパイル方法について説明する。さらに、実装したエンジンの性能調査のために、正規表現を用いてテキストマッチ処理を行う grep クローンを実装し、GNU grep 等の既存ソフト

ウェアとの比較を行った。

2. 正規表現

2.1 正規表現によるテキストマッチ

正規表現は与えられた文字列を受理するかどうかを判定できるパターンマッチングの機構であり、sed, grep, awk を始めとするテキスト処理ツールに広く利用されている。正規表現には定められた文法によって記述され、例えば、正規表現 a^*b は a の 0 回以上の繰り返し直後、 b で終わる文字列 (b , ab , $aaaab$) を受理し、 $a(b|c)$ は a で始まり、直後が b または c で終わる文字列 (ab , ac) を受理する。

2.2 正規表現の演算

本論文では、以下に定義された演算をサポートする表現を正規表現として扱う。

- (1) 接続 二つの言語 L と M の接続 (LM) は、 L に属する列を一つとり、そのあとに M に属する列を接続することによってできる列全体から成る集合である。
- (2) 集合和 二つの言語 L と M 集合和 ($L|M$) は、 L または M (もしくはその両方) に属する列全体からなる集合である。
- (3) 閉包 言語 L の閉包 (L^*) とは、 L の中から有限個の列を重複を許して取り出し、それらを接続

[†] 琉球大学

University of the Ryukyus

してできる列全体の集合である。

正規表現は、この3つの演算について閉じており、この3つの演算によって定義される表現は、数学的には正則表現と定義されている。本論文では、特に区別のない限り、正則表現と正規表現を同じものとして扱う。

3. 正規表現エンジンの実装

正規表現は等価な NFA に、また NFA は等価な DFA に変換することが可能である⁵⁾。以下にその変換方法を説明する。

3.1 正規表現から NFA への変換

NFA(Non-deterministic Finite Automaton) は、入力に対して複数の遷移先を持つ状態の集合であり、遷移先が非決定的(Non-deterministic)である。ここでは、NFA を 5 個組 $(Q, \Sigma, \delta, q_0, F)$ で定義する。ただし、

- (1) Q は状態の有限集合。
- (2) Σ は入力記号の有限集合。
- (3) q_0 は Q の要素で、開始状態と呼ぶ。
- (4) F は Q の部分集合で、受理状態と呼ぶ。
- (5) δ は、状態と入力記号に対して状態の集合を返す遷移関数。 $(\epsilon$ 遷移を許す)

正規表現が、等価な NFA に変換できるということを、2.2 で定義した 3 つの演算について対応する NFA に変換できることから示す。

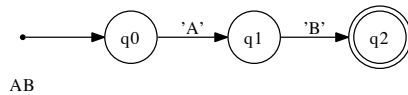


図 1 “A” と “B” の接続

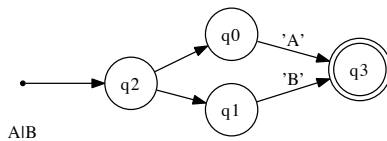


図 2 “A” と “B” の集合和

- (1) 接続 図 1 は正規表現 “AB” に対応する NFA.
- (2) 集合和 図 2 は正規表現 “A|B” に対応する NFA.
- (3) 閉包 図 3 は正規表現 “A*” に対応する NFA.

図 2, 3 において、ラベルのない矢印は無条件の遷移を現しており、 ϵ 遷移と呼ばれる。また、二重丸で囲まれた状態は受理状態を現しており、NFA において入力が終了した時点で、受理状態を保持している場合に限

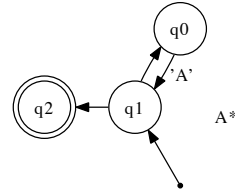


図 3 “A” の閉包

り、その文字列を受理したことになる。なお、NFA は同時に複数の遷移先をもつことがあるので、テキストのマッチング途中で複数の状態を保持することがある。

現在実装されている正規表現エンジンの多くは、正規表現を内部的に NFA に変換して評価を行っている²⁾。NFA による実装は、後述する後方参照や最長一致に対応しやすいが、同時に遷移する可能性のある複数の状態を保持する必要があるため、正規表現の複雑度によってマッチングの時間が多くなってしまう場合がある。文献²⁾では、“a?a?aaa” のような “a?a^n” のように表現 (“a?” は “a” “か空文字” を認識する拡張された正規表現の一つ) の評価において、NFA ベースの正規表現エンジンでは遷移する状態の数が増えてしまうのでマッチングにかかる処理時間が n の指数的に増加する問題をベンチマーク結果と共に指摘している。文献⁸⁾では正規表現から NFA ベースで効率的なマッチング処理を行うエンジンを IBM 7094 機械語で生成する例が紹介されている。

3.2 NFA から DFA への変換

非決定的な遷移を行う NFA から、決定的な遷移を行う DFA(Deterministic Finite Automaton) に変換する手法を説明する。なお、遷移が決定的であるということは、1 つの入力に対して、遷移する状態がただ 1 つであるということを示す。DFA は、NFA と同様な 5 個組で $(Q, \Sigma, \delta, q_0, F)$ 定義できる。ただし、DFA において δ において ϵ 遷移は認められず、また任意の状態 q と入力 σ について、 $\delta(q, \sigma) = q'$ となる q' は Q の要素となる。つまり、遷移先が決定的であるということに他ならない。

以下に ϵ 遷移を許す ϵ -NFA $E = (Q_E, \Sigma, \delta_E, q_0, F_E)$ から等価な DFA $D = (Q_D, \Sigma, \delta_D, q_D, F_D)$ を構成する手順を示す。

- (1) Q_D は Q_E の部分集合全から成る集合であり、おの中で D において到達可能な状態は、 ϵ 遷移に関して閉じている Q_E の部分集合 S に限られる。ここで、状態 q において ϵ 遷移に関して閉じている集合全体を $ECLOSE(q)$ と表す。 $ECLOSE$ を使って S を定義すると、

- $S = \bigcup_{q \in S} ECLOSE(q)$ を満たす S .
- (2) $q_D = ECLOSE(q_0)$. すなわち, E の開始状態の ε 閉包.
- (3) F_D は E の状態の集合で, 受理状態を少なくとも一つ含むもの全体からなる集合である. すなわち, $F_D = \{S | S \in Q_D \wedge S \cap F_E \neq \emptyset\}$
- (4) $\delta_D(S, a)$ は Q_D の要素 S と Σ の要素 a に対して次のように計算される.
- (a) $S = \{p_1, p_2, \dots, p_k\}$ とする.
- (b) $\bigcup_{i=1}^k \delta(p_i, a)$ を求め, その結果を $\{r_1, r_2, \dots, r_m\}$ とする.
- (c) このとき, $\delta_D(S, a) = \bigcup_{j=1}^m ECLOSE(r_j)$

この方法によって得られた DFA D は NFA E と同等の言語を認識し, また NFA の元となる正規表現と同等である.

3.3 DFA からの実行バイナリ生成

DFA からの実行バイナリ生成には, 生成するコードについて 3 種類の実装を行った.

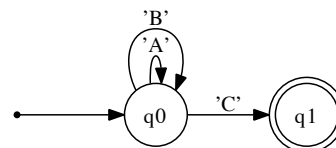
- (1) DFA \rightarrow Continuation based C \rightarrow GCC によるコンパイル
- (2) DFA \rightarrow C \rightarrow GCC によるコンパイル
- (3) DFA \rightarrow LLVM 中間表現 \rightarrow LLVM によるコンパイル

以下, Continuation based C, LLVM の説明と, それを利用した DFA からの実行バイナリ生成の方法を説明する.

3.3.1 Continuation based C

Continuation based C (以下 CbC) は, プログラミングの基本単位としてコードセグメントを持ち, コードセグメント間の軽量継続を基本とした C の下位言語である. 本研究室での先行研究により CbC コンパイラは, GNU C Compiler 上で実装されており⁹⁾, GCC の末尾再帰最適化を強制することで, 関数と同様の記述が可能で, かつ関数呼び出しに伴うリターンアドレスの保存や, スタックの成長のない, “引数付き goto”として継続を実装している. 本研究では gcc-4.5 上に実装された CbC コンパイラを用いた.

正規表現 “(A|B)*C” に対応する DFA と, DFA の各状態に対応する CbC のコードセグメントの例を載せる.



(A|B)*C

図 4 正規表現 “(A|B)*C” に対応する DFA

```

typedef unsigned char * UCHARP;
--code state_0(UCHARP beg, UCHARP buf, UCHARP end) {
switch(*buf++) {
case 65: /* 'A' */
goto state_0(beg, buf, end);
case 66: /* 'B' */
goto state_0(beg, buf, end);
case 67: /* 'C' */
goto state_1(beg, buf, end);
default: goto reject(beg, buf, end);
}
}

```

```

--code state_1(UCHARP beg, UCHARP buf, UCHARP end) {
goto accept(beg, buf, end);
}

```

図 4 の DFA に対応する CbC コードセグメント

DFA の遷移とは直接関係のない引数 (ファイル名やバッファへのポインタ等) が目立が, CbC では環境をコードセグメント間で引数として明示的に持ち運ぶ軽量継続を基盤としたプログラミングスタイルが望ましい. 今回コンパイラによって生成した CbC ソースコードでは, 大域変数は持たず, 必要な変数は全て引数に載せている. CbC の state_1, state_0 から呼ばれている accept, reject はそれぞれ受理と非受理を表す. accept ではテキスト行を出力して次の行へ, reject では次の文字へと処理を移すコードセグメントへ継続を行う.

生成した CbC ソースコードを, GCC 上に実装した CbC コンパイラによってコンパイルすることで実行バイナリを得る.

3.3.2 C

C による実装では, CbC のコードセグメントに代わり関数を用いて DFA を実装した.

```
typedef unsigned char * UCHARP;
void state_0(UCHARP beg, UCHARP buf, UCHARP end) {
    switch(*buf++) {
        case 65: /* 'A' */
            return state_0(beg, buf, end);
        case 66: /* 'B' */
            return state_0(beg, buf, end);
        case 67: /* 'C' */
            return state_1(beg, buf, end);
        default: return reject(beg, buf, end);
    }
}

void state_1(UCHARP beg, UCHARP buf, UCHARP end) {
    return accept(beg, buf, end);
}
```

図 4 の DFA に対応する C 関数

DFA による遷移を関数呼び出しで行っているため, 実行時のスタックの使用領域や, スタック操作によるオーバーヘッドが予想される.

3.3.3 LLVM

LLVM(Low Level Virtual Machine) は, さまざまなコード最適化/分析機能を備えた, モジュール単位で利用可能なコンパイラ基盤である⁶⁾.

CbC/C による実装では, DFA から CbC/C のソースコードに変換し, GCC によってコンパイルを行っている. LLVM による実装では, LLVM の中間表現である LLVM-IR を, 提供されている API を用いて直接操作を行い, コンパイルを経て実行バイナリを得る. Python から LLVM API の利用は, LLVM API の Python ラッパーである `llvm-py`^{*}を使用した.

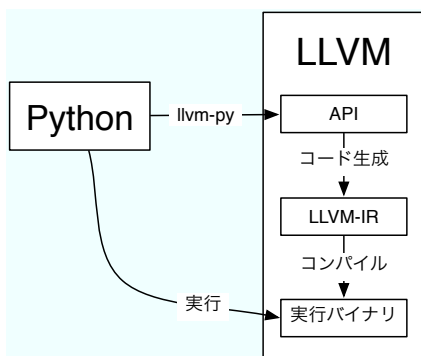


図 5 LLVM を用いた実装

LLVM による実装でも, C による実装と同様に, DFA の状態遷移を `switch` 文と関数呼び出しによって表現している.

4. grep

正規表現は, テキストのパターンをシンプルに記述できるという利点から, テキストファイルから, 任意のパターンにマッチするテキストを検索するなどの用途に使用される.

`grep` は, それを実現するソフトウェアの一つであり, 多くの実装が存在する. 引数として与えられたファイルから, 与えられた正規表現にマッチするテキストを含む行を出力する機能を持っている.

4.1 本実装により生成される `grep`

`grep` は, 与えられた正規表現にマッチするテキストを含む “行を出力する. ここで重要なのは, `grep` による正規表現マッチングは行単位で行なわれ, 行頭から行末までの完全一致ではなく, 基本的に部分一致で行なわれる^{**}.

部分一致を実現する方法として, 与えられた正規表現の先頭に正規表現 “.*” を付加して, 前方一致を行なう. ここで “.” は “任意の一文字” を表す特殊記号である.

正規表現 “(A|B)*C+” に対して, 本実装で生成される `grep` のコードを示す. 既に説明した通り, `grep` による部分一致を行なうために正規表現の先頭に “.*” を追加し, DFA に変換しコードを生成する.

```
void grep(UCHARP beg, UCHARP buf, UCHARP end) {
    state_1(beg, buf, end);
    return;
}
```

```
void state_1(UCHARP beg, UCHARP buf, UCHARP end) {
    switch(*buf++) {
        case 0: /* NUL */
            return reject(beg, buf, end);
        case 65: /* 'A' */
            return state_1(beg, buf, end);
        case 66: /* 'B' */
            return state_1(beg, buf, end);
        case 67: /* 'C' */
            return state_0(beg, buf, end);
        case 10: /* LF */
            return reject(beg, buf, end);
        case 13: /* CR */
            return reject(beg, buf, end);
        default: return state_1(beg, buf, end);
    }
}
```

^{*} <http://www.mdevan.org/llvm-py/>

^{**} GNU `grep` では完全一致を行なう “-x” オプションがある.

```
void state_0(UCHARP beg, UCHARP buf, UCHARP end) {
    return accept(beg, buf-1, end);
}
```

本実装では、“.”に対応する遷移を switch 文の default で行なっている。

本実装による grep では、検索対象となるテキストファイルを mmap でメモリにマップしている。メモリ上にマップされたテキストファイルを、(改行を含む)1つの巨大な文字列として関数 grep に与え、マッチングと出力を継続的なコード遷移で行う。

上記のコードでは、改行文字 (LF/CR) 及び終端文字 (NUL) に対して特殊な状態 *reject* に遷移しているが、この特殊状態 *reject* では行頭を表す引数 beg を更新、再度 grep に遷移する状態である。

```
void reject(UCHARP beg, UCHARP buf, UCHARP end) {
    if (buf >= end) return;
    beg = buf;
    return grep(beg, buf, end);
}
```

行頭を更新/保持することによって、マッチ行の出力時には行末のみを気にすれば良い。受理状態を表す *state_0* では、出力を行なう特殊状態 *accept* に直接遷移を行なう。ここで重要なのは、本来 *state_0* は与えられた正規表現の末尾 “C+” に対して “C” 字が入力された時点の状態であり、さらに “C” が入力文字として続く場合の遷移状態も本来保持している。しかし、本実装では、“最左最短一致” で出力を行なうよう実装を行なっているため、受理状態に遷移した時点で出力を行なう。

特殊状態 *accept* では、マッチングに成功した時点でのポインタ変数 buf を基準に、行末を memchr で探索し、出力を行なっている。

```
void accept(UCHARP beg, UCHARP buf, UCHARP end) {
    UCHARP ret = (UCHARP)memchr(buf, '\n', (buf-end));
    if (ret == NULL) {
        print_line(beg, end);
        return;
    }
    print_line(beg, ret);
    beg = buf = ret + 1;
    return grep(beg, buf, end);
}
```

出力を行なった時点で、まだファイルの終端に達していない場合は、出力した行の次の行頭から、再度 grep を呼び出している。このように、本実装では、状態遷移及び出力が末尾再帰的に行なわれており、コンパイラによって末尾呼び出しが適切に最適化された場合、スタック操作の無い高速な状態遷移を行なう実行バイナリに変換される。

Intel Compiler で本実装による grep を最適化オブ

ションを付加してコンパイル、実行したところ、末尾呼び出し最適化が行なわれず、実行して即座にスタックオーバーフローを起こしてしまった。C 言語では、末尾呼び出し最適化を明示的に記述することができないので、このような状態遷移ベースのプログラミングは Continuation based C による記述が望ましいといえる。

4.2 固定文字列フィルタリング

GNU grep を含むいくつかの grep 実装では、マッチングを高速化するために様々な高速化の工夫を行なっている。正規表現を DFA に変換してのマッチングや、固定文字列フィルタリングもその一つである。

与えられた正規表現が、固定文字列を含む場合、その固定文字列を含む行を探索し (フィルタリング)、続いて DFA によるマッチングを行なうことで、全行に対して DFA によるマッチングを行なう場合より高速なマッチングが可能となる。これは、探索する文字列が固定文字列の場合、DFA による探索よりも、Boyer-Moore 法等の固定文字列に特化した探索手法が高速なマッチングを行なえるからである。 n を被探索文字列の長さ、 m を探索文字列の長さとした場合、固定文字列探索を DFA で行なうと探索時間は単純に入力文字列 (被探索文字列) 長に比例するので $\Theta(n)$ となる。一方、Boyer-Moore 法等のアルゴリズムを用いて探索を行なった場合、最良で m 文字分探索をスキップできるので、探索時間は $\Omega(n/m), O(n)$ となる。

Boyer-Moore 法は、検索可能な固定文字列は 1 つのみであるが、これを複数の文字列に一般化したアルゴリズムとして、Commentz-Walter 法¹⁾ があり、GNU grep、及び本実装ではこれを採用している。

5. 評価

本実験で実装した正規表現エンジンの CbC/C/LLVM による三つの実装に対して、コンパイル時間及びマッチング時間の比較を行った。なお、GCC によるコンパイルには最適化オプション “-O3” を、LLVM のも同様の最適化オプションを用いてコンパイル/マッチングを行っている。

5.1 ベンチマーク: コンパイル

n 個の単語を正規表現の和集合演算 “|” で繋げたパターンに対し、各実装のコンパイル時間の比較を行った。

実験環境

- CPU : Core 2 Duo 950 @3.0GHz
- Memory : 16GB
- GCC : 4.5.0
- LLVM: 2.4

表 1 に結果を示す。

表 1 ベンチマーク:compile

n (単語数)	1	10	50	100
DFA 変換 [ms]	0.19	3.28	22.2	73.8
DFA の状態数	8	50	187	381
GCC-C [s]	0.34	0.78	2.18	4.27
GCC-CbC[s]	0.75	1.03	9.14	9.43
LLVM [s]	0.044	0.08	0.246	0.472

表 1 から、LLVM によるコンパイルが GCC に比べ 10 倍程高速に行われている。LLVM による実装では、API を通じて直接 LLVM の中間表現を操作することで、ファイル I/O やパース等のオーバーヘッドもない。

5.2 ベンチマーク: grep

マッチング時間の比較では、様々な正規表現を用いて比較を行った結果、3 つの実装でマッチング時間にあまり差が見られなかった。生成されるコードはコードセグメント/関数と、switch 文によるシンプルな実装であることから、コンパイルされたバイナリの性能にあまり差が出ていないものと思われる。

本実装の中で最もマッチングが高速だった GCC-C で生成した正規表現エンジンを用いて grep に相当するプログラムを実装し、実際にテキストファイルからのパターンマッチを行い、それぞれの評価を行う。本実装との比較対象として選択した grep について説明する。

5.3 GNU grep

GNU grep (<http://www.gnu.org/software/grep/>) は GNU が開発している OSS(Open Source Soft-

ware) で、非常にポピュラーなツールとして様々な OS に標準搭載され使用されている。GNU grep は POSIX 拡張正規表現に対応した egrep、固定文字列探索に対応した fgrep を内包しており、オプションで切り替えることができる。

C 言語による 1 万行程度の実装で、POSIX 規定の正規表現、後方参照等さまざまな多くの機能を有しながら、マッチング自体も高速である。

なお、今回の実験では GNU grep 2.6.3 及び GNU grep 2.5.4 を用いる。2 つのバージョンを用いる理由は、GNU grep 2.6 が 2010 年 3 月のリリースと比較的最近で、現行のほぼすべての OS では GNU grep 2.5.x が使用されているからである。さらに、GNU grep 2.6.x からは UTF-8 の対応が改善されている。

5.4 cgrep

cgrep (<http://sourceforge.net/projects/cgrep/>) は、Bill Tanenbaum が開発した OSS で、GNU grep に比べて 300-25%ほど高速なマッチングを行なう^{*}。C 言語による 2 万行程度の実装で、2004/6/30 にリリースされた cgrep 8.15 以降開発が止まっている。また、マルチバイト文字のサポートなどは行っていない。

5.5 Google RE2

Google RE2 (<http://code.google.com/p/re2/>) は Google の Russ Cox を中心に開発されている OSS な正規表現エンジンライブラリで、正規表現マッチングを行なう複数の API を利用できる。

C++ による 2 万行程度の実装で、UTF-8 に対応している。

本実験では、Google RE2 の API である FindAndConsume を使用して、grep クローンを実装した。FindAndConsume はマッチング対象文字列と正規表現、マッチした文字列を格納するポインタを引数にとり、マッチングの開始部分を内部的に更新していく API である。以下が、RE2 を用いた grep 実装のメインルーチンである。

```
void grep(string regexp, UCHARP beg, UCHARP end) {
    RE2 re(".*"+regexp+".*");
    re.ok(); // compile;
    string contents = string((const char *)beg, (end - beg));
    StringPiece input(contents);
    string word; // container of printable line.
    while (RE2::FindAndConsume(&input, re, &word)) {
        cout << word << "\n"; //print
    }
    return;
}
```

実験環境

^{*} cgrep の manpage より抜粋

- CPU : Core i7 950 @3.0GHz
 - Memory : 16GB
 - GCC : 4.4.1
 - Text : Wikipedia 日本語版全記事
(XML, UTF-8, 4.7GB, 8000 万行)
- 表 2 に結果を示す。

表 2 ベンチマーク:grep

テストケース	fixed-string	complex-regex
本実装 (GCC-C)[s] (コンパイル [s])	1.69 0.20	4.51 0.41
cgrep 8.15 [s]	1.85	6.48
GNU grep 2.6.3[s]	2.93	16.08
GNU grep 2.5.4[s]	2.96	188.51
Google RE2 [s]	30.10	16.92

以下に、それぞれのテストケースのパターン、grep にマッチした行数、および考察を記す。なお、ここで扱う正規表現の“複雑さ”とは、DFA に変換した時点の状態数、遷移規則の多さを基準としている。

fixed-string 固定文字列によるマッチングパターン : “Wikipedia”
マッチ行数: 348936 行

GNU grep では、与えられたパターン内に、確実に含まれる文字列 (固定文字列) が存在する場合は、Boyer-Moore 法等の高速な文字列探索アルゴリズムを用いてフィルタをかけることで、DFA によるマッチングを減らし、高速化している。本実装の grep でも、同様に固定文字列フィルタによる高速化を行っているが、単純な固定文字列の検索でも、コンパイルすることによる高速化が結果に出ている。

complex-regex 複雑な正規表現でのマッチングパターン :“(Python|Perl|Pascal|Prolog|PHP|Ruby|Haskell|Lisp|Scheme)”
マッチ行数: 15030 行

このパターンは、9 個の単語を和集合演算 “|” で並べたもので、確実に含まれる文字列は存在しないが、前述の Commentz-Walter 法によるフィルタリングが適用できる。GNU grep、cgrep 及び本実装において fixed-string のケースよりはマッチングに時間が

さらに、GNU grep 2.5.4 は 190 秒と、本実装及び GNU grep 2.6.3 に対して非常に低速な結果となっているが、これは後述する mbtowc(3) によるマルチバイト文字の変換処理によるオーバーヘッドによるものである。

2 つのテストケースの結果を見てみると、本実装はそれぞれ GNU grep と比べて高速な結果となってお

り、この規模のファイルに対する grep の場合、コンパイル時間はマッチングにかかる時間と比べて無視できるほど短い時間であることが分かる。

6. 本実装の特徴

本実験によって実装された正規表現評価器の特徴を、GNU grep との比較をさみながら説明する。

6.1 正規表現からの動的なコード生成

本実装による一番の特徴として、正規表現から変換を行うことで得られる等価な DFA を、C/CbC/LLVM に変換しコンパイルすることで正規表現に応じた実行バイナリを生成することが挙げられる。またコンパイラ理論におけるさまざまな最適化が期待される。grep による検索のような、与えられるパターン (正規表現) に対してマッチ対象のテキストファイルが十分に大きい場合、正規表現のコンパイルにかかる時間はマッチングにかかる時間によって隠される。

GNU grep では、本実装と同様に DFA ベースのマッチングを行うが、DFA の各状態は構造体によって表され、状態遷移は各状態毎に保持している遷移先ポイントによる配列を、1Byte 単位でテーブルルックアップを行うことで実装されている。

6.2 UTF-8 対応

本実装は、マルチバイト文字の代表的な符号化方式である UTF-8 に内部的に対応しており、正規表現の演算は 1Byte 単位ではなく、1 文字単位で適用される。マルチバイト文字を含む正規表現のサンプルとして、“(あ|い)*う” を DFA に変換した図 6 を載せる。図における ‘\x’ に始まる文字は 16 進数表記で、‘\x82’ は 16 進数 82 を表す。

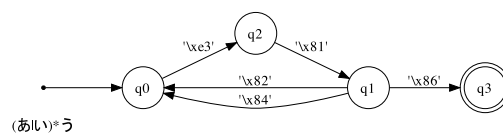


図 6 正規表現“(あ|い)*う”に対応する DFA

GNU grep 2.5.x では、マルチバイト文字に対応しているものの、プログラム内部で libc mbtowc(3) を用いて固定サイズであるワイド文字に変換して処理を行っており、テストケース complex-regex ではそのオーバーヘッドが顕著に現れている。2010 年 3 月にリリースされた GNU grep 2.6 から、UTF-8 に対して本実装と同様に内部的に対応することで、mbtowc(3) による変換コストが無くなっている。

6.3 柔軟な実装

本実験で実装した正規表現評価器は、Python によって実装されており、全体で 3000 行程度の軽量なプログラムとなっている。今回比較対象として利用した GNU grep, cgrep, Google RE2 はいずれも 12 万行の規模のソフトウェアに比べ機能の追加やリファクタリングが用意であり、本実装の優れた点と言えるだろう。

ここで重要なのは、本実装から動的に生成されるコードも、コードセグメント/関数と switch を基準としたシンプルな記述で高い可読性を持ちつつ、細かい最適化を GCC/LLVM の最適化技術を利用することで実行速度も GNU grep に比べ 2-4 倍高速であり、他の grep 実装よりも高い性能であることが実験結果から分かった。

7. 今後の課題

本研究では、現段階で正規表現をコードセグメント/関数による状態遷移を行うコードに変換する手法で正規表現エンジン及び grep クローンを実装し、他の grep 実装との比較を行い、良好な結果が得られた。

今後の課題として、正規表現マッチングのデータ並列な並列アルゴリズムの実装を考えており、Cell B.E. 等の並列環境で実行可能なコードの生成を目指している。

参考文献

- 1) Commentz-Walter, B: A String Matching Algorithm Fast on the Average, Proc. 6th International Colloquium on Automata, Languages, and Programming (1979)
- 2) Cox, R : Regular Expression Matching Can Be Simple And Fast, <http://swtch.com/~rsc/regexp/regexp1.html> (2007)
- 3) Cox, R : Regular Expression Matching: the Virtual Machine Approach, <http://swtch.com/~rsc/regexp/regexp2.html> (2009)
- 4) Cox, R : Regular Expression Matching in the Wild, <http://swtch.com/~rsc/regexp/regexp3.html> (2010)
- 5) Hopcroft, J, E. Motowani, R. Ullman, J : オートマトン言語理論計算論 I (第二版), pp.39-90.
- 6) Lattner, Chris. Adve, Vikram : The LLVM Compiler Framework and Infrastructure, <http://llvm.org/pubs/2004-09-22-LCPCLLVMTutorial.pdf> (2004)
- 7) 新屋 良磨 : 動的なコード生成を用いた正規表現マッチャの実装, 日本ソフトウェア科学会第 27 回大会 (2010)
- 8) Thompson, K : Regular Expression Search Algorithm, Communications of the ACM 11(6) (1968).
- 9) 与儀 健人, 河野 真治 : Continuation based C コンパイラの GCC-4.2 による実装, 情報処理学会システムソフトウェアとオペレーティング・システム研究会 (OS) (2008)

(平成 2010 年 12 月 03 日受付)

(平成 0 年 0 月 0 日採録)

新屋 良磨 (正会員)

1988 年生. 2007 年琉球大学工学部情報工学科入学

河野 真治 (正会員)

1959 年生. 1989 年東京大学大学院情報工学課程修了 (工学博士) 同年 Sony Computer Science Laboratory, Inc. 入社. 1996 年より琉球大学工学部准教授工学博士. ACM

会員.