

学位論文最速文法マスター

A DISSERTATION
SUBMITTED TO THE DEPARTMENT OF PHYSICS,
ELECTRICAL AND COMPUTER ENGINEERING
AND THE COMMITTEE ON GRADUATE STUDIES
OF YOKOHAMA NATIONAL UNIVERSITY
IN PARTIAL FULFILLMENT OF THE REQUIREMENTS
FOR THE DEGREE OF
MASTER OF ENGINEERING

1264904 MAI MAI CUONG

February 2016

© Copyright by MAI MAI CUONG 2016
All Rights Reserved

I certify that I have read this dissertation and that, in my opinion, it is fully adequate in scope and quality as a dissertation for the degree of Master of Engineering.

(KURAMITSU Kimio) Principal Adviser

I certify that I have read this dissertation and that, in my opinion, it is fully adequate in scope and quality as a dissertation for the degree of Master of Engineering.

(HAMAGAMI Tomoki)

I certify that I have read this dissertation and that, in my opinion, it is fully adequate in scope and quality as a dissertation for the degree of Master of Engineering.

(HOGE Fuga)

Approved for the University Committee on Graduate Studies

Abstract

Acknowledgements

本論文は私が横浜国立大学大学院工学府物理情報工学専攻電気電子ネットワークコースに在籍中の研究成果をまとめたものである。本論文を執筆するにあたり、横浜国立大学工学部電子情報工学科准教授である倉光君郎先生には、指導教官として研究の方針をご指導頂きました。ここに深く感謝の意を表します。

(stab)

加えて、私の学生生活を長きにわたって見守って下さった私の家族に深く感謝致します。

Contents

List of Figures

Chapter 1

はじめに

近年、IoTの発展につれ、様々な機器への組み込みやすさやより高度なパケット処理が必要になった。また、ビッグデータ時代になった今は、より高速な構文解析技法が求められている。そこで、組み込みやすさと高速な処理の面から、FPGA() というデバイスが注目されている。FPGA を用いて、構文解析器を実現する研究がいくつかある。例えば、研究1では正規表現を使って、FPGA上で構文解析器を作る。正規表現からFPGA合成に必要なHDLファイル、ネットリスト又はコンフィグレーションデータに変換する。しかし、パーサごとに回路を合成しなければならず、手間がかかる。また、開発者にはある程度FPGAの知識が必要である。また、研究2で文脈自由文法を使ってFPGAでパーサを実現する研究がある。こちらはサブ回路を作ってたくさん並ぶ。ただし、 n 個のサブ回路が $n+1$ 文字しか受理できず、性能的にはよいとは言えない。

FPGAで構文解析器を実現するのは非常に手間がかかる。例えば、ソフトウェア開発と異なり、FPGAの回路では、遅延が発生しているため、開発者は常に遅延時間をいしきしなければならない。また、FPGA開発の際には、合理合成、配線などという段階あり、検証時間が大きい。さらに、新しいプロトコルが続々誕生されている今、それに対応するには、構文解析器を書き換える必要がある。

そこで、本研究の目的は自動生成可能、高速な処理、そして組み込みやすいパーサを実現する。方法としては、解析表現文法からFPGA上でパーサが自動生成する。ユーザはFPGAの知識が必要なく、文法を定義するpegファイルとテキストを入力するだけで、結果が得られる。

本論文の構成は次の通りである。

第1章：初めに第2章：FPGA及び解析表現文法PEG第3章：実装方針第4章：デザイン第5章：実装第6章：実験第7章：関連研究第8章：結論

Chapter 2

関連知識

この章では、FPGA及びPEGについて説明する。

2.1 FPGA

FPGAとは、ユーザー側で回路を変更することができるデバイスのことである。FPGA（Field Programmable Gate Array）とは、デジタル回路を構築する大規模集積回路の一種である。通常集積回路の場合、出荷された回路を変更することができないが、FPGAの場合、出荷された回路も使用者側で任意の論理機能を実装できる。FPGAは並列度の高い電子回路や機能の専門化によって機能向上した電子回路を容易に構成できる。FPGAは並列処理に擦れており、並列化によって高速化することが可能である。また、ユーザー側が自由に回路を変更することが可能であるため、ある問題に特化した専用回路を作り、性能を向上することができる。

さらに、稼働中の機器のハードウェアも変更できるため、処理の対象や内容に合わせて、状況に応じた回路に書き換えることで、最小限の回路で、多様な処理を実行できる。

一般にFPGA開発の際、ハードウェア記述言語（HDL）が用いられる。本研究では、HDLの一種であるVHDLを使用する。

2.2 PEG

構文解析とは、定義された文法に従ってテキストの構造を解析することである。本研究では、構文解析を実現するために、形式文法を使って、FPGA上でパーサを自動生成する。形式文法はたくさんの種類がある。例えば、解析表現文法（PEG）や文脈自由文法（CFG）などがある。

本研究では、パーサを構築するために、解析表現文法を使用する。PEGは $A_i = e$ というルール集合である。解析表現は表 () にある値と演算子を組み合わせた式である。

PEGのメリットは書きやすく、強力な表現力をもっている。また、Pactrack を使って線形時間で処理することができる。

Chapter 3

実装方針

本研究では、Virtual Machine を採用している。PEG ファイルから、実行すべき命令列を生成する。命令セットは以下となる。命令では、文字の操作の命令と制御用の命令がある。

文字操作の命令は BYTE, SET, ANY がある。これらの命令は実際に文法定義ファイルの文字とテキストからの文字を比較し、文字を消費するかを判断する。BYTE は文字リテラルを処理するための命令である。例えば、 $A = 'a'$ の PEG である場合、`[BYTE 'a']` という命令が生成される。この命令は、現在評価しているテキストの文字が 'a' であれば成功でそうでなければ失敗である。SET は文字クラスを処理するための命令である。例えば、 $A = [1-9]$ の場合、`[SET [1-9]]` という命令が生成される。ANY は任意の文字を処理するための命令である。例えば、 $A = .$ の場合、`[ANY]` という命令がある。非終端記号は CALL 命令を使って、非終端記号の定義されたところにジャンプする。ジャンプする前に、現在命令の次の命令を Return stack に push する。RET という命令の時、Return stack からジャンプ先を取得する。例えば、以下の PEG ファイルはこのような命令列が生成される。 $Add = Value (+/-) Value$ $Value = [1-9]$

生成された命令列 : `Add L1 Call L5 L2 Set [+/-] L3 Call L5 L4 Ret Value L5 Set[1-9] L6 Ret`

まず非終端記号があるため、その非終端記号 (L5) へジャンプする。L 5 にジャンプする前に、L 2 を Return stack に push される。ここで L5 が成功したら、L6 に進み、L6 は RET であるため、Return stack からジャンプ先を pop される。先ほど L2

を Return stack に push されたため、L2 にジャンプされる。次に L2 が成功したら、L3 に進む。L3 で Value を呼び出し、L4 を Return stack に push される。Value の命令が終ったら、L4 に戻る。L4 で戻る先がないため、プログラムを終了する。どこから失敗した場合、プログラムが終了される。

A L T はジャンプ先を Fail stack に push する命令である。PEG では、Option, R, 選択という演算子があり、どこかの命令が失敗しても、処理が続けることがある。どこかが失敗した時、Fail stack からジャンプ先を取得する。A L T の使い道はグルーピング、非終端記号、選択を処理する時に使う。具体的に下記に述べる。

Option を処理するため、2つの方法がある。対象となる Option は上記の基本要素であれば、その命令と組み合わせて一つの命令として実行する。例えば、'a'? である場合、[O B Y T E 'a'] という Option と B Y T E を組み合わせた命令が生成される。対象となる Option はグルーピングまたは非終端記号の場合、A L T という命令を使って飛び先を stack に push する。例えば、A = ('a' 'b')? 'c' の場合、'a' を実行する前に、ALT で 'c' を評価するための位置を stack に push する。この場合、以下の命令列が生成される。

```
L1 ALT L4 L2 BYTE 'a' L3 BYTE 'b' L4 BYTE 'c'
```

'a', 'b' でマッチした場合、L2,L3,L4 の順に実行される。一方、例えば 'a' で失敗した場合、Fail が発生し、stack から飛び先を取得する。先 ALT で L4 を stack に push されたため、L4 が pop されて、次に実行する命令が L4 となる。O p t i o n 対象が非終端記号の場合、ALT の後に C A L L が実行される。このように、ALT によって Option を処理することができる。

0 回以上の場合も処理が 2 つに分かれる。対象が B Y T E, S E T, A N Y の場合、R B Y T E, R S E T, R A N Y として組み合わせた命令を実行する。対象はグルーピング又は非終端記号の場合、A L T 及び J U M P を使って処理する。まず A L T でこの処理が終ったときのジャンプ先を Fail Stack に push する。次に、グルーピング・非終端記号の要素を順番に実行し、成功したら、J U M P でグルーピング・非終端記号の最初の要素に戻る。この処理を F a i l が発生するまで、繰り返される。F a i l が発生したら、先 A L T で push した Fail stack から、ジャンプ先を取得する。例えば、以下の P E G はこのような命令列が生成される。Value = '1' (' + ' '2') * ' + ' '3'

```
L1 BYTE '1' L2 ALT L6 L3 BYTE '+' L4 BYTE '2' L5 JUMP L3 L6 BYTE '+'
L7 BYTE '3' L8 RET
```

まず、A L TでL 6をFail stackにpushする。つぎに順番に1 L3, L4を実行し、終わったらJUMP命令でL3に戻る。L3, L4のどこかで失敗するまで、繰り返しが続く。一方、失敗すると、Fail stackにL6をpopされて、L 6が実行される。1回以上の場合、基本命令と0回以上にみなして実行する。例えば、A = 'a'+ の場合は、[BYTE 'a'] 及び [RBYTE 'a'] のように命令が生成される。

Choiseの場合、backtrackを削減するために、選択枝の最初の文字によって、最初にどの選択にするかを定める。この処理をする命令はFIRSTという。複数の選択枝の最初の文字が同じである場合、各選択枝を評価する前に、この選択枝が失敗である場合、別の選択枝に移すためのジャンプ先をALT命令でstackに保存する。例えば、以下のP E Gはこのような命令が生成される。Value = ('1' '1')/('1' '2')/'3'

```
L1 FIRST '1' -j L3 '3' -j L8 default -j L2 L2 FAIL L3 BYTE '1' L4 ALT L7 L5
BYTE '1' L6 RET L7 BYTE '2' jump L6 L8 BYTE '3' jump L6
```

まず、最初の文字によって、jump先が分かれている。最初の文字が'1'であればL3に、'3'であればL8に、そうでなければL2に移ってF a i処理が発生する。L 3にジャンプした場合、まず'1'を消費する。この場合、選択枝が'1' '1'又は'1' '2'の2つがあるため、'1' '1'の選択枝を評価する前に、この選択枝が失敗した場合のジャンプ先をA L T命令でFail stackにpushする。もし、この選択枝が成功した場合、L6のRETでプログラムを終了させる。一方、失敗した場合、選択枝'1' '2'を評価し、つまりL7が実行される。

否定先読みの場合、処理が2つに分かれる。対象がB Y T E, S E T, A N Yの場合、N B Y T E, N S E T, N A N Yとして組み合わせた命令を実行する。対象はグルーピング又は非終端記号の場合、まずグルーピング・非終端記号を評価する前に、失敗した時のジャンプ先をFail stackにpushする。その後、グルーピング・非終端記号の要素を順序に評価する。どこかで失敗が発生したら、A L Tで保存したジャンプ先にジャンプする。この場合は次の解析表現に進むことになる。ただし、グルーピング・非終端記号が成功した場合、つまり否定先読みが失敗した場合、もしFail処理を発生させるとしたら、A L Tで保存されたジャンプ先に飛ぶことになるため、次の解析表現に進んでしまう。そのため、まずSUCC命令でA L Tで保存したジャンプ先を消してから、Fail処理を発生させる。例えば、以下のP E Gはこのような命令列が生成される。Value = !('1'+) [1-9]+ L1 ALT L6 L2 BYTE '1' L3 RBYTE '1' L4 SUCC L5 FAIL L6 SET [1-9] L7 RSET [1-9] L8 RET

肯定先読みの場合、対象の表現を順序に評価することになるが、文字を消費しないため、評価する前に、POS命令で今の文字の位置を保存しておく。評価が終わったら、BACK命令で、保存された位置に戻る。例えば以下のPEGはこのような命令列が生成される。

```
L1 POS L2 BYTE '1' L3 BACK L4 SET [1-9] L5 RSET [1-9] L6 RET
```

シーケンスは連続な命令で実行される。

FPGAに情報を渡す前に、前処理としてPEGファイルから命令列を生成する。この処理は本研究室が開発しているNezライブラリを使う。FPGAは命令列を受け取り、テキストを解析する。

Chapter 4

設計

まず、各命令に対応するのサブ回路を作る。制御部の信号によってどのサブ回路が実行されるのか決められる。

4.1 命令用回路

4.1.1 文字操作命令

文字操作命令は BYTE, SET, ANY がある。

BYTE 命令は、テキストの文字と形式文法の文字が一致するかどうかを判断する。BYTE 命令用回路の入力は TRG, PEG からの文字、テキストの文字がある。一致する場合、match 信号が'1' となり、逆に一致しない場合は fail 信号が'1' となる。Trg がある場合だけ、作動する。T r g がいないときは match 及び fail が'0' となる。match が'1' になった場合、次の命令に進む処理をする。fail が'1' の場合、Alt 命令で保存されたスタックからとび先をもらう。スタックにデータが存在しない場合、プログラムを終了し、パース失敗となる。

ANY 命令は、1 文字を消費するという命令である。ANY 命令用の回路の入力は TRG で、出力は nextchar 信号である。TRG がある場合だけ、nextchar が'1' となり、文字を消費する処理を始める。

SET 命令用回路の入力は、TRG, PEG からの 2 の文字、テキストテキストの文字がある。文字の範囲か、文字の選択かによって、比較条件が異なる。マッチした

場合、match が'1' となり、マッチしない場合、fail が'1' となる。

4.1.2 組み合わせ命令

組み合わせ命令は Option, 0 個以上、否定先読みを BYTE, SET, ANY を組み合わせた命令である。

Option の場合、0 個マッチしてもよいため、TRG がある場合、match が'1' となり、この回路から fail 信号が発生されない。しかし、文字を消費するかを決めるため、1 個マッチした場合、nextchar が'1' となり、文字を消費する。

0 個以上の場合、マッチしたら、nextchar 信号が'1' となり、文字を消費する。また、nextchar 信号が TRG とつながっており、マッチした場合、命令をもう一度実行する (図)。ただし、新たな文字をロードするには、1 クロックが必要であるため、nextchar を 1 クロック遅らせて TRG に入る。このように、マッチしなくなるまで、テキストを消費しながら、回路を動作する。一方、マッチされなくなったら、nextcmd 信号が'1' になり、次の命令に進む。

否定先読みの場合、条件を逆にするだけである。

4.1.3 その他の命令

CALL 命令用回路は、TRG があるとき、まず Return stack に命令がもっている保存用の命令アドレスを追加し、Return stack の先頭をインクリメントする。同時に命令が持っているジャンプ用の命令アドレスをプログラムレジスタにロードする。

ALT 命令用回路は、TRG があるとき、Fail stack に命令が持っている保存用アドレスを追加し、Fail stack の先頭をインクリメントする。

SUCC 命令用回路は、TRG があるとき、Fail stack から先頭の命令アドレスを消し、Fail stack の先頭をデクリメントする。

FAIL 命令用回路は、TRG があるとき、fail を'1' となり、Fail 処理を発生させる。

4.2 制御部

4.3 制御部

実効ステップは以下となる。ステップ1:PRの値がMARに転送ステップ2:メモリからデータをIRに転送ステップ3:命令デコード及ぶ命令実行

ステップ3が終わったら、またステップ1に戻

Chapter 5

実装

Chapter 6

実験

Chapter 7

関連研究

Chapter 8

結論

Bibliography

- [1] KonohaScript: A Static Scripting Language. <http://konohascript.org>.
- [2] K. Tim and W. Rob. The goal structuring notation-a safety argument notation. In *Proc. of Dependable Systems and Networks 2004 Workshop on Assurance Cases*, 2004.