

スキヤナ・ジェネレータ S G E N の開発

3M-2

山口 裕子 久島 伊知郎 神野 俊昭

(株)日立製作所 システム開発研究所

1. はじめに

コンパイラの自動生成の研究の一環としてスキヤナ・ジェネレータ S G E N (Scanner G E N e r a t o r ) を開発した。コンパイラの中でスキヤナは入力されるソース・プログラム中の全ての文字を読む唯一のフェーズであり、最も時間のかかる場所である。生成するスキヤナの実行性能がよく、入力仕様の書き易さ等の点で使い勝手の良いものを開発することを目的とした。

2. 概要

S G E N は、プログラミング言語の字句の正規表現に基づく形式的仕様を入力とし、決定性有限オートマトンモデルに基づくスキヤナを生成する(図1参照)。スキヤナは字句を認識したら対応する字句番号を返す。字句解析と同時に行う処理は、意味動作として入力仕様中に記述できる。生成されるスキヤナの記述言語にはCとPascalがある。

3. スキヤナの処理

S G E N が生成するスキヤナは文字を読み込み、文字と現在の状態から、オートマトンを表す遷移表上の次の状態をひく遷移表検索手続きを呼び出して次の動作を決めていく。

文字をみるだけで済む場合は、なるべく遷移表検索手続きを呼ばない方法をとる。これにより、全て検索手続きを呼び出す場合に比べ、5~20%高速化している。

4. S G E N への入力仕様

S G E N への入力例の一部を図2に示す。字句の定義部(「#token」以下)には、字句の仕様と対応する字句番号を記述する(何も書いてないものには直前のものに1足した値を自動設定する)。これにより、スキヤナの動作として字句番号を返すという意味動作を書かずに済み、記述が簡潔になる。また、字句の出現回数を数えるといった通常のスキヤナ以外の用途にも適用できるように、字句番号を返す動作をしないようにも指定できる。

キーワードは字句定義部でも定義できるが、遷移表を小さくするために専用のキーワード定義部(「#keyword」以下)を設けた。

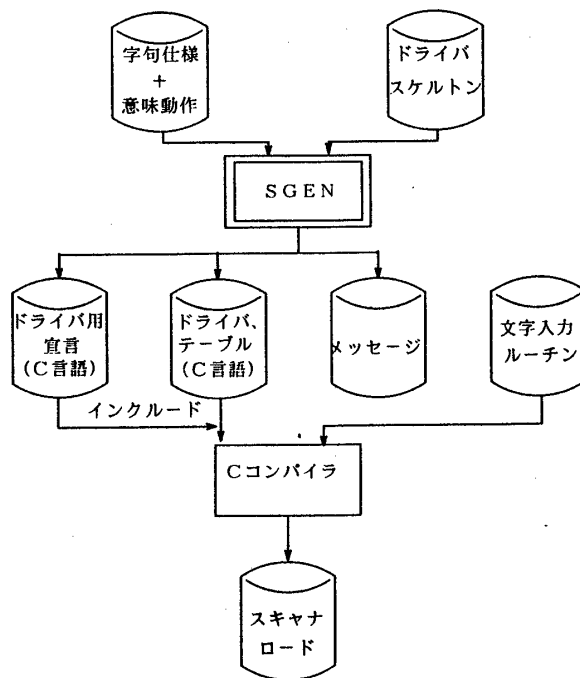


図1 SGEN C版によるスキヤナの生成

```

#keyword
  'sizeof' 17 #permit 1;  --①
  :
  'while'   #recover 'whle';
#lexicon      --②
  ID = letter(letter|digit)*;
#token
  ';'      1 #refact action;
  ','      #refact action;
  :
  ID       #refact IsKeyword
          #limit 14;  --③

```

図2 SGENへの入力例の一部

#### 4. 1. エラー訂正のための仕様

字句解析時に字句エラーを訂正することは一般に困難であり、誤った訂正をする恐れもある。安全に字句エラーを訂正するために、入力仕様で字句毎にどの程度のつづり誤りを訂正するかを指定し処理する方式を考案した。つづり誤りを訂正するか、あるいはエラーメッセージ出力に留めるかはスキヤナ実行時にユーザが指定する。

通常、スキヤナは文字列を切り出したら、定義されたキーワードを登録したキーワードテーブルを検索して該当するものがあればこの字句番号を返し、なければ識別子とする。このためキーワードのつづり誤りは識別子とみなされることが多い。

SGENでは、実際の字句のつづり誤りは1文字違いという場合が多いことに着目して、キーワードテーブルを検索する際に1文字違いのものがあつたら、つづり誤りとみてこの字句番号を返すことにした。このとき問題になるのは、キーワードと1文字違う字句がつづりを誤ったものではなく、識別子だった場合である。このような不正な誤り訂正を避けるため、1文字違いをつづり誤りとみなし

てよいキーワードを入力仕様中で指定させることにした。この指定をしたときユーザは自分で指定したキーワードと紛らわしい識別子を、字句解析するソース中で用いないようにしなければならない。

1文字違いのつづり誤りを許すキーワードに対しては図2の①のように「#permit 1」と書く。字数の多いキーワードは2文字違いまで許すとすれば、「#permit 2」と指定する(3文字以上違っている場合は、つづり誤りというより識別子とみるべきなので、「#permit」の後には2までしか書けない)。

また、図2の②のようにユーザがよく間違えるパターンを指定する仕方も用意した。「#recover 'whle」と指定してあるとスキヤナは「whle」を「while」の誤りとみなす(キーワードテーブルには「whle」を登録してあるので、テーブル検索は「#permit 1」と指定した場合より、高速になる)。これは指定した文字列「whle」を本来の字句「while」に置き換える機能という見方もできる。ユーザは自分用のスキヤナを生成して自分好みの字句を用いたプログラムが処理できるようになる。

#### 4. 2. その他

一般に識別子や文字列リテラルは字数の上限値が定められているが、正規表現ではこれを表すことができない。字句の字数の上限値をチェックするための指定方法を用意した。図2の③の「#limit 14」は識別子(ID)の上限値を14文字としてチェックすることを指定している。

#### 5. おわりに

エラー訂正機能とスキヤナの高性能性を特徴とするスキヤナ・ジェネレータSGENを開発し適用に供している。

#### 参考文献

- H.Mossenbock, Alex--A Simple and Efficient Scanner Generator, SIGPLAN Notices 21(12),139-148(1986)