

## ClamAV の AC 法をダブル配列により拡張したアンチウイルスエンジンの提案

松浦 寛生<sup>†</sup> 蔵満 琢麻<sup>†</sup> 望月 久稔<sup>†</sup><sup>†</sup>大阪教育大学

## 1. はじめに

コンピュータウイルスによる被害が増加し、対策として様々なアンチウイルスソフトが開発されている。アンチウイルスソフトには、ウイルスパターンを登録して対象のデータを照合するソフトや、ヒューリスティックな学習によりウイルスを検知するソフト、プログラムの動作を監視し、動作の異常を検知するソフトなどがある。オープンソースのアンチウイルスソフト ClamAntiVirus(以下、ClamAV)は、登録したウイルスパターンと対象ファイルを照合するソフトである[1]。

本論文では、ClamAV の照合に用いられる AC 法[2]をダブル配列[3]により拡張し、照合手法に合わせてアンチウイルスエンジンを最適化したシステムを提案する。

## 2. アンチウイルスソフト ClamAV

## 2.1 ウイルスパターンの照合方法

AC 法はマシン AC と呼ばれる複数パターンマッチングマシンを用いて照合する手法で、ClamAV においても使用される。

ClamAV に登録するウイルスパターンには、部分的に不定なパターンを表現するためにワイルドカードを用いる。しかし、不定なパターンはマシン AC に登録することができない。そのため ClamAV は、ウイルスパターンにおいて不定ではない部分の 3 バイト以下をパターンキーとしてマシン AC に登録し、各節点において検出するキーに該当するウイルスパターンを各節点の Output 集合として登録する。ワイルドカードは、キー検出時に該当するウイルスパターンと対象データを 1 バイトずつ照合することで処理される。以下、Output 集合として登録したウイルスパターンを V パターンと呼ぶ。V パターンの照合はキー検出時の処理であるため、マシン AC における対象データの照合ポインタは進まない。よって、V パターン照合回数の増加により ClamAV における照合速度は低下する。

また、ClamAV におけるマシン AC は配列構造により Goto 関数を実現しているため、マシン AC における各節点にはすべての遷移種に対応した遷移先を保持する空間が必要となる。そのため、すべてのウイルスパターンをマシン AC に登録すると、節点数の増加により照合時に必要な空間が膨大になる。そこで、ClamAV はワイルドカードを用いないウイルスパターンをハッシュ法により処理し、マシン AC に登録するキー数を制限する。これにより、ClamAV は 1 つの照合対象データに対して AC 法とハッシュ法による 2 度の照合を要する。

## 2.2 ウイルスデータベースの構築

ClamAV は、マシン AC における節点間のリンクや、ハッシュ法における節点間のリンクをポインタ管理により実現する。そのため、ClamAV はシステムを起動する度にファイルからウイルスパターンを読み込み、ウイルスデータベースを構築する必要がある。

The Improvement of ClamAV using MachineAC extended with Double-Array Structure

<sup>†</sup> Nobutaka MATSUURA, Takuma KURAMITSU,  
Hisatoshi MOCHIZUKI  
Osaka Kyoiku University

## 3. ダブル配列により拡張した AC 法

パターンの効率的な照合を実現してシステムの性能を向上させるために、マシン AC における Goto 関数をダブル配列により実現し、Goto 関数に決定性をもたせたマシン AC を提案する。

## 3.1 データ構造

ダブル配列は配列 Base, Check を用いてトライ構造を実現する手法で、節点  $s$  から遷移種  $a$  による遷移先節点  $t$  を式(1)により定義できる[4]。以下、節点  $s$  における Base 値、Check 値をそれぞれ  $B[s]$ ,  $C[s]$  とする。

$$\begin{cases} t \leftarrow B[s] + a \\ C[t] = a \end{cases} \quad (1)$$

提案マシンにおける Goto 関数、Output 関数を以下のように定義する。ただし、Output 関数内で用いる '#' はパターンに使用する遷移種の最大値 +1 として定義し、Goto 関数と Output 関数を拡張するためにダミー節点を付加する。

ダブル配列における遷移は式(1)により定義できるため、Base 値を写像することで節点の遷移情報を写像できる。よって、Goto 遷移が存在しない節点においては Failure 遷移先の Base 値を写像することで遷移を定義する。ここで、式(1)による遷移先が未定義である場合、Goto 関数により遷移すべき節点  $t$  をトライ構造上の深さが 1、もしくは 2 である節点に限定し、Goto 関数に決定性をもたせることができる。式(1)による遷移先が未定義である場合、節点  $s$  が初期節点であるか、遷移すべき節点  $t$  のトライ構造上の深さが 3 以上であるときに Goto 関数を拡張するためのダミー節点を作成する。ダミー節点の節点番号  $d$  を  $B[s] + a$  とし、 $B[d]$  に遷移すべき節点  $t$  の Base 値を写像する。ここで、ダミー節点  $d$  は節点  $t$  の遷移情報をもつため、節点  $s$  から節点  $t$  への遷移を擬似的に定義できる。そこで、Goto 関数を以下のように拡張することで、提案マシンは Failure 関数を用いることなく対象データを照合することができる。

関数 Goto(節点  $s$ , 遷移種  $a$ ) 戻り値：遷移先節点

式(1)に示す節点  $t$  が定義されていれば  $t$  を返す。未定義であれば  $t \leftarrow B[B[\text{初期節点}]] + C[s] + a$  とし、 $C[t] = a$  であれば  $t$  を返し、そうでなければ  $B[\text{初期節点}] + a$  を返す。

関数 Output(節点  $s$ ) 戻り値：Output 集合のインデクス

ダミー節点  $s_0 \leftarrow s + \#$  とし、 $C[s_0] = \#$  であれば Output 集合のインデクス  $B[s_0]$  を返し、そうでなければ Output 集合が存在しないことを示す False を返す。

例1：図 1 に、パターン集合  $P = \{\text{"ACBAC"}, \text{"BA"}, \text{"BAB"}\}$  を登録した提案マシンを示す。ここで、遷移種 'A', 'B', 'C', '#' の内部表現値をそれぞれ 0, 1, 2, 3 とし、Output 集合  $\{\text{"ACBAC"}\}, \{\text{"BA"}\}, \{\text{"BAB"}\}$  のインデクスをそれぞれ 1, 2, 3 とする。また、図 1 のマシン AC において、節点 1 を初期節点とし、Base 値、Check 値がともに空欄である節点を未使用節点とする。

図 1 における節点 4, 14 は Goto 関数を拡張するためのダミー節点に該当し、それぞれ初期節点 1, 節点 13 における遷移情報

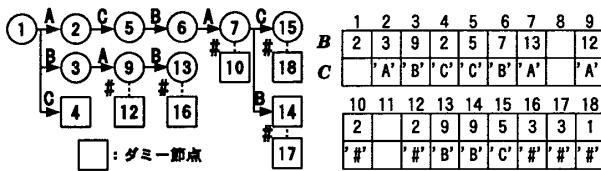


図 1 パターン集合 P を登録した提案マシン

をもつ。これにより節点 1 から節点 1, 節点 7 から節点 13 への Goto 関数による遷移を擬似的に定義できる。同様に、節点 10, 12, 16, 17, 18 は Output 関数を拡張するためのダミー節点に該当し、それぞれ節点 7, 9, 13, 14, 15 における Output 集合のインデックスを管理する。これにより Output 集合のインデックスを、ダミー節点における Base 値を参照することで得られる。

図 1 のマシン AC を用いて、対象データ “ACBABC” を照合する場合について考える。まず初期節点 1において、Goto(1, ‘A’)=2 より、節点 2 へ遷移する。ここで、Output(2)=False より、節点 2 において検出する Output 集合は存在しない。同様に、節点 5, 6, 7 と遷移し、 $C[7+\#(10)] = \#$  より、 $B[10] = 2$  に該当する Output 集合 {“BA”} を検出する。次に、Goto(7, ‘B’)=14 より節点 14 へ遷移し、Output(14)=3 に該当する Output 集合 {“BAB”} を検出する。従来のマシン AC ではダミー節点 14 が存在しないため、節点 7 から Failure 関数により節点 9 に遷移する必要がある。その後、Goto(14, ‘C’)=4 より、節点 4 へ遷移し、対象データの末尾まで照合したので終了する。

#### 4. AC 法の拡張に伴ったアンチウイルスエンジンの最適化

##### 4.1 AC 法のみによる実装

提案マシンは、1 節点あたりに必要な空間が Base 値、Check 値を保持する空間のみとなるため、ClamAV のマシン AC よりも照合時に必要な空間が小さい。そこで、すべてのウイルスパターンのキーを提案マシンに登録し、AC 法のみにより照合するアンチウイルスエンジンを提案する。また、提案システムに登録するウイルスパターンのキー長を伸ばすことで Output 集合を保持する節点を分散させ、V パターン照合回数を抑制する。これにより提案システムは 1 度の照合でウイルスを検知でき、V パターン照合回数を抑制できるため高速である。しかし、キー長を伸ばすことにより照合時に必要な空間が増加するため、システムを使う環境に合わせてキー長を決定する必要がある。

##### 4.2 ウイルスデータベースの復元

ダブル配列はカーソル管理によりデータ構造を実現する。そのため、1 度作成したウイルスデータベースをファイルに保存することで、データベースをファイルから直接コア上に復元できる。これによりデータベースの構築時間を短縮できる。

#### 5. 実験評価

提案システムの有効性を示すため、ClamAV との比較実験を Intel Pentium Dual 1.60GHz, Fedora 7 上で行った。実験では、2007 年 12 月 26 日時点のウイルスパターン [1] をシステムに登録し、ファイル数 22、総バイト数 9.38MB の EXE ファイルを照合した。実験における、システム全体の稼働時間、ウイルスパターンの照合時間、ウイルスデータベースの読み込み時間、ウイルスパターンの照合時に必要な空間、マシン AC の節点数、V パターンの照合回数を表 1 に示す。また、実験において、提案システムに

表 1 実験結果

	ClamAV	提案システム
稼働時間 (s)	46.07	35.05
照合時間 (s)	43.09	33.06
読み込み時間 (s)	2.81	1.89
使用空間 (MB)	26.27	25.76
マシン AC の節点数	9,764	876,976
V パターン照合回数	601,457,898	509,169,419

登録するウイルスパターンのキー長を 5 とした。

表 1 に示すように提案システムはマシン AC の Output 関数により発生する V パターンの照合回数が ClamAV よりも少ない。これは、マシン AC に登録するウイルスパターンのキー長を伸ばしたことにより、Output 集合をもつ節点が分散し、各節点において発生する照合回数を抑制できるためである。さらに、ClamAV がハッシュ法と AC 法の併用により、1 つの照合対象データに 2 度の照合を要するのに対し、提案システムは AC 法による 1 度の照合しか行わないため、表 1 に示すように、ClamAV に対し約 77% の時間で照合可能となる。

また表 1 より、提案システムにおけるウイルスデータベースの読み込み時間は ClamAV の約 67% であり、提案システムはデータベースの読み込みにおいても ClamAV より高速である。これは、提案システムにおけるマシン AC をカーソル管理により実現したことで、提案システムが構築済みのデータベースをファイルから直接読めるためである。

提案システムは AC 法を拡張して実装したことで、表 1 に示すように ClamAV よりマシン AC の節点数が増加する。しかし、提案システムのマシン AC における 1 節点あたりに必要な空間は ClamAV よりもはるかに小さいため、キー長 5 における提案システムの照合時に必要な空間は、キー長 3 における ClamAV とほぼ同等であった。

以上より、提案システムは ClamAV と比較して、照合時に必要な空間を増加させることなくシステム全体の稼働時間を短縮できることを示した。

今後の課題として、システムにヒューリスティックな検出機能を追加し、未知のウイルスにも対応したエンジンに拡張することがあげられる。

#### 参考文献

- [1] Clam antivirus. <http://www.clamav.org/>.
- [2] M.J.Corasick A.V.Aho. Efficient string matching: An aid to bibliographic search. *Communications of the ACM*, Vol. 18, No. 6, pp. 333–340, 1975.
- [3] J.Aoe. An efficient digital search algorithm by using a double-array structure. *IEEE Trans. on Software Engineering*, Vol. 15, No. 9, pp. 1066–1077, 1989.
- [4] 矢田晋, 森田和宏, 泷田正雄, 平石亘, 青江順一. ダブル配列におけるキャッシュの効率化. *FIT2006*, pp. 71–72, 2006.