

# コンパイラ

湯浅太一

# 講義資料のダウンロード

URL:

[www.yuasa.kuis.kyoto-u.ac.jp/~yuasa/index\\_J.html](http://www.yuasa.kuis.kyoto-u.ac.jp/~yuasa/index_J.html)

「コンパイラ」講義資料（パスワードが必要）

ユーザー名：compiler

パスワード：konpaira

「コンパイラ」講義資料

PostScript 版

1up 2up 4up

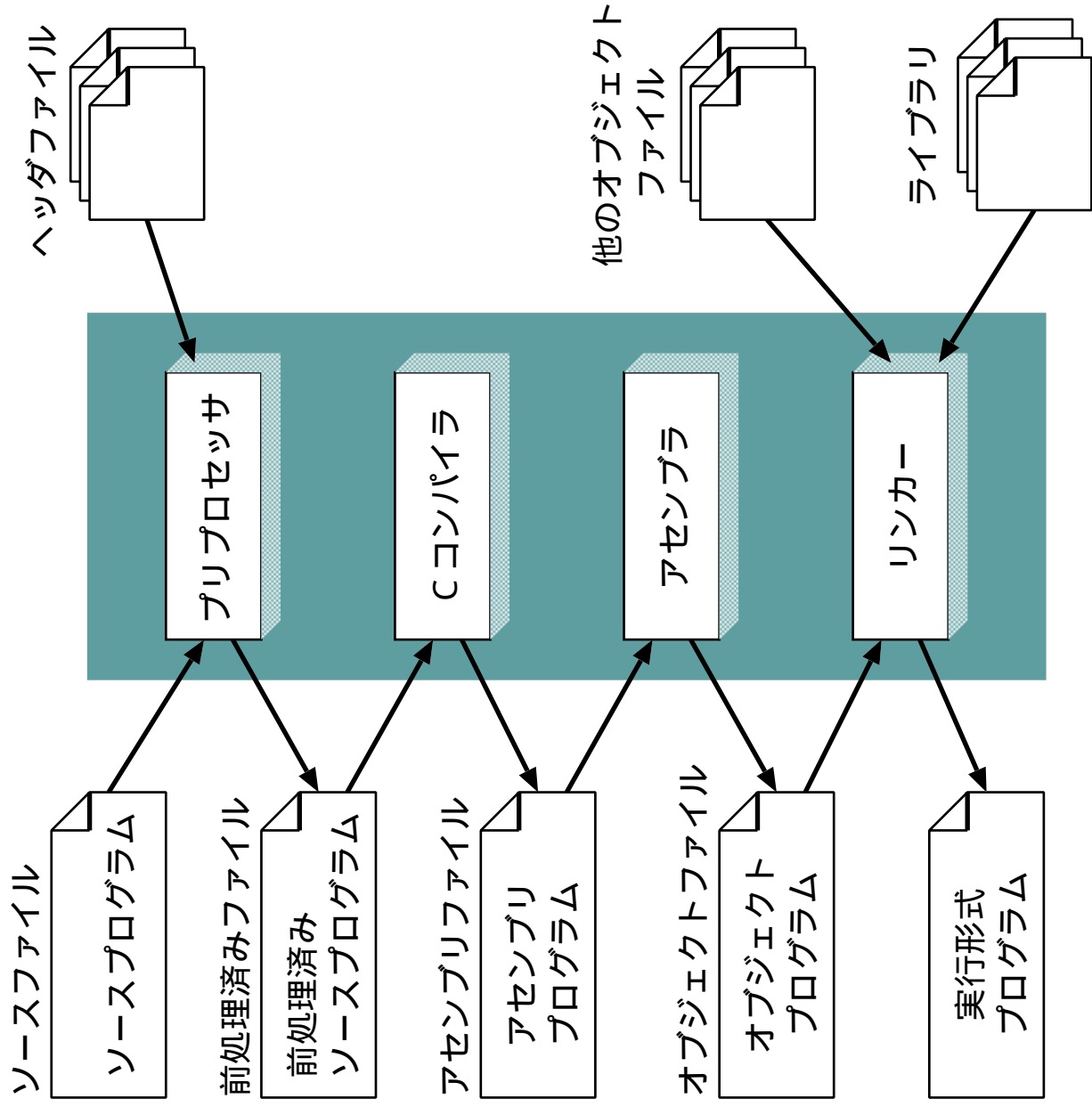
PDF 版

1up 2up 4up

## 第1章

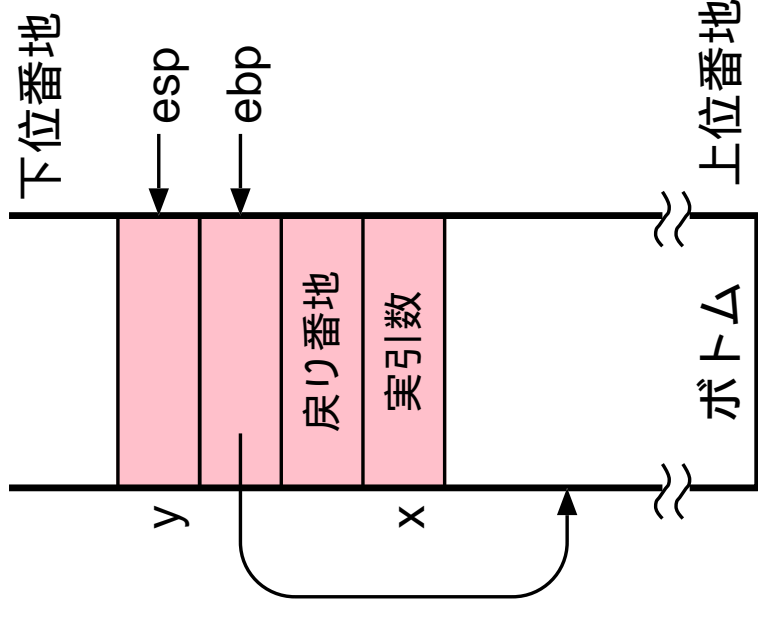
# コンパイラの概要

# Cコンパイラ



# コード生成例

```
int foo(int x) {  
    int y = x*x;  
    return y+2;  
}  
  
1  _foo: push  ebp  
2      mov  ebp, esp  
3      sub  esp, 4  
4      mov  eax, 8[ebp]  
5      imul eax, 8[ebp]  
6      mov  -4[ebp], eax  
7      mov  eax, -4[ebp]  
8      add  eax, 2  
9      mov  esp, ebp  
10     pop  ebp  
11     ret
```



# 機械語コード

0101 0101 1000 1001 1110 0101 1000 0011  
1110 1100 0000 0100 1000 1011 0100 1101  
0000 1000 0000 1111 1010 1111 0100 1101  
0000 1000 1000 1001 0100 1101 1111 1100  
1000 1011 0101 0101 1111 1100 1000 0011  
1100 0010 0000 0010 1000 1001 1101 0000  
1100 1001 1100 0011

1 0101 0101  
2 1000 1001 1110 0101  
3 1000 0011 1110 1100 0000 0100  
4 1000 1011 0100 1101 0000 1000  
5 0000 1111 1010 1111 0100 1101 0000 1000  
6 1000 1001 0100 1101 1111 1100  
7 1000 1011 0101 0101 1111 1100  
8 1000 0011 1100 0010 0000 0010  
9 1000 1001 1101 0000  
10 1100 1001  
11 1100 0011

## mov 命令

1000 1001 11 $x_1x_2$   $x_3y_1y_2y_3$

レジスタ  $x_1x_2x_3$  からレジスタ  $y_1y_2y_3$  へ移動する .

1000 1011 01 $x_1x_2$   $x_3y_1y_2y_3$   $i_1i_2i_3i_4$   $i_5i_6i_7i_8$

レジスタ  $y_1y_2y_3$  の指す位置から相対的に  $i_1i_2i_3i_4i_5i_6i_7i_8$  番地の位置にあるデータを , レジスタ  $x_1x_2x_3$  へ移動する .

1000 1001 01 $x_1x_2$   $x_3y_1y_2y_3$   $i_1i_2i_3i_4$   $i_5i_6i_7i_8$

レジスタ  $y_1y_2y_3$  の指す位置から相対的に  $i_1i_2i_3i_4i_5i_6i_7i_8$  番地の位置へ , レジスタ  $x_1x_2x_3$  のデータを移動する .

## AT&T形式

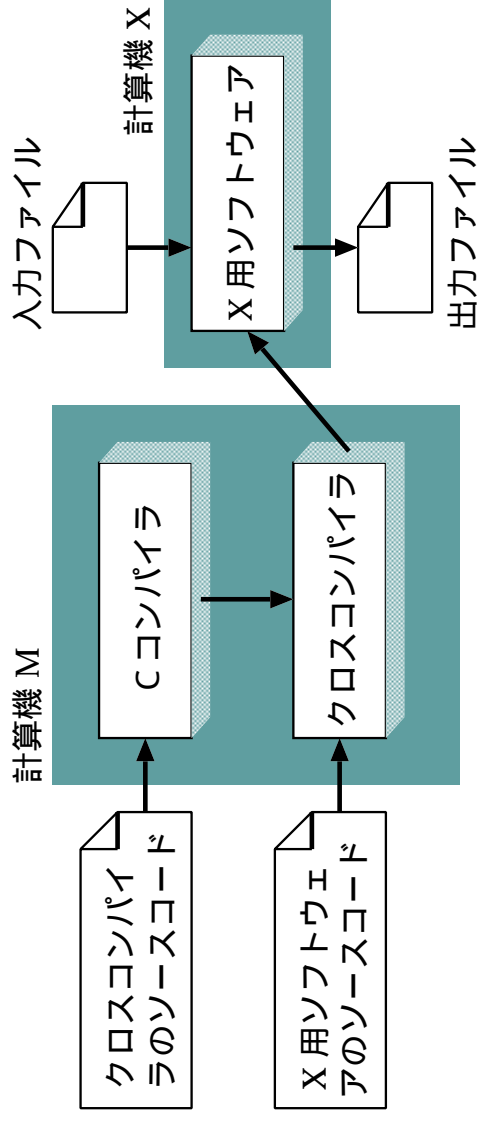
```
1  _foo:  pushl  %ebp
2         movl  %esp,%ebp
3         subl  $4,%esp
4         movl  8(%ebp),%eax
5         imull 8(%ebp),%eax
6         movl  %eax,-4(%ebp)
7         movl  -4(%ebp),%eax
8         addl  $2,%eax
9         movl  %ebp,%esp
10        popl  %ebp
11        ret
```



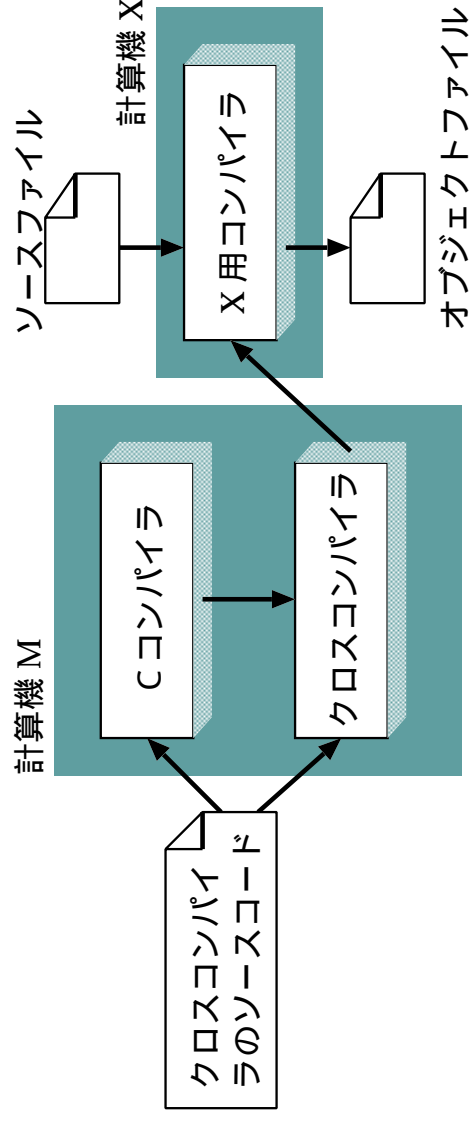
# コンパイラとは

ソース言語から目的言語への変換

## クロスコンパイラ



## ブートストラップ



# コンパイラの構造

1. 字句解析
2. 構文解析
3. 意味解析
4. コード生成

## 字句解析

字句要素を切り出して，トークン列を生成．

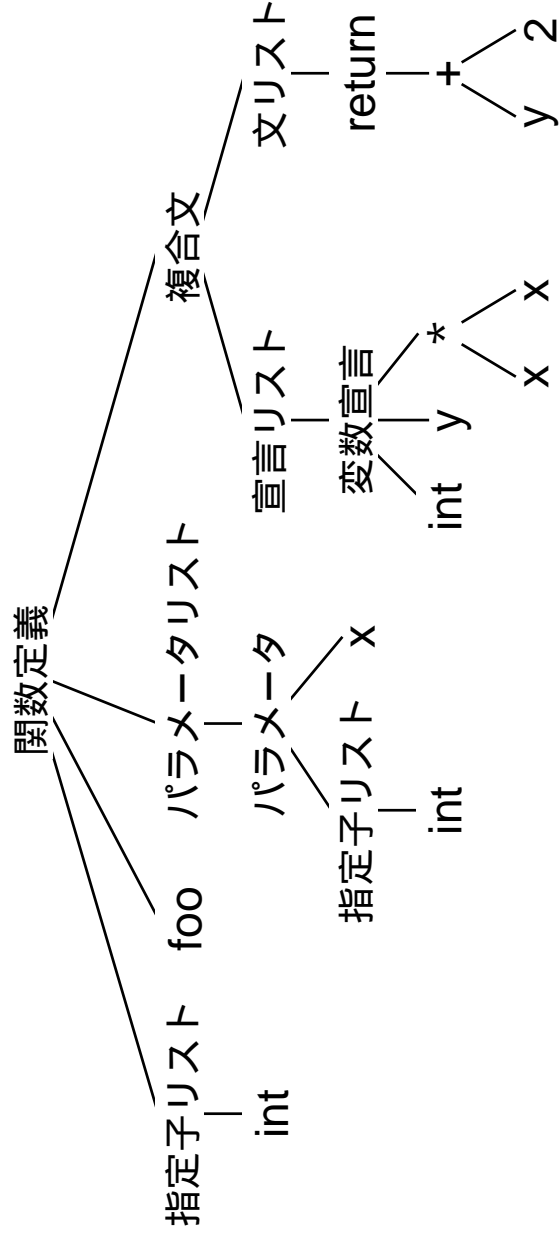
## C言語の字句要素

1. 識別子
2. キーワード
3. 定数（整数定数，浮動小数点定数，文字定数）
4. 文字列リテラル
5. 演算子
6. 区切り記号

```
int foo ( int x ) { int y = x * x ; return y + 2 ; }
```

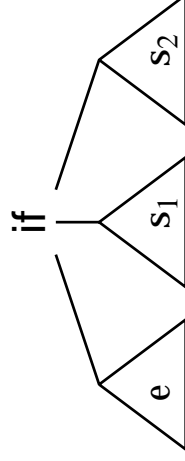
# 構文解析

文法を手エックし，構文木を生成．

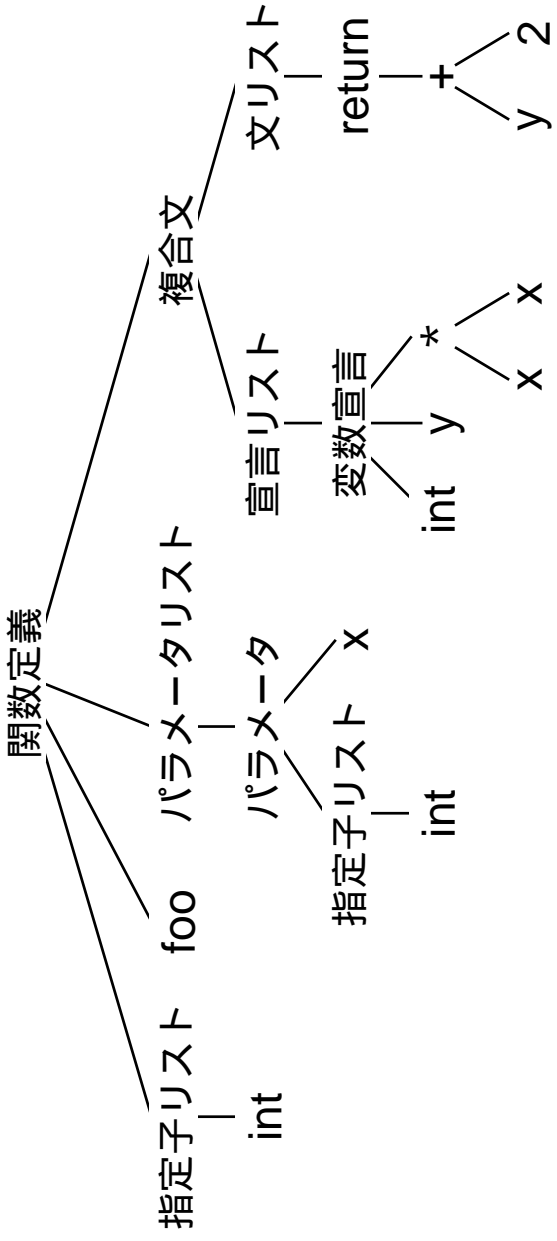


必要な情報のみ保持

if ( 条件 ) 文1 else 文2



# 意味解析



## コード生成

### 局所変数等の割り付け

### 最適化

```
1 _foo: mov  eax, 4[esp]
2      imul eax, eax
3      add  eax, 2
4      ret
```

## 第2章

### 字 句 解 析

## 文字列の計算

アルファベット：言語を構成する文字の（有限）集合

文字列：文字の並び

空の文字列：文字を一つも含まない文字列  $\varepsilon$

文字列の連結  $s \cdot t$

例：aa · bc = aabc

結合則  $s \cdot (t \cdot u) = (s \cdot t) \cdot u$  が成り立つ

→  $s \cdot t \cdot u$  と書いてよい

→  $stu$  と略記することもある

文字列のべき乗

$$s^0 = \varepsilon$$

$$s^n = s \cdot s^{n-1} \quad (n = 1, 2, 3, \dots)$$

公式

$$s^1 = s \cdot s^0 = s \cdot \varepsilon = s$$

$$s^{i+j} = s^i \cdot s^j$$

## 文字列集合の連結

$$A \cdot B = \{s \cdot t \mid s \in A, t \in B\}$$

$$\text{例: } \{a, bc\} \cdot \{a, c\} = \{aa, ac, bca, bcc\}$$

結合則  $A \cdot (B \cdot C) = (A \cdot B) \cdot C$  が成り立つ

→  $A \cdot B \cdot C$  と書いてよい

→  $ABC$  と略記することができる

## 文字列集合のべき乗

$$A^0 = \{\varepsilon\}$$

$$A^n = A \cdot \underbrace{A \cdot \dots \cdot A}_n \quad (n = 1, 2, 3, \dots)$$

例:  $A = \{a, bc\}$  のとき

$$A^3 = A \cdot A \cdot A = \{aaa, aabc, abca, abcabc, bcaa, bcabc, bcbca, bcbcbc\}$$

## 公式

$$A^1 = A \cdot A^0 = A \cdot \{\varepsilon\} = A$$

$$A^{i+j} = A^i \cdot A^j = \underbrace{A \cdot A \cdot \dots \cdot A}_i \cdot \underbrace{A \cdot A \cdot \dots \cdot A}_j$$

## $\{\varepsilon\}$ と空集合

$$\{\varepsilon\} \cdot A = A \cdot \{\varepsilon\} = A$$

$$\phi \cdot A = A \cdot \phi = \phi$$

文字列集合の正の閉包 (positive closure)

$$A^+ = A^1 \cup A^2 \cup A^3 \cup \dots = \bigcup_{i=1}^{\infty} A^i$$

例：

$$\{a\}^+ = \{a, aa, aaa, aaaa, \dots\}$$

$$\{a, b\}^+ = \{a, b, aa, ab, ba, bb, aaa, \dots\}$$

文字列集合の閉包 (closure)

$$A^* = \{\varepsilon\} \cup A^+ = A^0 \cup A^1 \cup A^2 \cup A^3 \cup \dots = \bigcup_{i=0}^{\infty} A^i$$

例：

$$\{a\}^* = \{\varepsilon, a, aa, aaa, aaaa, \dots\}$$

公式

$$A^* \cdot A = A \cdot A^* = A^+$$

$$A^* \cdot A^+ = A^+ \cdot A^* = A^+$$

$$A^* \cdot A^* = A^*$$

例：C 言語における 10 進定数

$$\{1, 2, 3, 4, 5, 6, 7, 8, 9\} \cdot \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}^*$$



# 正規表現

アルファベット  $\Sigma = \{a_1, \dots, a_n\}$  上の正規表現

1.  $\varepsilon$
2.  $a_i$  ( $a_i \in \Sigma$ )
3.  $r_1 | r_2$  ( $r_1, r_2$  は正規表現)
4.  $r_1 \cdot r_2$  ( $r_1, r_2$  は正規表現)
5.  $r^*$  ( $r$  は正規表現)

$L(r)$  : 正規表現  $r$  が表す文字列集合

$$L(\varepsilon) = \{\varepsilon\}$$

$$L(a_i) = \{a_i\}$$

$$L(r_1 | r_2) = L(r_1) \cup L(r_2)$$

$$L(r_1 \cdot r_2) = L(r_1) \cdot L(r_2)$$

$$L(r^*) = L(r)^*$$

例 :  $a \cdot ((a|b)^*)$  は ,  $\{a, b\}$  上の正規表現 .  $a(a|b)^*$  と略記

$$L(a(a|b)^*)$$

$$= L(a) \cdot L((a|b)^*) = L(a) \cdot (L(a|b))^* = L(a) \cdot (L(a) \cup L(b))^* = \{a\} \cdot (\{a\} \cup \{b\})^*$$

$$= \{a\} \cdot \{a, b\}^* = \{a\} \cdot \{\varepsilon, a, b, aa, ab, ba, bb, aaa, \dots\}$$

$$= \{a, aa, ab, aaa, aab, aba, abb, aaaa, \dots\}$$

$\varepsilon \mid r$  :  $r$  が省略可能

$a_i$  :  $a_i$  自体

$r_1 \mid r_2$  :  $r_1$  または  $r_2$

$r_1 \cdot r_2$  :  $r_1$  のあとに  $r_2$

$r^*$  :  $r$  が 0 個以上

## C 言語の字句要素

識別子:  $\alpha(\alpha \mid \delta)^*$

キーワード: auto | break | case | char | ...

演算子: + | & | = | += | ...

区切り記号: , | ; | ( | ) | ...

10 進定数: (1 | ... | 9)  $\delta^*$   $\langle$  オプション  $\rangle$

浮動小数点定数: (( $\delta^+ \cdot \delta^*$  |  $\delta^+$ )  $\langle$  指数部  $\rangle_{opt}$  |  $\delta^+$   $\langle$  指数部  $\rangle$ ) (F | f | L | l)  $_{opt}$

文字定数: L  $_{opt}$  ' ( \  $\langle$  文字  $\rangle$  | \  $\langle$  と ' と改行以外  $\rangle$  )  $^+$  ,

ここで ,

$\alpha = \_ \mid a \mid b \mid \dots \mid z \mid A \mid B \mid \dots \mid Z$

$\delta = 0 \mid 1 \mid 2 \mid \dots \mid 9$

$\langle$  オプション  $\rangle = ((U \mid u) (L \mid l)  $_{opt}$  | (L | l) (U | u)  $_{opt}$ )  $_{opt}$$

$\langle$  指数部  $\rangle = (E \mid e) -  $_{opt}$   $\delta^+$$

# 有限オートマトン (finite automaton)

有限個の状態を持ち，文字を1文字ずつ読み込むことによって，自分自身の状態を変化させる，一種の仮想計算機

FAは五つ組  $\langle \Sigma, Q, \Delta, s, F \rangle$  で定義

$\Sigma$ : アルファベット

$Q$ : 状態全体の集合

$\Delta$ : 状態遷移  $\langle q_1, a, q_2 \rangle$  の集合

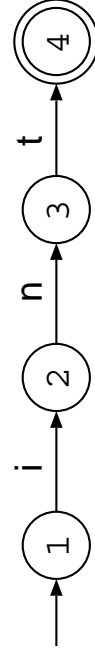
状態  $q_1$  で  $a$  を読んだら状態  $q_2$  に遷移

$s$ : 初期状態 ( $s \in Q$ )

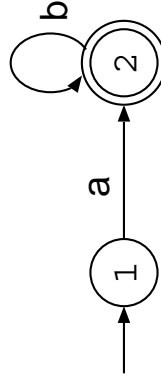
$F$ : 受理状態の集合 ( $F \subset Q$ )

## 状態遷移図

$M1 = \langle \Sigma, \{1, 2, 3, 4\}, \{ \langle 1, i, 2 \rangle, \langle 2, n, 3 \rangle, \langle 3, t, 4 \rangle \}, 1, \{4\} \rangle$



$M2 = \langle \Sigma, \{1, 2\}, \{ \langle 1, a, 2 \rangle, \langle 2, b, 2 \rangle \}, 1, \{2\} \rangle$



$M = \langle \Sigma, Q, \Delta, s, F \rangle$ : FA

$\alpha = a_1 a_2 \cdots a_n$ : 入力文字列

次の条件を満たす状態列  $q_0, q_1, \dots, q_n$  が存在するとき,  $M$  は  $\alpha$  を受理 (accept) する .

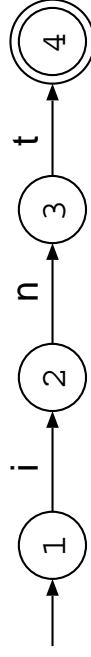
$$s = q_0 \xrightarrow{a_1} q_1 \xrightarrow{a_2} \cdots \xrightarrow{a_n} q_n \in F$$

$L(M)$ :  $M$  が受理する文字列全体

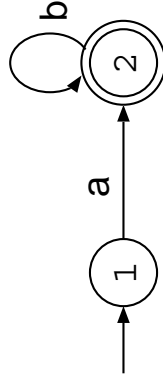
$\rightarrow M$  が受理する言語 (language)

例 :

$L(M1) = \{\text{int}\}$



$L(M2) = \{a, ab, abb, abbb, \dots\} = L(ab^*)$



$L(M) = L(M')$  のとき,  $M$  と  $M'$  は等価 (equivalent)

## 決定性 (deterministic) 有限オートマトン

任意の  $q \in Q$  と  $a \in \Sigma$  に対して,  $\langle q, a, q' \rangle$  の形の状態遷移が高々一つ

DFA  $\langle \Sigma, Q, \Delta, s, F \rangle$  の状態遷移関数

$$\delta(q, a) = \begin{cases} q' & (\langle q, a, q' \rangle \in \Delta \text{ のとき}) \\ \perp & (\text{その他の場合}) \end{cases}$$

$\delta$  が与えられれば,

$$\Delta = \{ \langle q, a, \delta(q, a) \rangle \mid q \in Q, a \in \Sigma, \delta(q, a) \neq \perp \}$$

DFA を  $\langle \Sigma, Q, \delta, s, F \rangle$  と表記することができる.

適用例 (1): 受理するかどうかの判定

入力文字列:  $a_1 a_2 \cdots a_n$

1.  $i \leftarrow 1, q \leftarrow s$
2.  $i \leq n$  かつ  $\delta(q, a_i) \neq \perp$  ならば,  
 $q \leftarrow \delta(q, a_i), i \leftarrow i + 1$

とし, このステップ 2 を繰り返す.

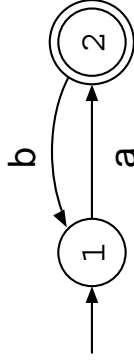
3.  $i > n$  かつ  $q \in F$  であれば  $\text{yes}$  と表示し, そうでなければ  $\text{no}$  と表示する.

## 適用例 (2) : 文字列の切り出し

入力文字列:  $a_1 a_2 \cdots a_n$

1.  $k \leftarrow 0, i \leftarrow 0, q \leftarrow s$  (初期状態)
2.  $q \in F$  であれば,  $k \leftarrow i$
3.  $i \leftarrow i + 1$
4.  $\delta(q, a_i) \neq \perp$  なら,  $q \leftarrow \delta(q, a_i)$  とし, ステップ 2 へ戻る. そうでなければ,  $k$  の値を返す.

例:



入力文字列が ababb のとき

$$1 \xrightarrow{a} 2_{(k=1)} \xrightarrow{b} 1 \xrightarrow{a} 2_{(k=3)} \xrightarrow{b} 1 \xrightarrow{b} \perp$$

# 非決定性有限オートマトン



状態  $q$  に達したときに，入力文字列と無関係に， $q'$  へ遷移してよい． $q$  に留まってもよい．

非決定性 (nondeterministic) 有限オートマトン (NFA)

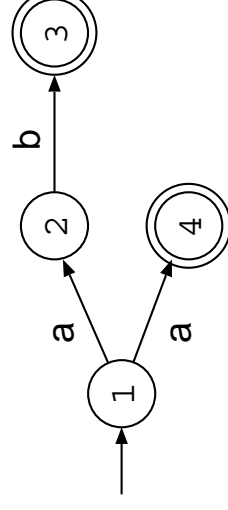
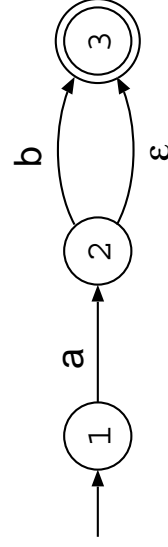
次のいずれかを満たす FA  $\langle \Sigma, Q, \Delta, s, F \rangle$

1.  $\Delta$  が  $\varepsilon$  遷移を含む．
2. ある  $q \in Q$  と  $a \in \Sigma$  に対して， $\langle q, a, q' \rangle$  の形の状態遷移が， $\Delta$  の中に複数存在する．

例

M4 =  $\langle \Sigma, \{1, 2, 3\}, \{ \langle 1, a, 2 \rangle, \langle 2, b, 3 \rangle, \langle 2, \varepsilon, 3 \rangle \}, 1, \{3\} \rangle$

M5 =  $\langle \Sigma, \{1, 2, 3, 4\}, \{ \langle 1, a, 2 \rangle, \langle 2, b, 3 \rangle, \langle 1, a, 4 \rangle \}, 1, \{3, 4\} \rangle$



次の条件をすべて満たす経路

$$q_0 \xrightarrow{x_1} q_1 \xrightarrow{x_2} \dots \xrightarrow{x_m} q_m$$

が存在するとき，NFA  $\langle \Sigma, Q, \Delta, s, F \rangle$  は，文字列  $\alpha = a_1 a_2 \dots a_n$  を受理する．

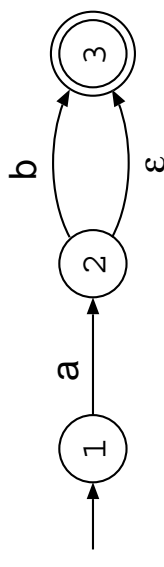
1.  $q_0 = s$
2.  $\langle q_{i-1}, x_i, q_i \rangle \in \Delta$  ( $i = 1, 2, \dots, m$ )
3.  $x_1, x_2, \dots, x_m$  から  $\varepsilon$  を除去すると  $\alpha$  と一致する．
4.  $q_m \in F$

例:

$$L(M4) = \{a, ab\}$$

$$a \text{ に対して: } 1 \xrightarrow{a} 2 \xrightarrow{\varepsilon} 3$$

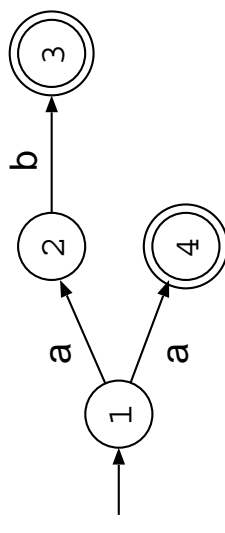
$$ab \text{ に対して: } 1 \xrightarrow{a} 2 \xrightarrow{b} 3$$



$$L(M5) = \{a, ab\}$$

$$a \text{ に対して: } 1 \xrightarrow{a} 4$$

$$ab \text{ に対して: } 1 \xrightarrow{a} 2 \xrightarrow{b} 3$$





# 正規表現から DFA への変換

$r$  : 正規表現

$M$  : FA

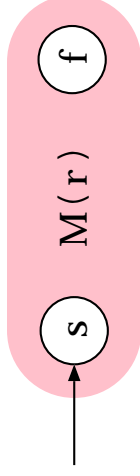
$L(r) = L(M)$  のとき,  $r$  と  $M$  は等価

正規表現  $r$  を, 等価な DFA に変換

1.  $r$  を, 等価な NFA に変換
2. NFA を, 等価な DFA に変換
3. DFA を, 状態数最小の等価な DFA に変換

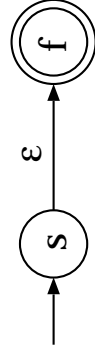
# 正規表現からNFAへの変換

$M(r)$ : 正規表現  $r$  に対する標準形NFA

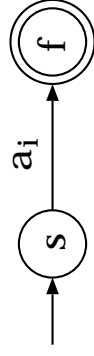


1. 初期状態への遷移は存在しない
2. 受理状態は一つだけ
3. 受理状態からの遷移は存在しない

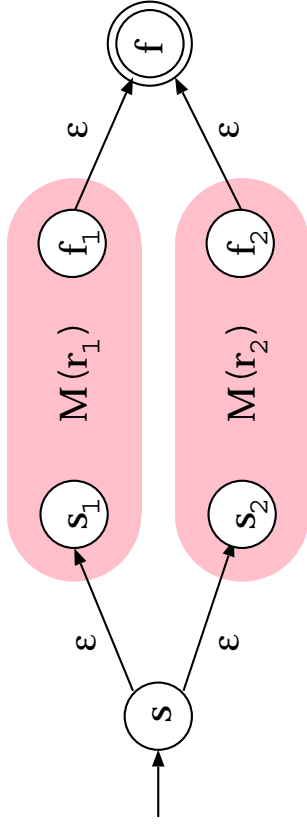
(1)  $M(\varepsilon)$



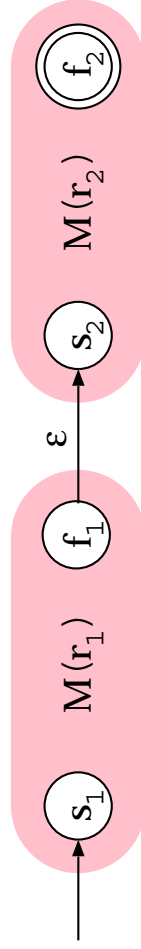
(2)  $M(a_i)$



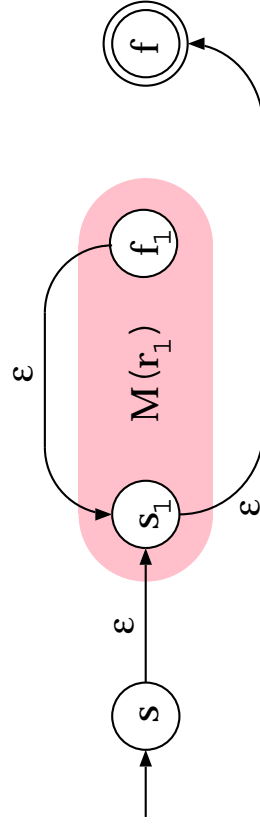
(3)  $M(r_1 \mid r_2)$



(4)  $M(r_1 \cdot r_2)$



(5)  $M(r_1^*)$

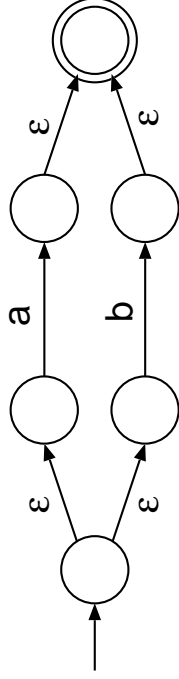


例:  $M((a|b)^*ab)$

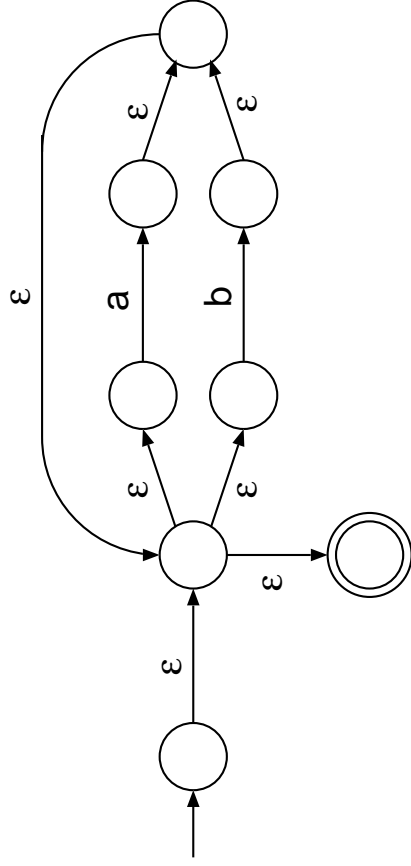
(1)  $M(a) \subsetneq M(b)$



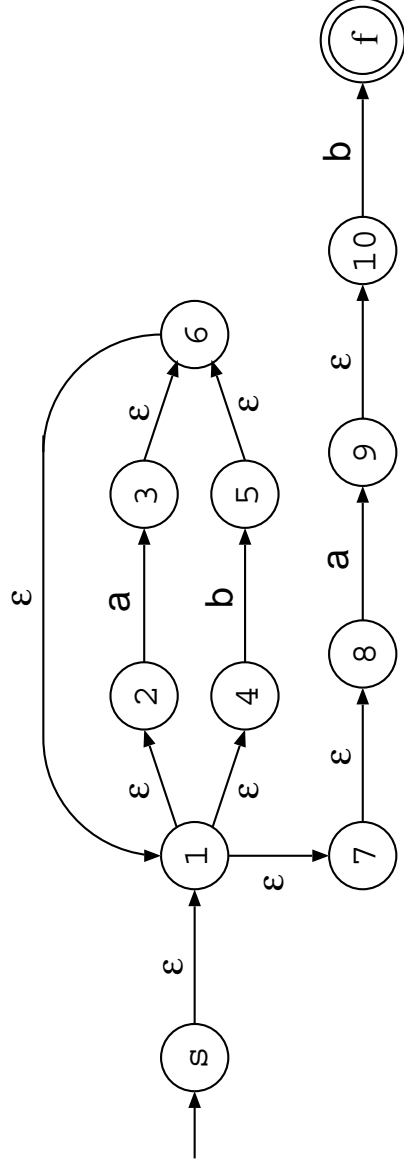
(2)  $M(a|b)$



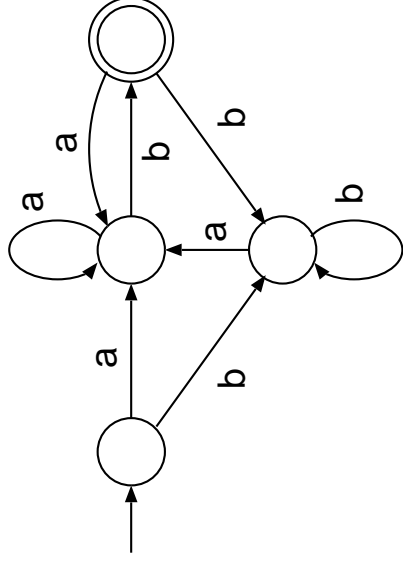
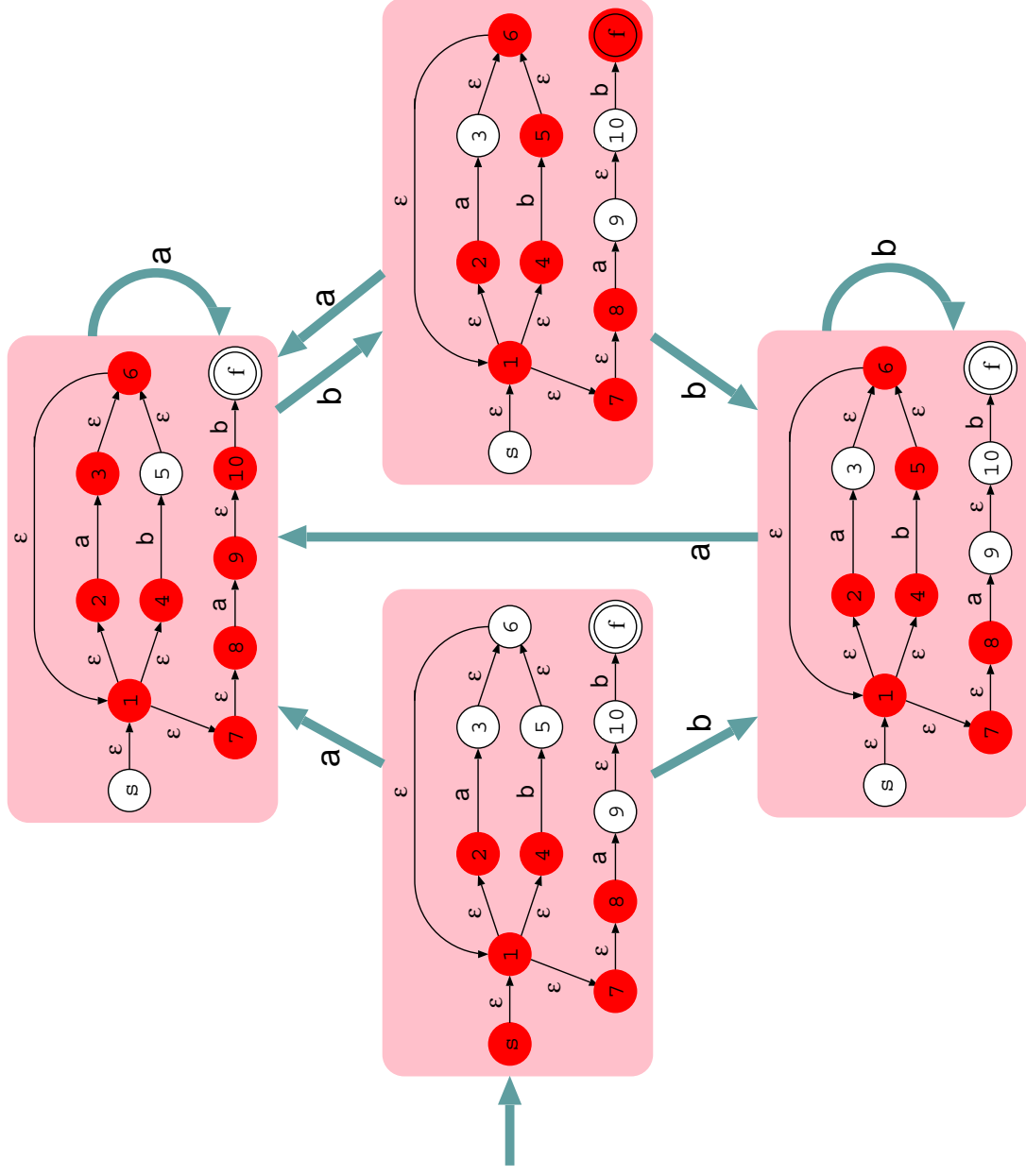
(3)  $M((a|b)^*)$



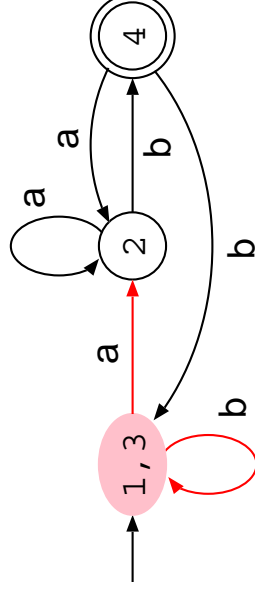
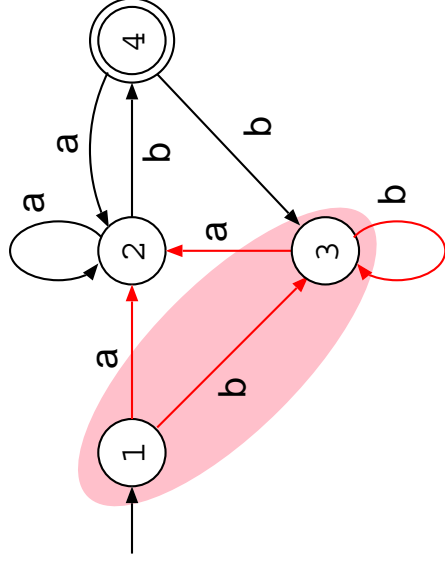
(4)  $M((a|b)^*ab)$



# NFAからDFAへの変換

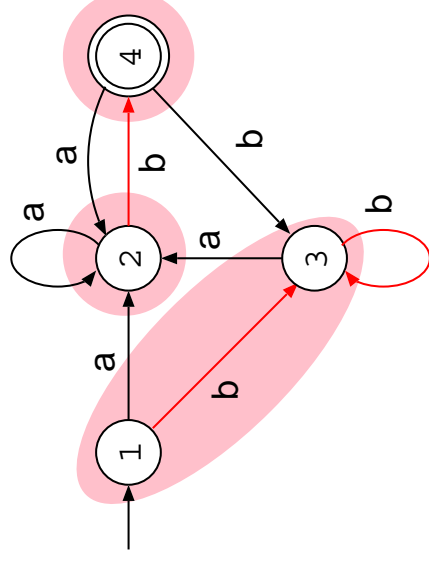
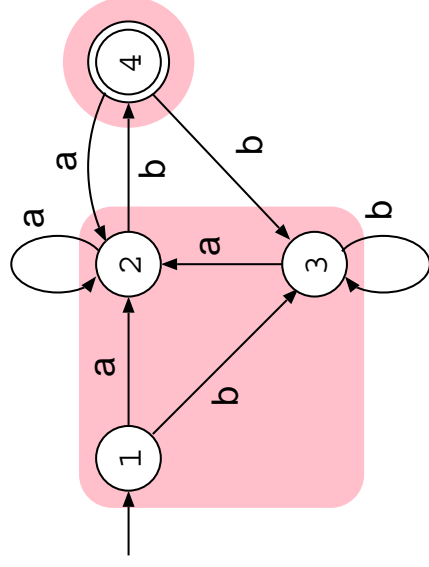


# 状態数最小のDFA



## アルゴリズム

1. 受理状態とその他の2グループに分ける .
2. 同一文字に対する遷移先が異なる状態を , 別グループに分ける .



# 字句解析プログラム

## 正規表現以外の規則

1. 最も長い字句要素を切り出す .

例:  $abc = a + bc = ab + c$

2. 字句要素の種類ごとに優先度を設ける .

例: `if` は識別子ではなく, キーワード

# 字句構造全体の DFA

$r_1, r_2, \dots, r_n$ : 字句要素の種類ごとの正規表現

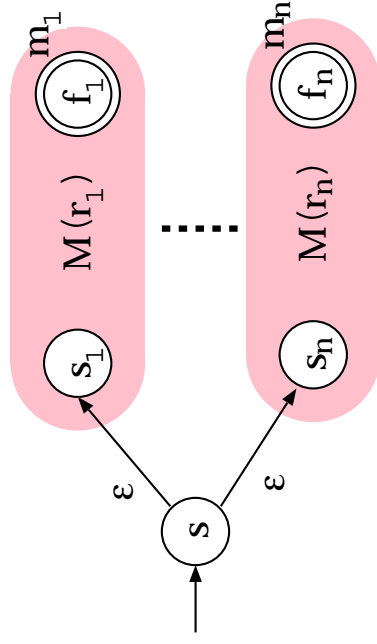
$r_1 | r_2 | \dots | r_n$  の DFA では種類がわからない

$m_1, m_2, \dots, m_n$ : 種類ごとのマーカー → 各受理状態に適切なマーカーをつける

1. 種類ごとに NFA  $M(r_i)$  を求め, 受理状態にマーカーをつける



2. 新しい初期状態  $s$  を用意し, 全体の NFA  $M_{all}$  を作成.

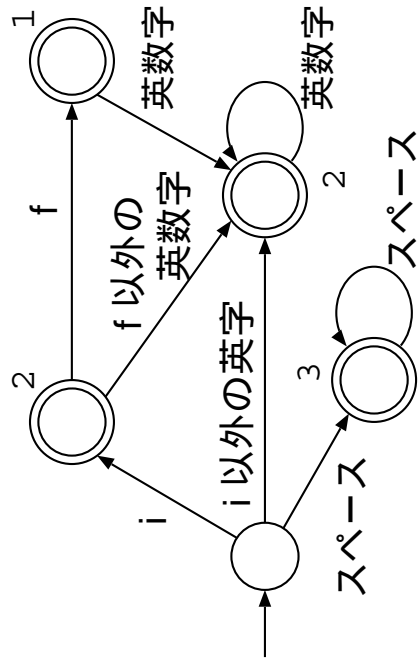
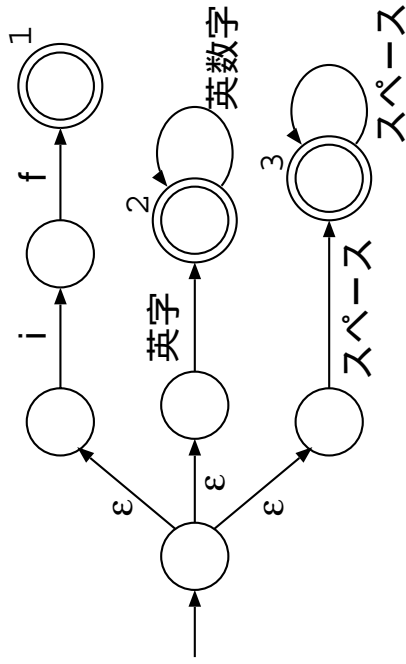


3.  $M_{all}$  を DFA  $\widehat{M}_{all}$  に変換.  
個々の受理状態には, 優先度最大のマーカーを残す.
4. 状態数最小の DFA を作成.  
ただし, マーカーの異なる受理状態は別グループ.



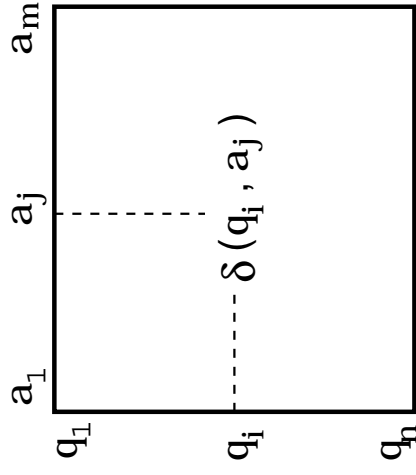
# 例

1. キーワード:  $if$
2. 識別子:  $\langle \text{英字} \rangle \langle \text{英数字} \rangle^*$
3. 空白:  $\langle \text{スペース} \rangle^+$



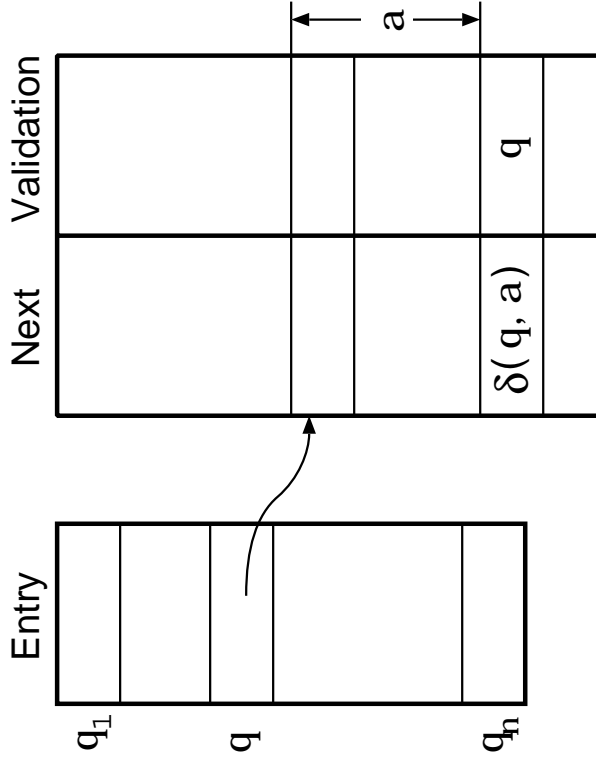
# 状態遷移表

状態遷移関数  $\delta : Q \times \Sigma \rightarrow Q \cup \{\perp\}$  を表にしたもの



ほとんどの要素が  $\perp$  (未定義)

# 状態遷移表の圧縮



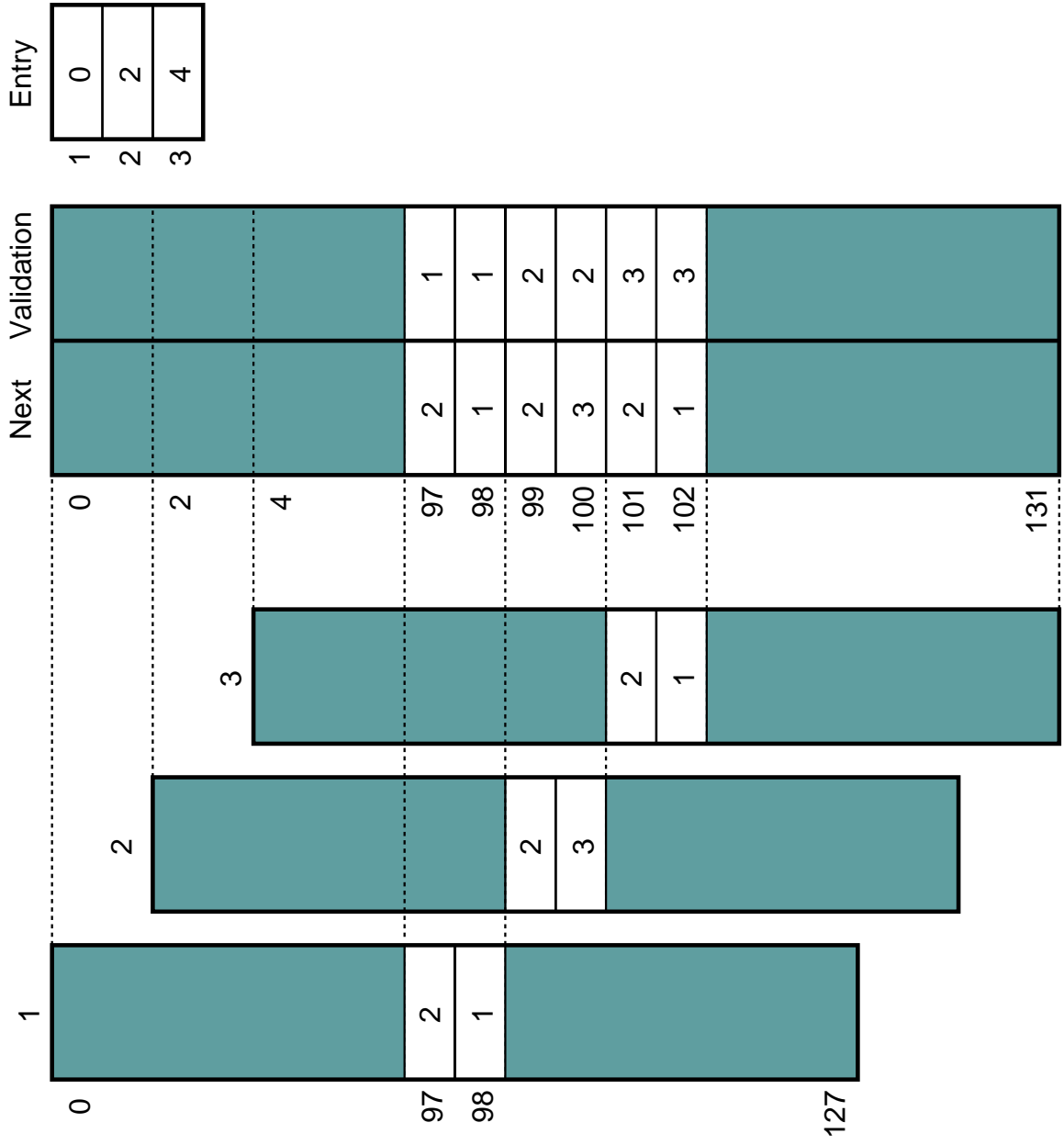
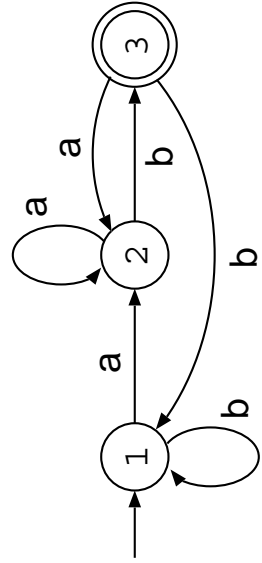
$$\text{Validation}[\text{Entry}[q]+a] = \begin{cases} q & (\delta(q, a) \neq \perp) \\ q \text{ 以外} & (\delta(q, a) = \perp) \end{cases}$$

```

next_state(q, a) {
    int i = Entry[q] + a;
    if (Validation[i] == q)
        return Next[i];
    else
        return UNDEFINED;
}

```

# 各ベクタの設定方法

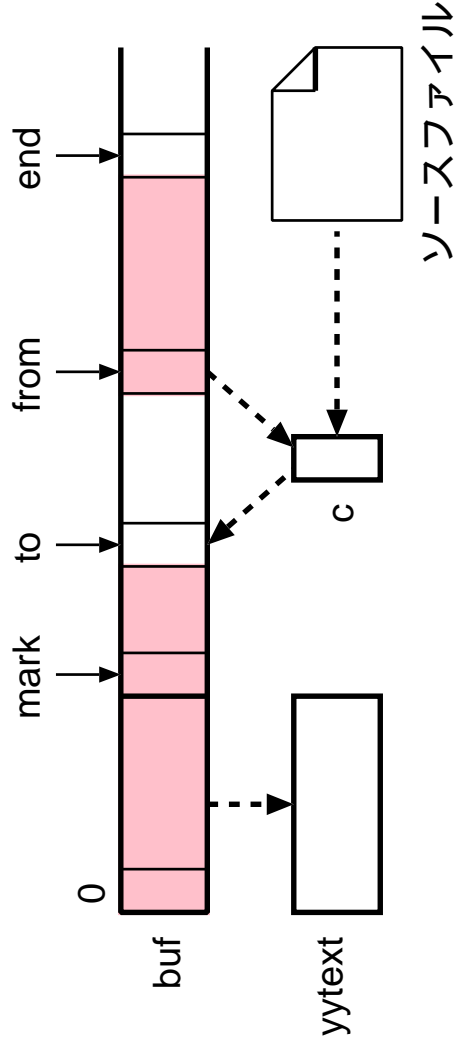


ASCIIコードでは, 'a' = 97 ; 'b' = 98

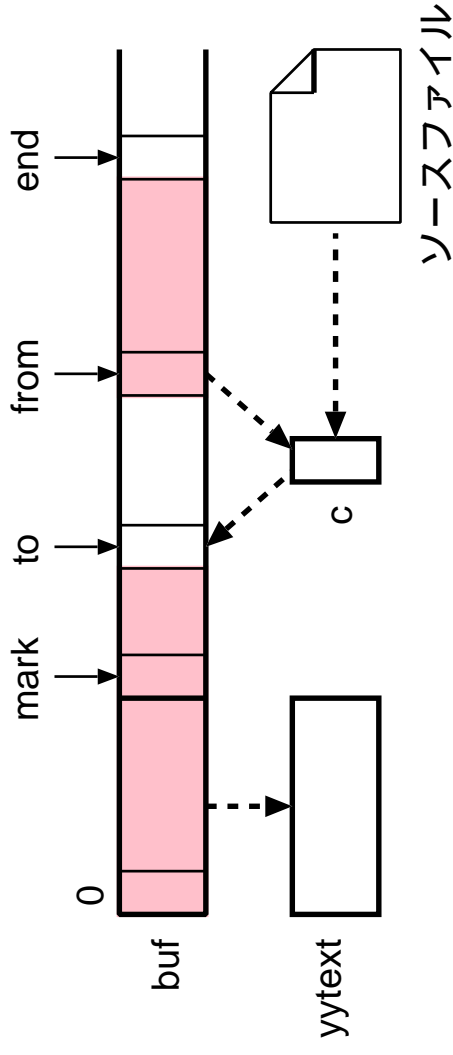
# 字句要素の切り出し

関数 scan:

最も長い字句要素を切り出して yytext に格納し，マーカーの値を返す．入力の終りに達していれば，ENDOFFILE を返す．



- `buf [0] ~ buf [to-1]`:  
これまでに読み込んだ文字列
- `buf [0] ~ buf [mark-1]`:  
最も長い字句要素
- `buf [from] ~ buf [end-1]`:  
前回読み込みすぎた文字列



## cへの文字設定

```

advance_c() {
    buf[to++] = c;
    if (from < end)
        c = buf[from++];
    else
        c = getc(in);
}

```

## バッファの初期化

```

init_buf() {
    mark = to = from = end = 0;
    c = getc(in);
}

```

## 字句要素の切り出し

```
cutout_lexeme() {
    for (i = 0; i < mark; i++)
        yytext[i] = buf[i];
    yytext[i] = '\0';
    if (mark < to) {
        buf[to++] = c;
        while (from < end)
            buf[to++] = buf[from++];
        c = buf[mark++];
        from = mark;
        end = to;
    }
    mark = to = 0;
}
```

## その他のサブ関数

`next_state( $q, a$ )`: 状態遷移関数  $\delta(q, a)$  の値

`final( $q$ )`: 状態  $q$  が受理状態かどうか

`marker( $q$ )`: 受理状態  $q$  のマーカー

`mark_buf()`: `to` の値を `mark` に保存

```
scan() {
    if (c == EOF) return ENDOFFILE;
    q = 1; qf = UNDEFINED;
    for (;;) {
        if (final(q)) {
            mark_buf(); qf = q;
        }
        if ((x = next_state(q,c)) != UNDEFINED) {
            q = x; advance_c();
        } else if (qf != UNDEFINED) {
            cutout_lexeme(); return marker(qf);
        } else {
            lexical_error();
            if (c == EOF) return ENDOFFILE;
            q = 1;
        }
    }
}
```



エラーリカバリ（次の空白まで読み飛ばす）

```
lexical_error() {  
    while (c != ' ' && c != '\n' && c != EOF)  
        advance_c();  
    mark_buf();  
    cutout_lexeme();  
    printf("lexical error: '%s'\n", yytext);  
}
```

# 字句解析の自動化

lexとflex

```
% cat scan.l
%%
[\\t ]+ {}
(a|b)*ab { printf("%s: OK.\\n", yytext);}
. { printf("%s: wrong.\\n", yytext);}
%%
% lex scan.l
% cc lex.yy.c -ll
% a.out
aabbab
aabbab: OK.

aba
ab: OK.
a: wrong.

%
```

```

% cat scan1.l
%%
[\\t ]+ {}
"+"|"-"|"*"|" "/"|"="
    printf("%s (operator) ", yytext);
";" printf("%s (separator) ", yytext);
[a-z][a-z0-9]*
    printf("%s (identifier) ", yytext);
[0-9\\.]+ {printf("%s (number) ", yytext);}
. printf("%s (don't know) ", yytext);
%%
% lex scan1.l
% cc lex.yy.c -ll
% a.out
abc=e*3.14+abc/e
abc (identifier) = (operator) e (identifier)
* (operator) 3.14 (number) + (operator)
abc (identifier) / (operator) e (identifier)

```

# DFAから正規表現への変換

$$\text{DFA } M = \langle \Sigma, \{q_1, q_2, \dots, q_n\}, \delta, q_1, F \rangle$$

$R_{ij}^k$ :  $q_i$  から  $q_j$  へ,  $q_1, q_2, \dots, q_k$  だけを通して遷移させる文字列全体の集合

$$q_i \rightarrow \boxed{q_1 \sim q_k} \rightarrow q_j$$

1.  $k > 0$  のとき

$$R_{ij}^k = R_{ik}^{k-1} (R_{kk}^{k-1})^* R_{kj}^{k-1} \cup R_{ij}^{k-1}$$

$$q_i \rightarrow \boxed{q_1 \sim q_{k-1}} \rightarrow q_k \rightarrow \dots \rightarrow q_k \rightarrow \boxed{q_1 \sim q_{k-1}} \rightarrow q_j$$

$$q_i \rightarrow \boxed{q_1 \sim q_{k-1}} \rightarrow q_j$$

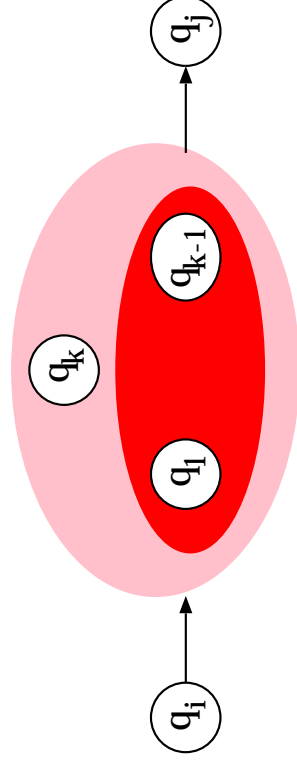
2.  $k = 0$  のとき

$$R_{ij}^0 = \begin{cases} \{a \mid \delta(q_i, a) = q_j\} & (i \neq j \text{ のとき}) \\ \{a \mid \delta(q_i, a) = q_j\} \cup \{\varepsilon\} & (i = j \text{ のとき}) \end{cases}$$

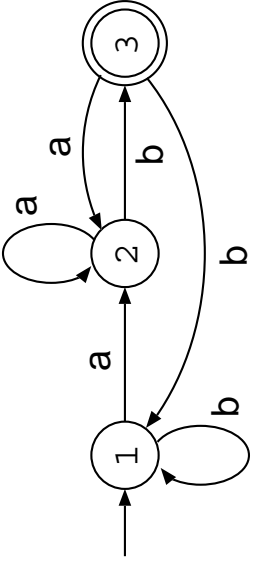
$n$ :  $M$  の状態数

$$L(M) = \bigcup_{q \in F} R_{sq}^n$$

この集合式を正規表現に変換する



例



$$R_{13}^3 = R_{13}^2(R_{33}^2)^*R_{33}^2 \cup R_{13}^2 = R_{13}^2(R_{33}^2)^+ \cup R_{13}^2 = R_{13}^2(R_{33}^2)^*$$

$$R_{13}^2 = R_{12}^1(R_{22}^1)^*R_{23}^1 \cup R_{13}^1$$

$$R_{33}^2 = R_{32}^1(R_{22}^1)^*R_{23}^1 \cup R_{33}^1$$

$$\begin{aligned} R_{12}^1 &= R_{11}^0(R_{11}^0)^*R_{12}^0 \cup R_{12}^0 = (R_{11}^0)^+R_{12}^0 \cup R_{12}^0 \\ &= (R_{11}^0)^*R_{12}^0 = \{b, \varepsilon\}^*\{a\} = \{b\}^*\{a\} \end{aligned}$$

$$R_{22}^1 = R_{21}^0(R_{11}^0)^*R_{12}^0 \cup R_{22}^0 = \phi \cup \{a, \varepsilon\} = \{a, \varepsilon\}$$

$$R_{23}^1 = R_{21}^0(R_{11}^0)^*R_{13}^0 \cup R_{23}^0 = \phi \cup \{b\} = \{b\}$$

$$R_{13}^1 = R_{11}^0(R_{11}^0)^*R_{13}^0 \cup R_{13}^0 = \phi$$

$$R_{32}^1 = R_{31}^0(R_{11}^0)^*R_{12}^0 \cup R_{32}^0 = \{b\}\{b\}^*\{a\} \cup \{a\} = \{b\}^*\{a\}$$

$$R_{33}^1 = R_{31}^0(R_{11}^0)^*R_{13}^0 \cup R_{33}^0 = \{b\}\{b\}^*\phi \cup \{\varepsilon\} = \{\varepsilon\}$$

$$R_{13}^3 = R_{13}^2(R_{33}^2)^*$$

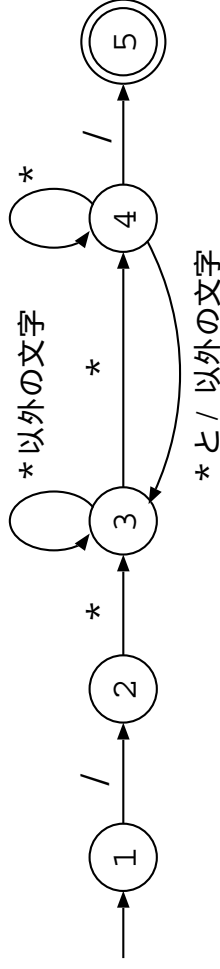
$$= (\{b\}^*\{a\}\{a, \varepsilon\}^*\{b\})(\{b\}^*\{a\}\{a, \varepsilon\}^*\{b\})^*$$

$$= (\{b\}^*\{a\}^+\{b\})(\{b\}^*\{a\}^+\{b\})^*$$

$$= (\{b\}^*\{a\}^+\{b\})^+$$

DFA と等価な正規表現は ,  $(b^*a^+b)^+$

## 例：コメントの正規表現



$$/* \langle * \text{以外の文字} \rangle^* * \\ (( \langle * \text{と} / \text{以外の文字} \rangle^* * | * )^* )^* /$$

状態番号の3と4を入れ替えると

$$/* ( \langle * \text{以外の文字} \rangle | ( * )^+ \langle * \text{と} / \text{以外の文字} \rangle )^* ( * )^* * /$$

## 正規表現の限界

正規表現では表せない文字列集合の例

$$X = \{a^i cb^i \mid i \geq 0\}$$

正規表現  $r$  で表せたと仮定する

$n$ :  $M(r)$  の状態数

$a^n cb^n$  ( $\in X$ ) を受理する  $M(r)$  の経路

$$q_0 \xrightarrow{a} \cdots \xrightarrow{a} q_n \xrightarrow{c} q_{n+1} \xrightarrow{b} \cdots \xrightarrow{b} q_{2n+1}$$

$q_i = q_j$  なる  $i, j$  ( $0 \leq i < j \leq n$ ) が存在する

$$q_0 \xrightarrow{a} \cdots \xrightarrow{a} q_i \xrightarrow{a} \cdots \xrightarrow{a} q_j \xrightarrow{a} \cdots \xrightarrow{a} q_n \xrightarrow{c} q_{n+1} \xrightarrow{b} \cdots \xrightarrow{b} q_{2n+1}$$

$q_i$  から  $q_j$  への経路をカットすると

$$q_0 \xrightarrow{a} \cdots \xrightarrow{a} q_i (= q_j) \xrightarrow{a} \cdots \xrightarrow{a} q_n \xrightarrow{c} q_{n+1} \xrightarrow{b} \cdots \xrightarrow{b} q_{2n+1}$$

$M(r)$  は  $a^{n-(j-i)} cb^n$  ( $\notin X$ ) も受理する  $\rightarrow$  矛盾

文脈自由文法なら

$$A \rightarrow aAb \mid c$$

## 第3章

# 文 法



## 構文，制約，意味

構文 (syntax) : プログラムの構造を定義

制約 (constraint) : 構文以外の規則

意味 (semantics) : 実行を定義

例 : if文

構文 : `if ( expr ) stat1 else stat2`

制約 : *expr* は真偽値を表す型

意味 : *expr* を計算し，その値が真なら *stat*<sub>1</sub> を，  
偽なら *stat*<sub>2</sub> を実行する

文法 : 構文 + 制約

構文解析は「文法解析」ではない。

構文解析の話をするときは，制約には言及しない。

→ 構文 = 文法

# 構文の記法

## ALGOL 60のBNF

$\langle \text{program} \rangle ::=$   
 $\text{program } \langle \text{identifier} \rangle ( \langle \text{identifier-list} \rangle ) ; \langle \text{block} \rangle .$   
 $\quad | \text{program } \langle \text{identifier} \rangle ; \langle \text{block} \rangle .$

$\langle \text{identifier} \rangle$  は非終端記号 (nonterminal symbol) あるいは構文変数 (syntax variable)

$\text{program}$  は、終端記号 (terminal symbol)

$\langle \text{identifier-list} \rangle ::= \langle \text{identifier} \rangle | \langle \text{identifier-list} \rangle , \langle \text{identifier} \rangle$

識別子リストとは:

1. 識別子であるか、
2. 識別子リストと識別子をカンマで区切ったもの

例：“x, y, z”

1. xは識別子なので識別子リスト
2. xが識別子リストなので“x, y”も識別子リスト
3. “x, y”が識別子リストなので，“x, y, z”は識別子リスト

生成規則 (production)

$::=$  と縦棒は、メタ記号 (meta-symbol)

拡張BNFの例：

$\langle \text{identifier-list} \rangle ::= \{ \langle \text{identifier} \rangle , \} \langle \text{identifier} \rangle$

C言語の構文記法：

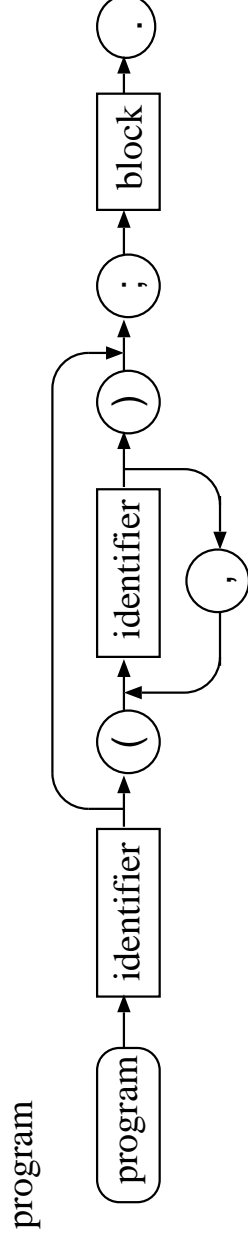
*identifier-list*:

*identifier*

*identifier-list* , *identifier*

**for** ( *expr opt* ; *expr opt* ; *expr opt* ) *stat*

Pascalの構文図式 (syntax diagram) の例：



## 文法

生成規則： $A \rightarrow \alpha$  ( $A$ は記号， $\alpha$ は記号列)

$\varepsilon$ 規則： $A \rightarrow \varepsilon$

文法：組 $\langle P, S \rangle$

$P$ ：生成規則の集合

$S$ ：出発記号 (start symbol)

語彙 (vocabulary)  $V$ ： $P$ に現れる記号全体の集合

非終端記号  $V_N$ ：生成規則の左辺に現れる記号

終端記号  $V_T$ ：非終端記号以外の記号

$$V = V_N \cup V_T$$

$$V_N \cap V_T = \phi$$

$$S \in V_N$$

例：文法  $G_1 = \langle P_1, E \rangle =$

$\langle V_N, V_T, P_1, E \rangle$

$$P_1 = \{ E \rightarrow E+T$$

$$E \rightarrow T$$

$$T \rightarrow T*F$$

$$T \rightarrow F$$

$$F \rightarrow (E)$$

$$F \rightarrow i \}$$

$$V_N = \{E, F, T\}$$

$$V_T = \{+, *, (, ), i\}$$

$$V = V_N \cup V_T$$

$$= \{E, F, T, +, *, (, ), i\}$$

$$P_1 = \{ E \rightarrow E+T \mid T$$

$$T \rightarrow T*F \mid F$$

$$F \rightarrow (E) \mid i \}$$

$v = xAy$ ,  $w = x\alpha y$ ,  
 $A \rightarrow \alpha \in P$  のとき,  $v \Rightarrow w$   
 $v$  は  $w$  を直接生成する.  
 $w$  は  $v$  に直接還元される.

$v \Rightarrow u_1 \Rightarrow u_2 \cdots \Rightarrow u_n \Rightarrow w$  のとき,  $v \stackrel{+}{\Rightarrow} w$

$v \stackrel{+}{\Rightarrow} w$  または  $v = w$  のとき,  $v \stackrel{*}{\Rightarrow} w$

$v$  は  $w$  を生成 (produce) する.

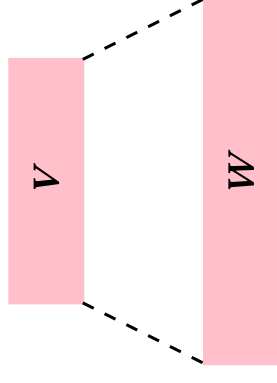
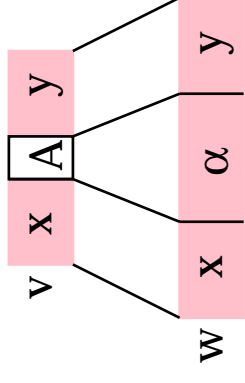
$w$  は  $v$  に還元 (reduce) される.

$S \stackrel{*}{\Rightarrow} x$  のとき,  $x$  は  $G$  の文形式 (sentential form)

特に  $x$  が終端記号列なら  $x$  は  $G$  の文 (sentence)

$L(G)$ :  $G$  の文全体の集合,  $G$  の定義する言語

例:  $P1 = \{ E \rightarrow E+T \mid T \quad E \Rightarrow E+T \Rightarrow T+T \Rightarrow T+T^*F \Rightarrow F+T^*F$   
 $T \rightarrow T^*F \mid F \quad T \Rightarrow F+T^*i \Rightarrow F+i^*i \Rightarrow F+i^*i+i^*i$   
 $F \rightarrow (E) \mid i \} \quad i+i^*i$  だけが文



導出 (derivation) :

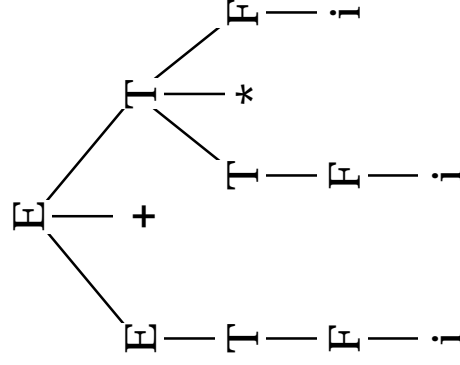
$v$  から  $v \stackrel{*}{\Rightarrow} w$  である記号列  $w$  を求めること

最左 (left-most) 導出

$$\begin{aligned}
 E &\stackrel{\text{lm}}{\Rightarrow} E+T \stackrel{\text{lm}}{\Rightarrow} T+T \stackrel{\text{lm}}{\Rightarrow} F+T \stackrel{\text{lm}}{\Rightarrow} i+T \stackrel{\text{lm}}{\Rightarrow} i+T*F \\
 &\stackrel{\text{lm}}{\Rightarrow} i+F*F \stackrel{\text{lm}}{\Rightarrow} i+i*F \stackrel{\text{lm}}{\Rightarrow} i+i*i
 \end{aligned}$$

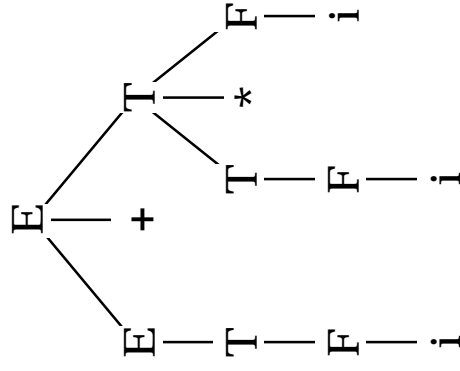
最右 (right-most) 導出

$$\begin{aligned}
 E &\stackrel{\text{rm}}{\Rightarrow} E+T \stackrel{\text{rm}}{\Rightarrow} E+T*F \stackrel{\text{rm}}{\Rightarrow} E+T*i \stackrel{\text{rm}}{\Rightarrow} E+F*i \\
 &\stackrel{\text{rm}}{\Rightarrow} E+i*i \stackrel{\text{rm}}{\Rightarrow} T+i*i \stackrel{\text{rm}}{\Rightarrow} F+i*i \stackrel{\text{rm}}{\Rightarrow} i+i*i
 \end{aligned}$$

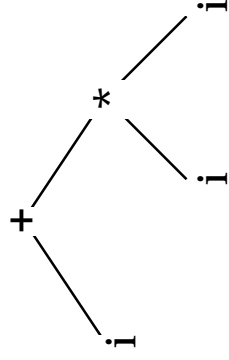


# 解析木とあいまい性

解析木



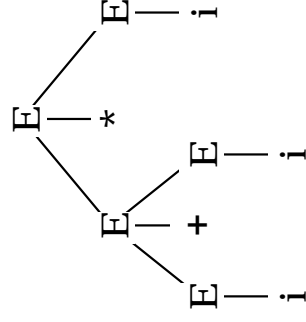
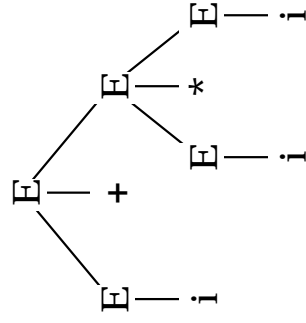
構文木



あいまいな文：二つの異なる解析木が考えられる文  
 あいまいな文法：あいまいな文を生成する文法

例：文法  $G_2 = \langle P_2, E \rangle$ ,  $P_2 = \{ E \rightarrow E + E \mid E * E \mid (E) \mid i \}$

文  $i + i * i$  に対する二つの解析木



## あいまいな英文

Time flies.

「主語 + 動詞」→「光陰矢の如し」

「動詞 + 目的語」→「蠅の速度計測をしる」

## あいまいな日本語文

ここではきものをぬぎなさい

「ここで履物を脱ぎなさい」

or 「ここでは着物を脱ぎなさい」?

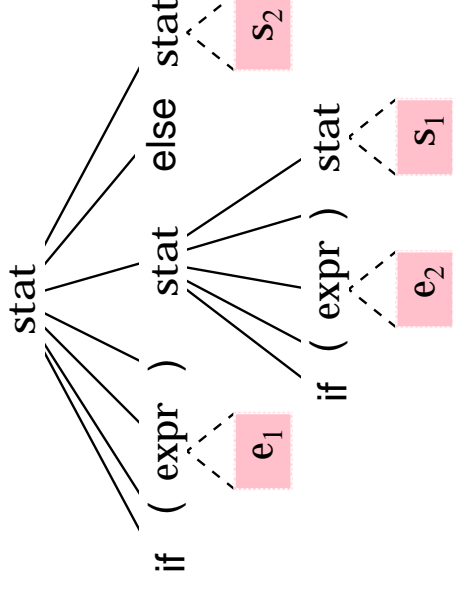
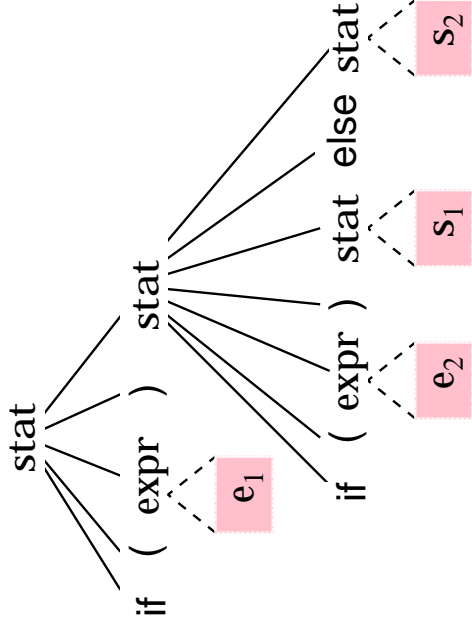


## if文の構文

```
stat:  
  if ( expr ) stat  
  if ( expr ) stat else stat  
  other
```

## あいまいなif文

```
if (  $e_1$  ) if (  $e_2$  )  $s_1$  else  $s_2$   
if (  $e_1$  ) { if (  $e_2$  )  $s_1$  } else  $s_2$ 
```



C言語の場合は左側

## if文のあいまい性除去

開いた ( open ) 文：直後に「 else + 文」をつけても

文になり得る文

閉じた ( closed ) 文：open でない文

```
stat:  
  closed-stat  
  open-stat  
closed-stat:  
  if ( expr ) closed-stat else closed-stat  
  other  
open-stat:  
  if ( expr ) stat  
  if ( expr ) closed-stat else open-stat
```

# 演算子の優先順位と結合性

## 演算子の優先順位

$$x + y * z = x + (y * z)$$

## 演算子の結合性

$$x - y - z = (x - y) - z$$

なので左結合的 (left associative)

$$x = y = z = x = (y = z)$$

なので右結合的

## 優先順位と結合性

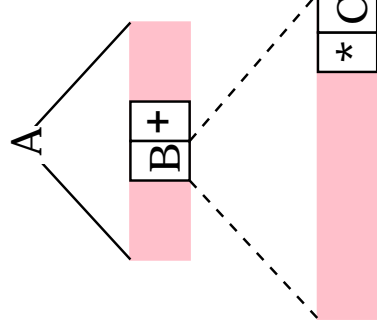
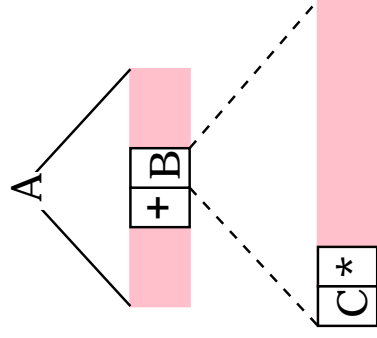
$$x + y - z = (x + y) - z$$

$op_1 \prec op_2 : \dots op_1 \text{ expr } op_2 \dots$  のとき  $op_2$  を先に計算

1. 生成規則  $A \rightarrow \dots op_1 B \dots$  が存在
2. ある非終端記号  $C$  に対して  $B \stackrel{+}{\Rightarrow} C op_2 \dots$

$op_1 \succ op_2 : \dots op_1 \text{ expr } op_2 \dots$  のとき  $op_1$  を先に計算

1. 生成規則  $A \rightarrow \dots B op_2 \dots$  が存在
2. ある非終端記号  $C$  に対して  $B \stackrel{+}{\Rightarrow} \dots op_1 C$



例：文法G1

1.  $E \rightarrow E+T$ が存在し，

2.  $E \Rightarrow E+T$

$\rightarrow + \triangleright +$

1.  $E \rightarrow E+T$ が存在し

2.  $T \Rightarrow T*F$

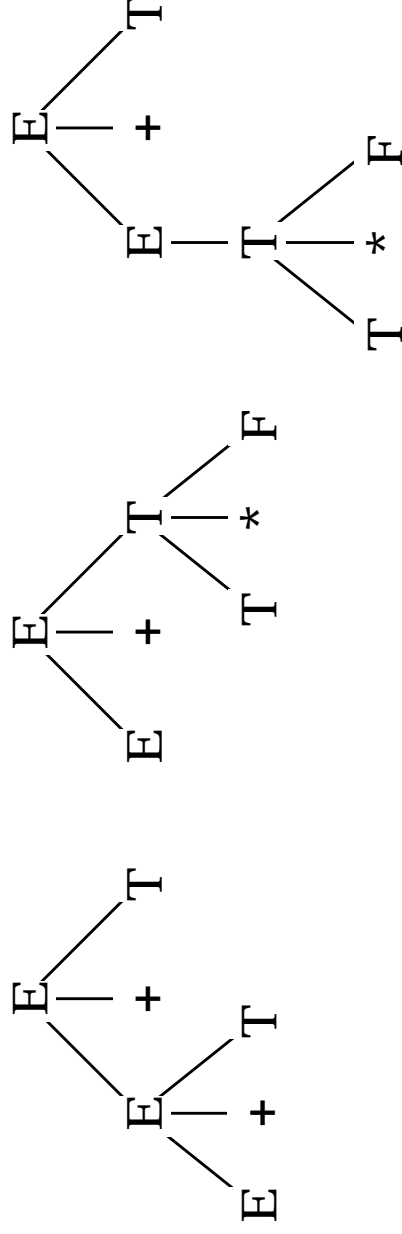
$\rightarrow + \triangleleft *$

1.  $E \rightarrow E+T$ が存在し，

2.  $E \Rightarrow T \Rightarrow T*F$

$\rightarrow * \triangleright +$

$+ \triangleleft * \text{かつ} * \triangleright +$ なので， $*$ は $+$ より優先順位が高い．



# 文脈自由文法とその限界

Aが文形式に現れていれば，Aの前後（文脈）に  
関係なく，Aを左辺とする生成規則を適用可能

文脈自由文法で表現できない記号列集合

$$Y = \{a^i b^i c^i \mid i \geq 1\} = \{abc, aabbcc, aaabbccc, \dots\}$$

$L(G) = Y$ なる文脈自由文法Gが存在したと仮定

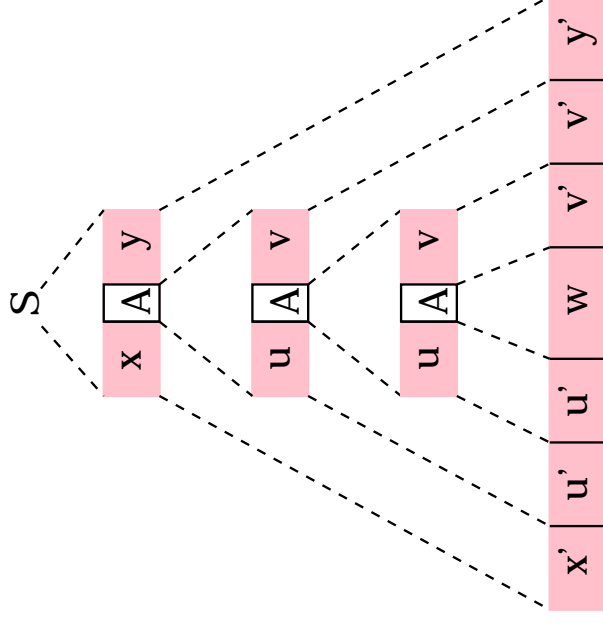
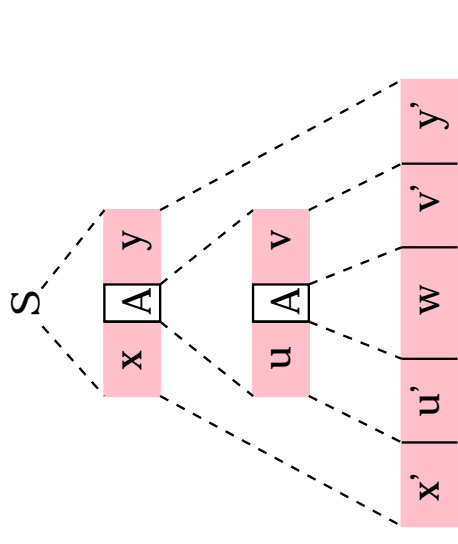
次の  $a^j b^j c^j$  が存在する .

$$S \stackrel{*}{\Rightarrow} xAy \stackrel{+}{\Rightarrow} xuAvy \stackrel{+}{\Rightarrow} a^j b^j c^j$$

Gが生成する記号列

$$S \stackrel{*}{\Rightarrow} xAy \stackrel{+}{\Rightarrow} xuAvy \stackrel{+}{\Rightarrow} xuuAvvy \stackrel{+}{\Rightarrow} x'u'u'wv'v'y'$$

は，Yの要素でない．  $\rightarrow$  矛盾



## Y を生成する文脈自由でない文法

$$S \rightarrow aBbc$$

$$B \rightarrow aBbC \mid \varepsilon$$

$$Cb \rightarrow bC$$

$$Cc \rightarrow cc$$

$$\begin{aligned} S &\Rightarrow aBbc \Rightarrow aaBbCbc \Rightarrow aaaBbCbCbc \\ &\Rightarrow aaabCbCbc \Rightarrow aaabCbCbc \Rightarrow aaabCbCbc \\ &\Rightarrow aaabbCbcc \Rightarrow aaabbCbcc \Rightarrow aaabbccccc \end{aligned}$$

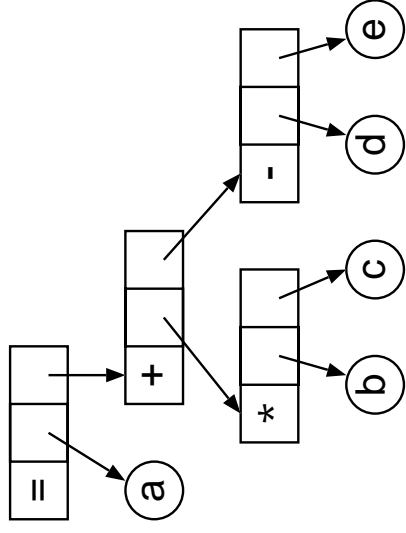
## 第4章

### 構文解析



# 構文木とその表現

N組 (N-tuple)



tuple : N組へのポインタの型

token : トークンへのポインタの型

tuple型は, token型を含む.

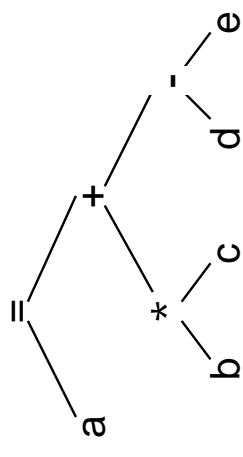
N組の生成

```
x1 = make_tuple("*", b, c);  
x2 = make_tuple("-", d, e);  
x3 = make_tuple("+", x1, x2);  
x4 = make_tuple("=", a, x3);
```

トークンの種類判定

```
int is_a(x)
```

$a = b * c + (d - e)$



# 再帰的下向き (recursive descent) 構文解析法

## 解析関数 (parsing function)

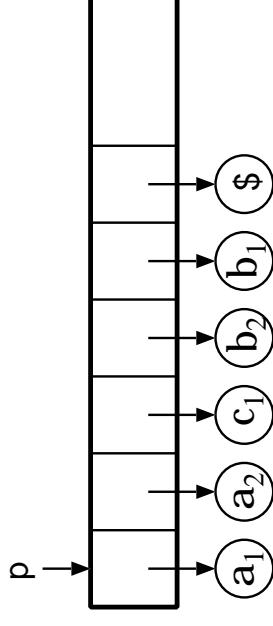
$A \rightarrow aAb \mid c$

```
tuple parse_A() {
    token *prev = p, x, z; tuple y;
    if (is_a(x = *p++) && (y = parse_A())
        && is_b(z = *p++))
        return make_tuple("A", x, y, z);
    p = prev;
    if (is_c(x = *p++)) return x;
    p = prev;
    return NULL;
}
```

\$は特殊な終端記号と見なす .

生成規則は再帰的

→ 再帰的な解析関数群



## 入力トークン列全体の解析

```
tuple parse() {  
    tuple x = parse_A();  
    if (x && is_eof(*p))  
        return x;  
    else  
        syntax_error();  
}
```

## 再帰的下向き構文解析法の問題点

左再帰 (left recursion)

バックトラック (backtracking)

## 左再帰

文法が左再帰を含む  $\rightarrow$  解析関数の再帰呼出しが止らない.

G1の生成規則:

$$E \rightarrow E+T \mid T$$

直接の左再帰 (immediate left recursion)

```
tuple parse_E() {
    token *prev = p; tuple x, y;
    if ((x = parse_E()) && is_plus(*p++))
        && (y = parse_T()))
        return make_tuple("+", x, y);
    p = prev;
    if (x = parse_T()) return x;
    p = prev;
    return NULL;
}
```

間接的な左再帰

$$A \rightarrow Ba \mid c$$
$$B \rightarrow Ab \mid d$$

## C 言語における左再帰

*identifier-list*:  
*identifier*  
*identifier-list*, *identifier*

右再帰に置き換え可能 .

*identifier-list*:  
*identifier*  
*identifier*, *identifier-list*

単純に右再帰に置き換えられない例

*additive-expression*:  
*multiplicative-expression*  
*additive-expression* + *multiplicative-expression*  
*additive-expression* - *multiplicative-expression*

左結合的が右結合的になってしまう .

$$a - b + c = (a - b) + c \neq a - (b + c) = a - b - c$$

C 言語の文法には , 間接的な左再帰は存在しない .

# 「直接の左再帰」の除去

$$A \rightarrow Aa \mid b$$

$$A \xRightarrow{\text{lm}} Aa \xRightarrow{\text{lm}} Aaaa \xRightarrow{\text{lm}} \dots \xRightarrow{\text{lm}} Aa^n \xRightarrow{\text{lm}} ba^n$$

$$A \rightarrow bA'$$

$$A' \rightarrow aA' \mid \varepsilon$$

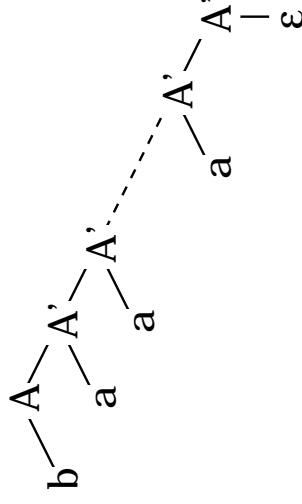
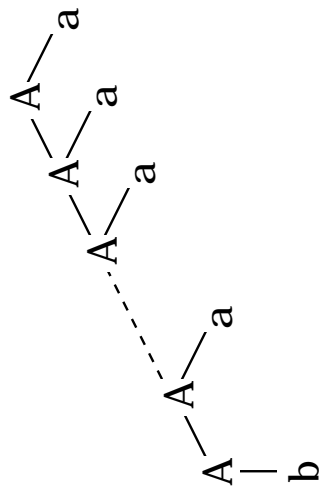
$$A \xRightarrow{\text{rm}} bA' \xRightarrow{\text{rm}} baA' \xRightarrow{\text{rm}} baaA' \xRightarrow{\text{rm}} \dots \xRightarrow{\text{rm}} ba^n A' \xRightarrow{\text{rm}} ba^n$$

```

tuple parse_A() {
    token *prev = p, y;
    if (is_b(y = *p++)) return parse_A1(y);
    p = prev;
    return NULL;
}

tuple parse_A1(tuple x) {
    token *prev = p, y;
    if (is_a(y = *p++))
        return parse_A1(make_tuple("A", x, y));
    p = prev;
    return x;
}

```



## 一般の場合

$$A \rightarrow A\alpha_1 \mid \cdots \mid A\alpha_n \mid \beta_1 \mid \cdots \mid \beta_m$$

ここで

$\beta_j \in V^*$  ,  $\beta_j$  は  $A$  で始まらない .

$\alpha_i \in V^+$

$$A \rightarrow \beta_1 A' \mid \cdots \mid \beta_m A'$$

$$A' \rightarrow \alpha_1 A' \mid \cdots \mid \alpha_n A' \mid \varepsilon$$

例 : G1の生成規則

$$E \rightarrow E+T \mid T$$

$$T \rightarrow T*F \mid F$$

$$F \rightarrow (E) \mid i$$

から左再帰を除去した文法G3

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \varepsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' \mid \varepsilon$$

$$F \rightarrow (E) \mid i$$

## バックトラック

複数の選択肢があるとき，その一つを試してみて，失敗だったらもとの状態に戻してやり直す．

例

$A \rightarrow aBc$   
 $B \rightarrow b \mid bd$

```
tuple parse_A() {
    token *prev = p, x, z; tuple y;
    if (is_a(x = *p++) && (y = parse_B())
        && is_c(z = *p++))
        return make_tuple("A", x, y, z);
    p = prev;
    return NULL;
}
tuple parse_B() {
    token *prev = p, x, y;
    if (is_b(x = *p++) return x;
    p = prev;
    if (is_b(x = *p++) && is_d(y = *p++))
        return make_tuple("B", x, y);
    p = prev;
    return NULL;
}
```

正しく動作しない(例： $abcd$ )  $\rightarrow$  一般的には複雑な機構が必要



くくり出し (factoring) によるバックトラックの回避

数学では「 $ax + ay$ の $a$ をくくり出して $a(x + y)$ 」

$B \rightarrow b \mid bd$ の $b$ をくくり出す .

```
B → b (d | ε)      tuple parse_B() {
    token *prev = p, x;
    if (is_b(x = *p++)) {
        token *prev1 = p, y;
        if (is_d(y = *p++))
            return make_tuple("B", x, y);
        p = prev1;
        return x;
    }
    p = prev;
    return NULL;
}
```

生成規則が終端記号で始まらない場合は適用が困難 .

$A \rightarrow aBc \quad B \rightarrow C \mid D \quad C \rightarrow \dots \quad D \rightarrow \dots$

バックトラックのもう一つの大きな問題：  
適切なエラーメッセージ出力が難しい

## LL( $k$ )構文解析法

$k$ 個のトークンを先読みし，どの生成規則を適用するかを決定する

$B \rightarrow C \mid D$ の解析関数で，

次の入力トークン  $x$  を先読みし，

1. もし  $C \stackrel{+}{\Rightarrow} x \dots$  なら `parse_C` を呼び出す．
2. もし  $D \stackrel{+}{\Rightarrow} x \dots$  なら `parse_D` を呼び出す．
3. どちらでもなければ構文エラー

# LL(1) 文法

文法  $G = \langle V_N, V_T, P, S \rangle$

$\text{First}(\alpha)$  : 記号列  $\alpha$  の first 集合

$$\text{First}(\alpha) = \{a \mid a \in V_T, \alpha \xRightarrow{*} a \dots\}$$

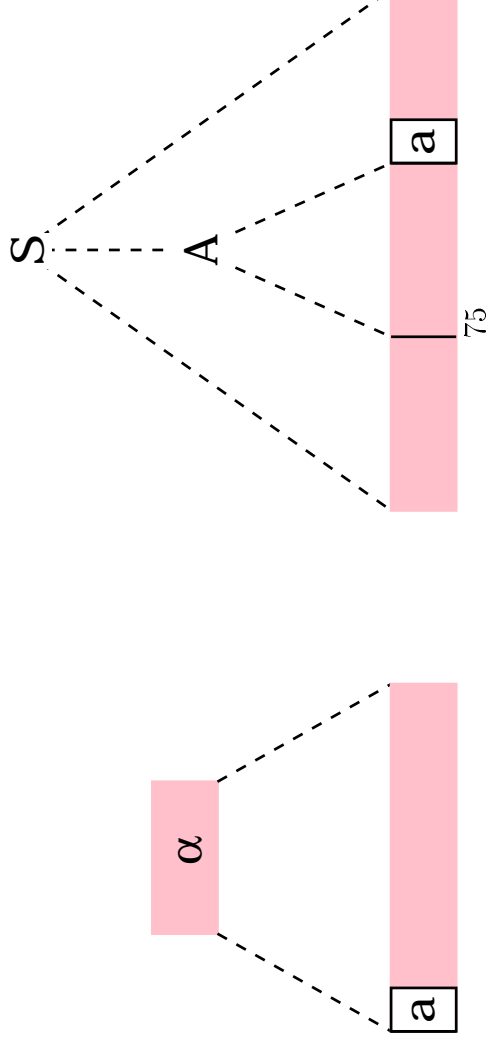
$\text{Follow}(A)$  : 非終端記号  $A$  の follow 集合

$$\text{Follow}(A) = \{a \mid a \in V_T, S \xRightarrow{*} \dots Aa \dots\}$$

$\$ \in \text{Follow}(A)$  なら ,

$$S \xRightarrow{*} \dots A$$

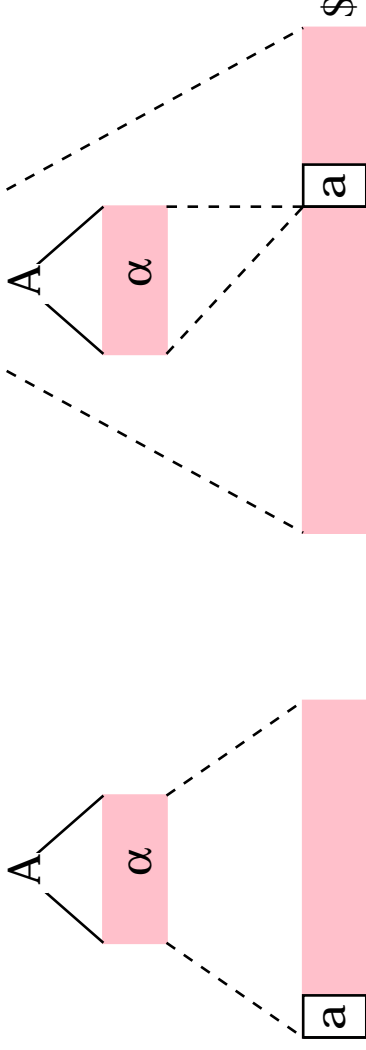
常に ,  $\$ \in \text{Follow}(S)$



Director( $A, \alpha$ ) : 生成規則  $A \rightarrow \alpha$  の director 集合

$$\text{Director}(A, \alpha) = \begin{cases} \text{First}(\alpha) & (\alpha \xrightarrow{*} \varepsilon \text{ でない}) \\ \text{First}(\alpha) \cup \text{Follow}(A) & (\alpha \xrightarrow{*} \varepsilon \text{ のとき}) \end{cases}$$

$A \rightarrow \alpha$  で還元するとき、入力トークン列の先頭に現れる可能性のある終端記号の全体



LL(1) 文法 : 右辺だけが異なる任意の生成規則

$A \rightarrow \alpha$  と

$A \rightarrow \beta$  に対して

必ず

$$\text{Director}(A, \alpha) \cap \text{Director}(A, \beta) = \phi$$

が成り立つ文法

## LL(1)文法のための計算

記号列  $\alpha$  が  $\varepsilon$  を生成する ( $\alpha \stackrel{*}{\Rightarrow} \varepsilon$ ) かどうかの判定

$\alpha = \varepsilon$  なら  $\alpha \stackrel{*}{\Rightarrow} \varepsilon$

$\alpha = a_1 a_2 \cdots a_n \neq \varepsilon$  なら  $\alpha \stackrel{*}{\Rightarrow} \varepsilon$  は,

すべての  $a_i$  に対して  $a_i \stackrel{*}{\Rightarrow} \varepsilon$

と同値

$\text{Null}(a)$  : 記号  $a$  に対して  $a \stackrel{*}{\Rightarrow} \varepsilon$  かどうか

$$\text{Null}(a) = \begin{cases} \text{true} & (a \stackrel{*}{\Rightarrow} \varepsilon \text{ のとき}) \\ \text{false} & (a \stackrel{*}{\Rightarrow} \varepsilon \text{ でないとき}) \end{cases}$$

1. すべての記号  $a$  に対して  $\text{Null}(a) \leftarrow \text{false}$
2.  $\varepsilon$  規則  $A \rightarrow \varepsilon$  に対して,  $\text{Null}(A) \leftarrow \text{true}$
3.  $\varepsilon$  規則以外の生成規則  $A \rightarrow a_1 \cdots a_n$  に対して,  
すべての  $a_i$  について  $\text{Null}(a_i) = \text{true}$  ならば,  
 $\text{Null}(A) \leftarrow \text{true}$
4. 変更がなくなるまで, ステップ 3 を繰り返す.

例：G1から左再帰を除去した文法G3

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' \mid \varepsilon \\ T &\rightarrow FT' \\ T' &\rightarrow *FT' \mid \varepsilon \\ F &\rightarrow (E) \mid i \end{aligned}$$

$\text{Null}(E') = \text{Null}(T') = \text{true}$

$\text{Null}(E) = \text{Null}(T) = \text{Null}(F) = \text{false}$

## first 集合の計算

$\alpha = \varepsilon$  なら  $\text{First}(\alpha) = \phi$

$\alpha = a_1 a_2 \cdots a_n \neq \varepsilon$  に対しては ,

$$\text{First}(\alpha) = \begin{cases} \text{First}(a_1) & (a_1 \stackrel{*}{\Rightarrow} \varepsilon \text{ でない}) \\ \text{First}(a_1) \cup \text{First}(a_2 \cdots a_n) & (a_1 \stackrel{*}{\Rightarrow} \varepsilon \text{ のとき}) \end{cases}$$

→ 記号についての first 集合から求められる .

1. 終端記号  $x$  に対して  $\text{First}(x) \leftarrow \{x\}$   
非終端記号  $A$  に対して ,  $\text{First}(A) \leftarrow \phi$
2. 生成規則  $A \rightarrow \alpha_1 a \alpha_2$  ( $a \in (V_N \cup V_T)$ ) に対して ,  
 $\alpha_1 \stackrel{*}{\Rightarrow} \varepsilon$  なら  $\text{First}(A) \leftarrow \text{First}(A) \cup \text{First}(a)$
3. 変更がなくなるまで , ステップ 2 を繰り返し返す .

## 例：文法G3のfirst集合

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' \mid \varepsilon \\ T &\rightarrow FT' \\ T' &\rightarrow *FT' \mid \varepsilon \\ F &\rightarrow (E) \mid i \end{aligned}$$

生成規則  $F \rightarrow (E)$  と  $F \rightarrow i$  から

$$\text{First}(F) = \{ (, i \}$$

生成規則  $T' \rightarrow *FT'$  から

$$\text{First}(T') = \{ * \}$$

$T \rightarrow FT'$  から (  $\text{Null}(F) = \text{false}$  なので )

$$\text{First}(T) = \text{First}(F) = \{ (, i \}$$

$\text{First}(T')$  ,  $\text{First}(T)$  と同様に ,

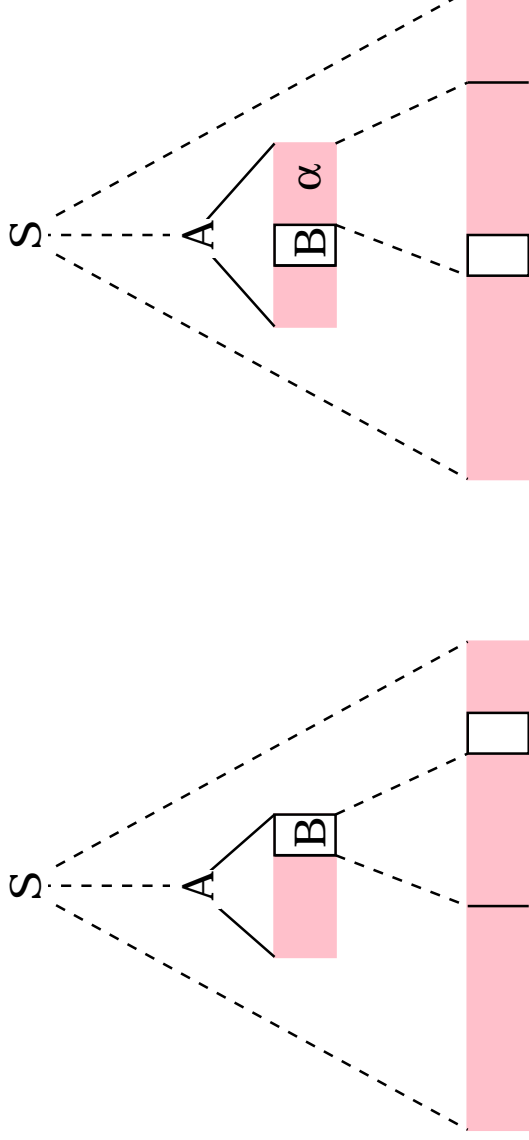
$$\text{First}(E') = \{ + \}$$

$$\text{First}(E) = \text{First}(T) = \{ (, i \}$$



## follow 集合の計算

1. 出発記号  $S$  に対して  $\text{Follow}(S) \leftarrow \{\$ \}$   
 $S$  以外の非終端記号  $A$  に対して  $\text{Follow}(A) \leftarrow \phi$
2. 生成規則  $A \rightarrow \dots B$  ( $B \in V_N$ ) に対して,  
 $\text{Follow}(B) \leftarrow \text{Follow}(B) \cup \text{Follow}(A)$
3. 生成規則  $A \rightarrow \dots B\alpha$  ( $B \in V_N, \alpha \neq \epsilon$ ) に対して,  
 $\text{Follow}(B) \leftarrow \text{Follow}(B) \cup \text{First}(\alpha)$   
もし  $\alpha \stackrel{*}{\Rightarrow} \epsilon$  なら, さらに  
 $\text{Follow}(B) \leftarrow \text{Follow}(B) \cup \text{Follow}(A)$
4. 変更がなくなるまでステップ 2 と 3 を繰り返す.



## 例：文法 G3 の follow 集合

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \varepsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' \mid \varepsilon$$

$$F \rightarrow (E) \mid i$$

出発記号  $E$  が右辺に現れる規則は  $F \rightarrow (E)$  だけ

$$\text{Follow}(E) = \{ \$, ) \}$$

$E'$  が右辺に現れる生成規則は  $E \rightarrow TE'$  と  $E' \rightarrow +TE'$

$$\text{Follow}(E') = \text{Follow}(E) = \{ \$, ) \}$$

$T$  が右辺に現れる生成規則も  $E \rightarrow TE'$  と  $E' \rightarrow +TE'$

$T$  の直後は  $E'$  ,  $\text{Null}(E') = \text{true}$  なので ,

$$\text{Follow}(T) = \text{First}(E') \cup \text{Follow}(E') = \{ +, \$, ) \}$$

$\text{Follow}(E')$  ,  $\text{Follow}(T)$  と同様に ,

$$\text{Follow}(T') = \text{Follow}(T) = \{ +, \$, ) \}$$

$$\text{Follow}(F) = \text{First}(T') \cup \text{Follow}(T') = \{ *, +, \$, ) \}$$

## 文法G3における director 集合

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \varepsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' \mid \varepsilon$$

$$F \rightarrow (E) \mid i$$

$$\text{Null}(E') = \text{Null}(T') = \text{true}$$

$$\text{Null}(E) = \text{Null}(T) = \text{Null}(F) = \text{false}$$

$$\text{First}(E) = \text{First}(T) = \text{First}(F) = \{(, i\}$$

$$\text{First}(T') = \{*\}$$

$$\text{First}(E') = \{+\}$$

$$\text{Follow}(E) = \text{Follow}(E') = \{\$, \}$$

$$\text{Follow}(T) = \text{Follow}(T') = \{+, \$, \}$$

$$\text{Follow}(F) = \text{First}(T') \cup \text{Follow}(T') = \{*, +, \$, \}$$

$$\text{Director}(E', +TE') = \{+\}$$

$$\text{Director}(E', \varepsilon) = \text{First}(\varepsilon) \cup \text{Follow}(E') = \{\$, \}$$

$$\text{Director}(T', *FT') = \{*\}$$

$$\text{Director}(T', \varepsilon) = \text{First}(\varepsilon) \cup \text{Follow}(T') = \{+, \$, \}$$

$$\text{Director}(F, (E)) = \{(}$$

$$\text{Director}(F, i) = \{i\}$$

$$\text{Director}(E', +TE') \cap \text{Director}(E', \varepsilon) = \phi$$

$$\text{Director}(T', *FT') \cap \text{Director}(T', \varepsilon) = \phi$$

$$\text{Director}(F, (E)) \cap \text{Director}(F, i) = \phi$$

G3はLL(1)文法である。

## LL(1)構文解析法 (Left-to-right scanning, Left-most derivation)

### バックトラックのない再帰的下向き構文解析

#### G3のための解析関数

```
tuple parse_E() {
    if (is_lpar(*p)
        || is_i(*p)) {
        tuple x = parse_T();
        return parse_E1(x);
    } else
        syntax_error();
}

tuple parse_E1(tuple x) {
    if (is_plus(*p)) {
        tuple y;
        p++;
        y = parse_T();
        return parse_E1(make_tuple("+", x, y));
    } else if (is_eof(*p) || is_rpar(*p))
        return x;
    else
        syntax_error();
}
```

バックトラックがなくなったので

1. 適切な時点で構文エラーを検出できる
2. pの値を保存しておく必要がない。

LR構文解析 (Left-to-right scanning, Right-most derivation in reverse)

文法  $\langle \{E \rightarrow E+i \mid i\}, E \rangle$  における  $i+i+i$  の最右導出

$$E \xRightarrow{\text{rm}} E+i \xRightarrow{\text{rm}} E+i+i \xRightarrow{\text{rm}} i+i+i$$

注目点挿入

$$\begin{aligned} E \cdot &\xRightarrow{\text{rm}} E+i \cdot \rightarrow E+\cdot i \rightarrow E \cdot \cdot +i \\ &\xRightarrow{\text{rm}} E+i \cdot \cdot +i \rightarrow E+\cdot \cdot i+i \rightarrow E \cdot \cdot \cdot +i+i \\ &\xRightarrow{\text{rm}} i \cdot \cdot +i+i \rightarrow \cdot i+i+i \end{aligned}$$

$i+i+i$  の LR構文解析

$$\begin{aligned} \cdot i+i+i \$ &\rightarrow i \cdot \cdot +i+i \$ \rightarrow E \cdot \cdot +i+i \$ \rightarrow E+\cdot \cdot i+i \$ \rightarrow E \cdot \cdot \cdot +i \$ \\ &\rightarrow E+\cdot \cdot i \$ \rightarrow E+i \cdot \cdot \$ \rightarrow E \cdot \cdot \$ \end{aligned}$$

各ステップでは、

還元 (reduce) : 注目点直前の記号列を還元

シフト (shift) : 注目点を一つ右に移動

$\cdot i + i + i \$ \rightarrow i \cdot \cdot + i + i \$ \rightarrow E \cdot \cdot + i + i \$ \rightarrow E + i \cdot \cdot + i \$ \rightarrow E + i + i \$ \rightarrow E \cdot \cdot + i \$$   
 $\rightarrow E + \cdot i \$ \rightarrow E + i \cdot \$ \rightarrow E \cdot \$$

$A \rightarrow \alpha$  で還元するとき

$A$  の節を用意し,  $\alpha$  の各記号の節へ枝を張る.

上向き構文解析法 (bottom-up parsing)

初期状態:  $\cdot x \$$

$x$ : 入力トークン列

解析中の状態:  $u \cdot v \$$

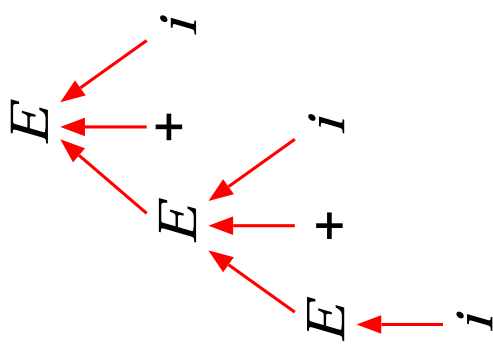
$u \in (V_T \cup V_N)^*$ : 構築した部分解析木の記号列

$v \in V_T^*$ : 未処理の入力トークン列

最終状態:  $S \cdot \$$

$S$ : 出発記号

最終状態に変換できれば, 構文エラー.



## LR(0) 項

生成規則の右辺にドットをつけたもの

例:  $E \rightarrow E+i$  の LR(0) 項:

$E \rightarrow \cdot E+i$  ( 導入項 )

$E+i$  を読み込めば  $E$  に還元できる .

$E \rightarrow E \cdot +i$

$E$  を読み込んだ .  $+i$  を読み込めば  $E$  に還元できる .

$E \rightarrow E+i \cdot$

$E+i$  を読み込んだ .  $i$  を読み込めば  $E$  に還元できる .

$E \rightarrow E+i \cdot$  ( 完全項 )

$E+i$  を読み込んだ .  $E$  に還元できる .

$E$ : 文法上の出発記号

$S$ : 新しい出発記号

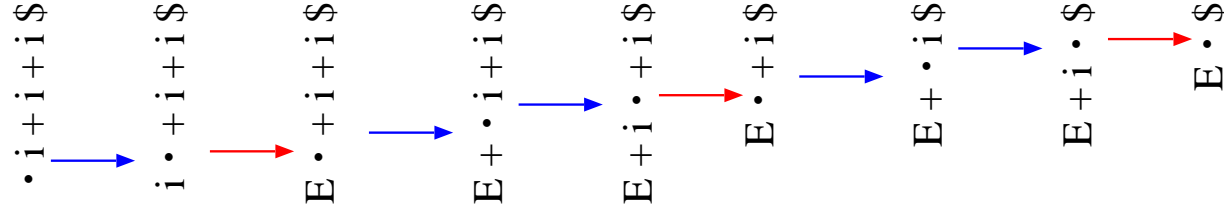
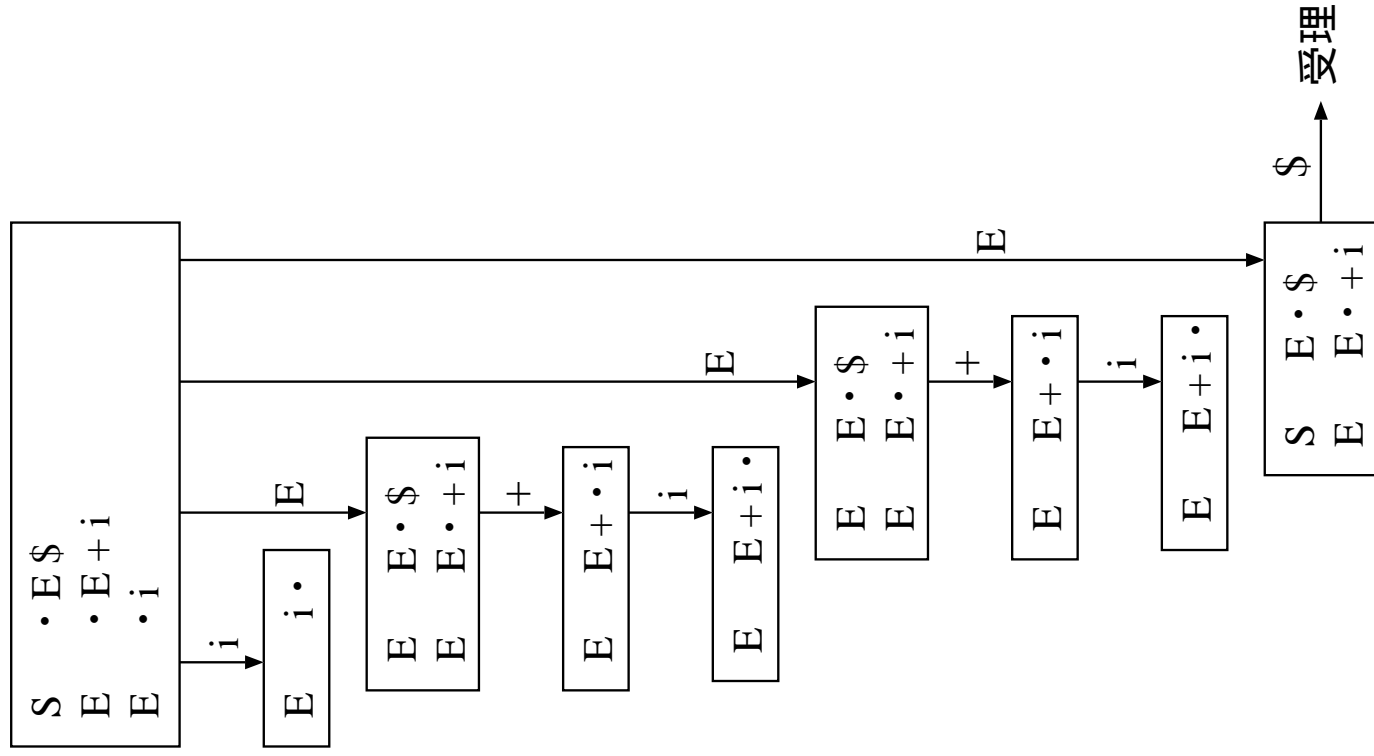
$S \rightarrow \cdot E\$$  ( 出発項 )

$E$  を読み込んで  $\$$  に達すれば受理できる .

$S \rightarrow E \cdot \$$  ( 受理項 )

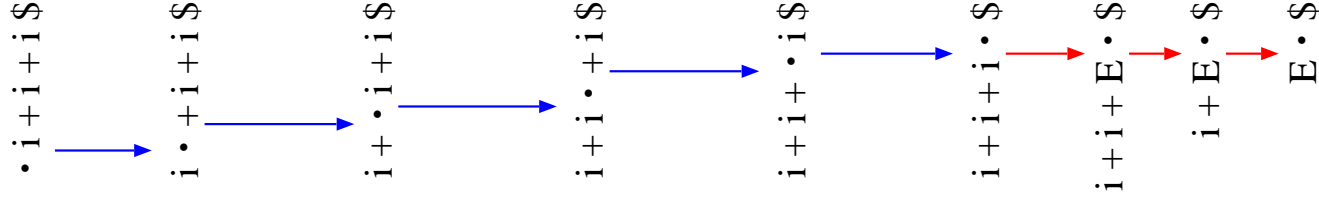
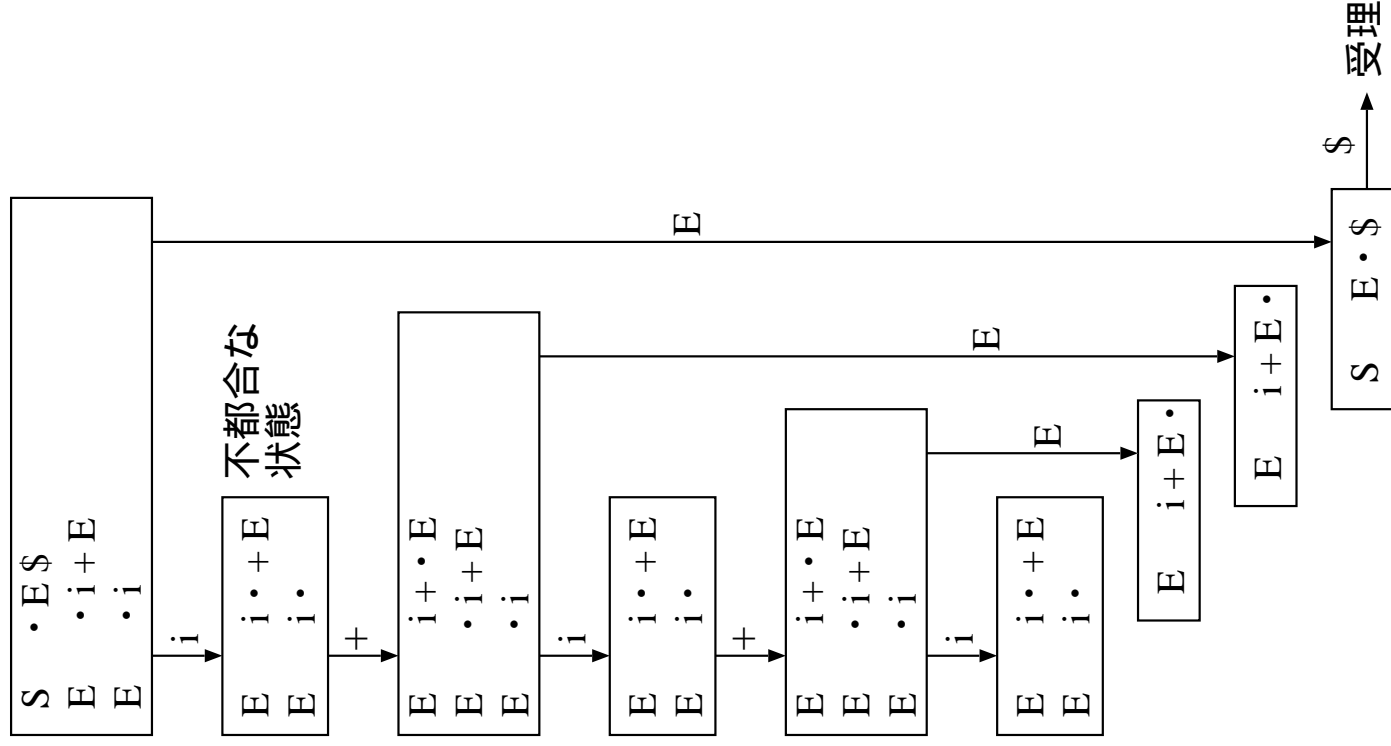
$E$  を読み込んだ .  $\$$  に達していれば受理できる .

# 文法 $\langle \{E \rightarrow E+i \mid i\}, E \rangle$ における LR 構文解析





# 文法 $\langle \{E \rightarrow i + E \mid i\}, E \rangle$ における LR 構文解析



$I$ : LR(0) 項集合

$\text{Closure}(I)$  :  $I$  の閉包 (closure)

- $\text{Closure}(I) \supseteq I$
- $A \rightarrow \alpha_1 \cdot B \alpha_2$  が  $\text{Closure}(I)$  の要素なら ,  
任意の導入項  $B \rightarrow \cdot \beta$  も ,  $\text{Closure}(I)$  の要素

例 : 文法  $\langle \{E \rightarrow i + E \mid i\}, E \rangle$

$$\begin{aligned} & \text{Closure}(S \rightarrow \cdot E \$) \\ &= \{S \rightarrow \cdot E \$, E \rightarrow \cdot i + E, E \rightarrow \cdot i\} \end{aligned}$$

$\text{Goto}(I, a)$

$$= \text{Closure}(\{A \rightarrow \alpha_1 a \cdot \alpha_2 \mid (A \rightarrow \alpha_1 \cdot a \alpha_2) \in I\})$$

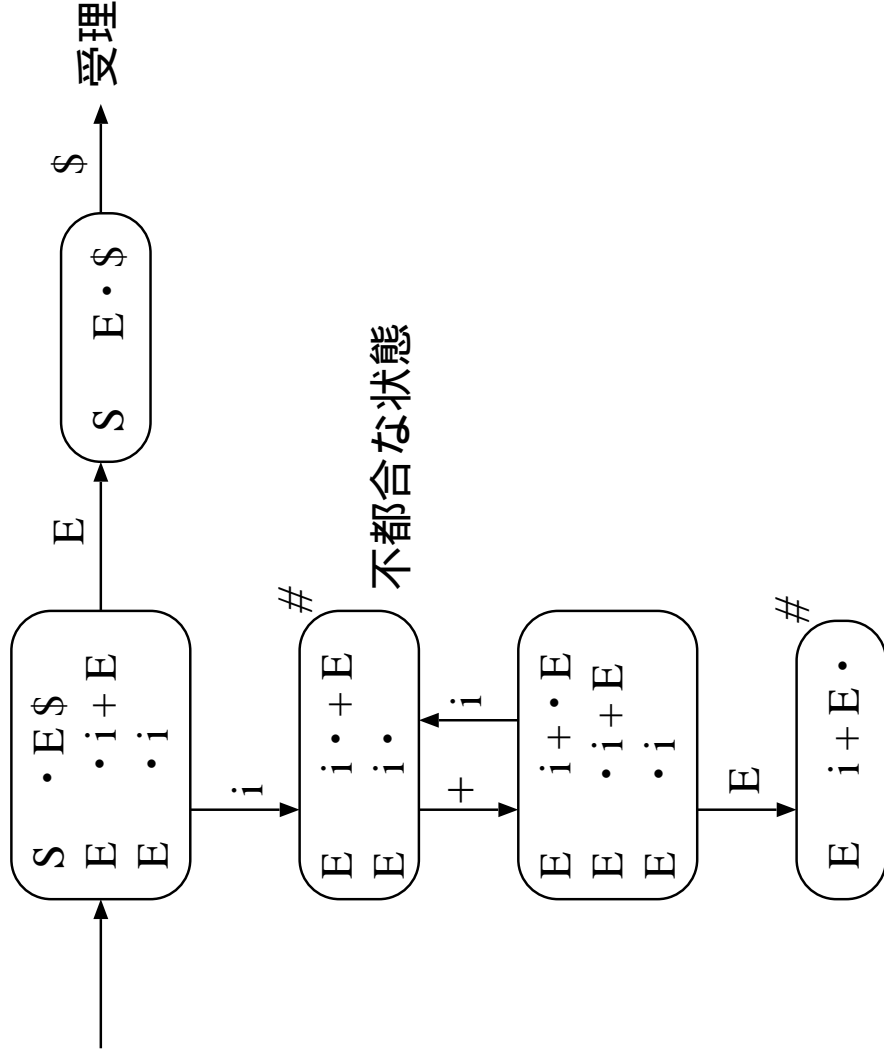
例 : 文法  $\langle \{E \rightarrow i + E \mid i\}, E \rangle$

$$\begin{aligned} & \text{Goto}(\{E \rightarrow i \cdot + E, E \rightarrow i \cdot\}, +) \\ &= \text{Closure}(\{E \rightarrow i + \cdot E\}) \\ &= \{E \rightarrow i + \cdot E, E \rightarrow \cdot i + E, E \rightarrow \cdot i\} \end{aligned}$$

# 正準オートマトン (canonical automaton)

初期状態  $\text{Closure}(\{S \rightarrow \cdot E \$\})$  から始め, 可能な遷移をすべて求める

例: 文法  $\langle \{E \rightarrow i + E \mid i\}, E \rangle$  の正準オートマトン



#: 還元状態 (完全項を含む状態)

正準集合: 正準オートマトンの状態集合

## 配置 (configuration)

LR構文解析のある時点

$$a_1 a_2 \cdots a_n \cdot x_k x_{k+1} \cdots x_m \$$$

における配置は ,

$$(q_0 a_1 q_1 a_2 q_2 \cdots a_n q_n, x_k x_{k+1} \cdots x_m \$)$$

$$q_0 : \text{初期状態}, q_i = \text{Goto}(q_{i-1}, a_i) = \text{Goto}(q_0, a_1 a_2 \cdots a_n)$$

LR構文解析開始時点の配置 :

$$(q_0, x_1 x_2 \cdots x_m \$)$$

$A \rightarrow a_i \cdots a_n$  による還元後 :

$$(q_0 a_1 q_1 \cdots a_{i-1} q_{i-1} A q, x_k x_{k+1} \cdots x_m \$)$$

$$q = \text{Goto}(q_{i-1}, A)$$

$x_k$  を読み込んでシフト後 :

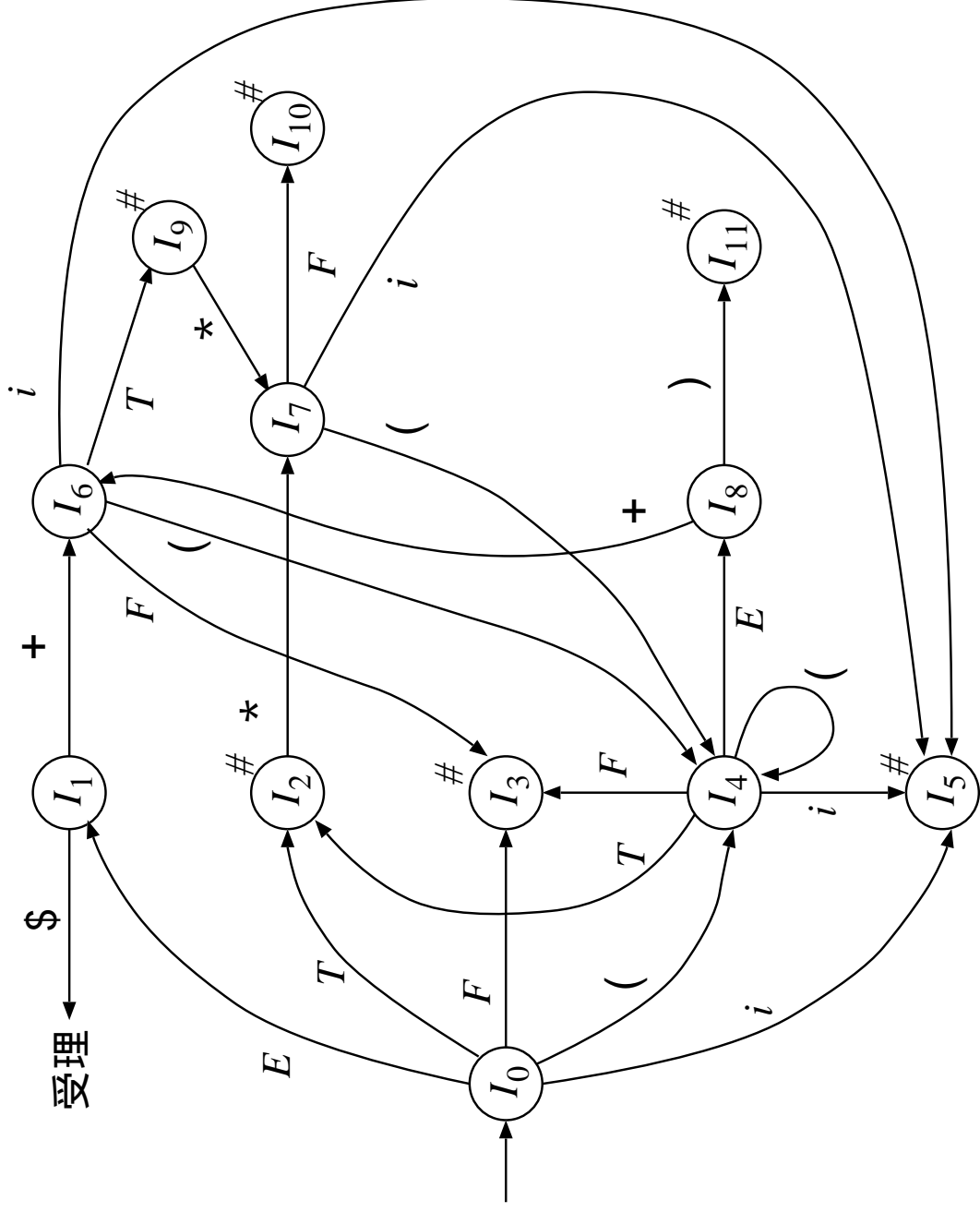
$$(q_0 a_1 q_1 \cdots a_n q_n x_k q, x_{k+1} \cdots x_m \$)$$

$$q = \text{Goto}(q_n, x_k)$$

最終的な配置 :

$$(q_0 S q, \$)$$

例：文法  $G1 = \langle P1, E \rangle$  の標準オートマトン



$$P1 = \{ E \rightarrow E+T \mid T \\ T \rightarrow T*F \mid F \\ F \rightarrow (E) \mid i \}$$

初期状態  $I_0$

$= \text{Closure}(\{S \rightarrow \cdot E \$\})$

$\#$ : 還元状態

不都合な状態 (とりあえずシフトを優先):

$$I_2 = \{ E \rightarrow T \cdot, T \rightarrow T \cdot * F \}$$

$$I_9 = \{ E \rightarrow E+T \cdot, T \rightarrow T \cdot * F \}$$

## G1の正準集合

$$I_0 = \{ S \rightarrow \cdot E \$, \\ E \rightarrow \cdot E + T, \\ E \rightarrow \cdot T, \\ T \rightarrow \cdot T * F, \\ T \rightarrow \cdot F, \\ F \rightarrow \cdot (E), \\ F \rightarrow \cdot i \}$$

$$I_1 = \{ S \rightarrow E \cdot \$, \\ E \rightarrow E \cdot + T \}$$

$$I_2 = \{ E \rightarrow T \cdot, \\ T \rightarrow T \cdot * F \}$$

$$I_3 = \{ T \rightarrow F \cdot \}$$

$$I_4 = \{ F \rightarrow (\cdot E), \\ E \rightarrow \cdot E + T, \\ E \rightarrow \cdot T \\ T \rightarrow \cdot T * F, \\ T \rightarrow \cdot F, \\ F \rightarrow \cdot (E), \\ F \rightarrow \cdot i \}$$

$$I_5 = \{ F \rightarrow i \cdot \}$$

$$I_6 = \{ E \rightarrow E + \cdot T, \\ T \rightarrow \cdot T * F, \\ T \rightarrow \cdot F, \\ F \rightarrow \cdot (E), \\ F \rightarrow \cdot i \}$$

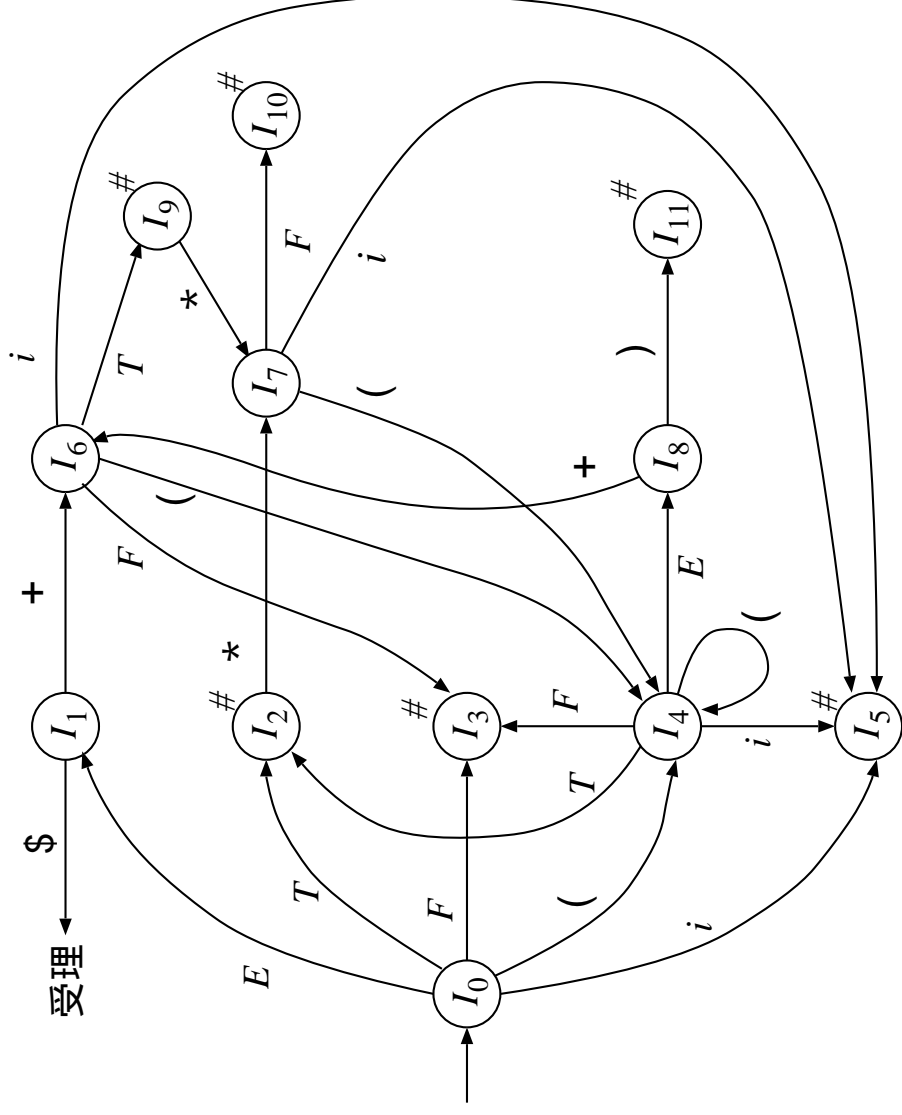
$$I_7 = \{ T \rightarrow T * \cdot F, \\ F \rightarrow \cdot (E), \\ F \rightarrow \cdot i \}$$

$$I_8 = \{ F \rightarrow (E \cdot), \\ E \rightarrow E \cdot + T \}$$

$$I_9 = \{ E \rightarrow E + T \cdot, \\ T \rightarrow T \cdot * F \}$$

$$I_{10} = \{ T \rightarrow T * F \cdot \}$$

$$I_{11} = \{ F \rightarrow (E) \cdot \}$$



$(I_0,$	$i+i*i\$)$	$i$ の読み込み	$(I_0 E I_1 + I_6 F I_3,$	$*i\$)$	$T \rightarrow F$ で還元
$(I_0 i I_5,$	$+i*i\$)$	$F \rightarrow i$ で還元	$(I_0 E I_1 + I_6 T I_9,$	$*i\$)$	$'*'$ の読み込み
$(I_0 F I_3,$	$+i*i\$)$	$T \rightarrow F$ で還元	$(I_0 E I_1 + I_6 T I_9 * I_7,$	$i\$)$	$i$ の読み込み
$(I_0 T I_2,$	$+i*i\$)$	$E \rightarrow T$ で還元	$(I_0 E I_1 + I_6 T I_9 * I_7 i I_5,$	$\$)$	$F \rightarrow i$ で還元
$(I_0 E I_1,$	$+i*i\$)$	$'+'$ の読み込み	$(I_0 E I_1 + I_6 T I_9 * I_7 F I_{10},$	$\$)$	$T \rightarrow T * F$ で還元
$(I_0 E I_1 + I_6,$	$i*i\$)$	$i$ の読み込み	$(I_0 E I_1 + I_6 T I_9,$	$\$)$	$E \rightarrow E + T$ で還元
$(I_0 E I_1 + I_6 i I_5,$	$*i\$)$	$F \rightarrow i$ で還元	$(I_0 E I_1,$	$\$)$	受理

# SLR(1) 構文解析

Simple LR,

1 文字先読み (lookahead) し ,  
follow 集合で不都合な状態を解消

SLR(1) 文法 : SLR(1) 構文解析が可能な文法

例 : G1 の不都合な状態  $I_2$

$$I_2 = \{ E \rightarrow T \cdot, T \rightarrow T \cdot * F \}$$

次の入力トークンが

1.  $\text{Follow}(E) = \{ +, ), \$ \}$  の要素なら ,  $E \rightarrow T$  で還元
2. ‘\*’ なら , ‘\*’ を読み込んでシフト
3. どちらでもなければ , 構文エラー

$$I_9 = \{ E \rightarrow E + T \cdot, T \rightarrow T \cdot * F \}$$

$$\text{Follow}(E) = \{ +, ), \$ \} \neq *$$

→ G1 は SLR(1) 文法



	$i$	$+$	$*$	$($	$)$	$\$$
$I_0$	s5			s4		
$I_1$		s6				受理
$I_2$		r2	s7		r2	r2
$I_3$		r4	r4		r4	r4
$I_4$	s5			s4		
$I_5$		r6	r6		r6	r6
$I_6$	s5			s4		
$I_7$	s5			s4		
$I_8$		s6			s11	
$I_9$		r1	s7		r1	r1
$I_{10}$		r3	r3		r3	r3
$I_{11}$		r5	r5		r5	r5

動作表

動作表 (action table)

$s_j$  : 状態  $I_j$  にシフト

$r_j$  :  $j$  番目の生成規則で還元

受理 : 構文解析終了

空欄 : 動作未定義  $\rightarrow$  構文エラー

	$E$	$T$	$F$
$I_0$	$I_1$	$I_2$	$I_3$
$I_4$	$I_8$	$I_2$	$I_3$
$I_6$		$I_9$	$I_3$
$I_7$			$I_{10}$

行先表

1:  $E \rightarrow E+T$

2:  $E \rightarrow T$

3:  $T \rightarrow T*F$

4:  $T \rightarrow F$

5:  $F \rightarrow (E)$

6:  $F \rightarrow i$

行先表 (goto table)

還元後の状態 ( $\text{Goto}(I, a)$ )

## SLR(1) 構文解析プログラム ( 文法に依存せず )

```
tuple parse() {
    int state = 0;
    for (;;) {
        action act = lookup_action(state, *p);
        if (act.todo == SHIFT) {
            push_state(state);
            push_value(*p++);
            state = act.num;
        } else if (act.todo == REDUCE) {
            int nt = reduce(act.num);
            state = lookup_goto(top_state(), nt);
        } else if (act.todo == ACCEPT)
            return pop_value();
        else /* if (act.todo == UNDEFINED) */
            syntax_error();
    }
}
```

配置  $(q_0 a_1 q_1 \cdots q_{n-1} a_n q_n, x_k x_{k+1} \cdots x_m \$)$  の表現 :

変数 state :  $q_n$

state スタック :  $q_0, q_1, \cdots, q_{n-1}$     value スタック :  $a_1, \cdots, a_n$

## 還元関数 ( 文法に依存 )

```
int reduce(int production) {
    tuple x, y;
    switch (production) {
        case 1: /* E -> E + T */
            x = pop_value(); pop_value();
            y = pop_value();
            pop_state(); pop_state();
            push_value(make_tuple("+", y, x));
        case 2: /* E -> T */
            return E;
        case 3: /* T -> T * F */ ...
        case 4: /* T -> F */ ...
        case 5: /* F -> ( E ) */ ...
        case 6: /* F -> i */ ...
    }
}
```

$E \rightarrow E + T$  による還元 :

$$(q_0 a_1 q_1 \cdots a_{n-3} q_{n-3} E_1 q_{n-2} + q_{n-1} T_1 q_n, \cdots) \rightarrow (q_0 a_1 q_1 \cdots a_{n-3} q_{n-3} E_2 q, \cdots)$$

$E_1$  の構文木は  $T_1$  と同じ

LR(1) や LALR(1) でも , プログラムは同じ . 動作表と行先表が異なる .

# LR(1) 構文解析

follow 集合より正確な先読み記号を使用

$A \rightarrow \alpha$  による還元

$$S \cdot \$ \xrightarrow{*}_{\text{rm}} uA \cdot v\$ \xrightarrow{\text{rm}} u\alpha \cdot v\$$$

次の入力  $x$  は,  $x \in \text{First}(v\$)$  のはず .

LR(1) 項 : 次の入力  $x$  なら  $A \rightarrow \alpha$  で還元可能

$$[A \rightarrow \alpha \cdot, x]$$

LR(1) 項の一般形 :

$$[A \rightarrow \alpha_1 \cdot \alpha_2, x]$$

完全項を核 (core) とする LR(1) 項

$$[A \rightarrow \alpha_1 \alpha_2 \cdot, x]$$

になって次の入力  $x$  なら,  $A \rightarrow \alpha_1 \alpha_2$  で還元可能

LR(1) 構文解析の初期状態 :

$$I_0 = \text{Closure}(\{ [S' \rightarrow \cdot S, \$] \})$$

$$\text{Goto}(I, a) = \text{Closure}(\{ [A \rightarrow \alpha_1 a \cdot \alpha_2, c] \mid [A \rightarrow \alpha_1 \cdot a \alpha_2, c] \in I \})$$

LR(1)項集合  $I$  の閉包  $I' = \text{Closure}(I)$

1.  $I' \leftarrow I$
2.  $I'$  の要素  $[A \rightarrow \alpha_1 \cdot B \alpha_2, x]$  ,  
生成規則  $B \rightarrow \beta$  ,  
 $x' \in \text{First}(\alpha_2 x)$  であるすべての  $x'$  について ,  
 $[B \rightarrow \cdot \beta, x']$  を  $I'$  に追加 .
3.  $I'$  に追加がなくなるまで , ステップ 2 を繰り返す .

例 : G1 を LR(1) 構文解析するときの初期状態

$$\begin{aligned} I_0 &= \text{Closure}(\{ [S \rightarrow \cdot E, \$] \}) \\ &= \{ [S \rightarrow \cdot E, \$], \\ &\quad [E \rightarrow \cdot E+T, \$], [E \rightarrow \cdot T, \$], [E \rightarrow \cdot E+T, +], [E \rightarrow \cdot T, +], \\ &\quad [T \rightarrow \cdot T*F, \$], [T \rightarrow \cdot F, \$], [T \rightarrow \cdot T*F, +], [T \rightarrow \cdot F, +], \\ &\quad [T \rightarrow \cdot T*F, *], [T \rightarrow \cdot F, *], [F \rightarrow \cdot (E), \$], [F \rightarrow \cdot i, \$], \\ &\quad [F \rightarrow \cdot (E), +], [F \rightarrow \cdot i, +], [F \rightarrow \cdot (E), *], [F \rightarrow \cdot i, *] \} \\ &= \{ [S \rightarrow \cdot E, \$], \\ &\quad [E \rightarrow \cdot E+T, \$/+], [E \rightarrow \cdot T, \$/+], \\ &\quad [T \rightarrow \cdot T*F, \$/+/*], [T \rightarrow \cdot F, \$/+/*], \\ &\quad [F \rightarrow \cdot (E), \$/+/*], [F \rightarrow \cdot i, \$/+/*] \} \end{aligned}$$

LR(1) 構文解析のための正準オートマトンの作り方：

SLR(1) の場合とまったく同じ .

例：G1 を LR(1) 構文解析するためのオートマトン

$$I_2 = \text{Goto}(I_0, T) = \text{Closure}(\{ [E \rightarrow T \cdot, \$ / +], [T \rightarrow T \cdot * F, \$ / + / *] \})$$

$$= \{ [E \rightarrow T \cdot, \$ / +], [T \rightarrow T \cdot * F, \$ / + / *] \}$$

$$I_4 = \text{Goto}(I_0, () = \text{Closure}(\{ [F \rightarrow (\cdot E), \$ / + / *] \})$$

$$= \{ [F \rightarrow (\cdot E), \$ / + / *],$$

$$[E \rightarrow \cdot E + T, ) / +], [E \rightarrow \cdot T, ) / +],$$

$$[T \rightarrow \cdot T * F, ) / + / *], [T \rightarrow \cdot F, ) / + / *],$$

$$[F \rightarrow \cdot ( E ), ) / + / *], [F \rightarrow \cdot i, ) / + / *] \}$$

$$I_2' = \text{Goto}(I_4, T) = \text{Closure}(\{ [E \rightarrow T \cdot, ) / +], [T \rightarrow T \cdot * F, ) / + / *] \})$$

$$= \{ [E \rightarrow T \cdot, ) / +], [T \rightarrow T \cdot * F, ) / + / *] \}$$

$$\neq I_2$$

SLR(1) のときは  $\text{Goto}(I_4, T) = I_2$  だった .

$\rightarrow$  LR(1) は , より詳細な構文解析が可能 .

SLR(1) 文法はすべて LR(1) 文法 .

LR(1) 文法で , SLR(1) 文法でないものが存在する .

## LALR(1) 構文解析

LR(1) 構文解析法は，状態数が多過ぎて実用的でない．核が同じである LR(1) 項を同一視して状態数を減らす．

例：

$$I_2 = \{ [E \rightarrow T \cdot, \$ / +], [T \rightarrow T \cdot * F, \$ / + / *] \}$$

$$I'_2 = \{ [E \rightarrow T \cdot, ) / +], [T \rightarrow T \cdot * F, ) / + / *] \}$$

を融合し，

$$I_2 = \{ [E \rightarrow T \cdot, \$ / + / ), [T \rightarrow T \cdot * F, \$ / + / * / )] \}$$

さらに

$$\text{Goto}(I_4, T) = I_2$$

LALR(1) は，正確さは LR(1) より劣るが，

SLR(1) よりは正確．

状態数は SLR(1) と同じ．

LALR(1)文法で，SLR(1)文法でない例：

$$P4 = \{ S \rightarrow E + E \mid i \\ E \rightarrow i \}$$

$$I_0 = \{ [S' \rightarrow \cdot S, \$], [S \rightarrow \cdot E + E, \$], [S \rightarrow \cdot i, \$], [E \rightarrow \cdot i, +] \}$$

$$I_1 = \{ [S' \rightarrow S \cdot, \$] \}$$

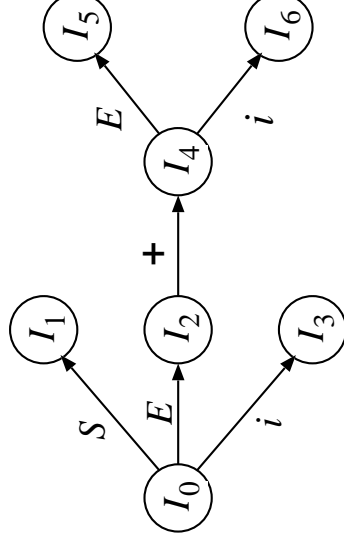
$$I_2 = \{ [S \rightarrow E \cdot + E, \$] \}$$

$$I_3 = \{ [S \rightarrow i \cdot, \$], [E \rightarrow i \cdot, +] \}$$

$$I_4 = \{ [S \rightarrow E \cdot E, \$], [E \rightarrow \cdot i, \$] \}$$

$$I_5 = \{ [S \rightarrow E + E \cdot, \$] \}$$

$$I_6 = \{ [E \rightarrow i \cdot, \$] \}$$



この文法はLALR(1)文法（LR(1)文法でもある）

SLR(1)用の  $I_3 = \{ S \rightarrow i \cdot, E \rightarrow i \cdot \}$  について

$$\text{Follow}(S) = \{ \$ \}, \text{Follow}(E) = \{ \$, + \}$$

なので，SLR(1)文法でない。

ポイント：

$I_3$ は， $I_0$ から  $i$ を読み込んだ直後の状態。

$E \rightarrow i$ で還元すれば， $S \rightarrow E + E$ の最初の  $E$ 。

だから，次の入力は ‘+’

二つ目の  $E$ のために， $\text{Follow}(E)$ が  $\$$ を含む。

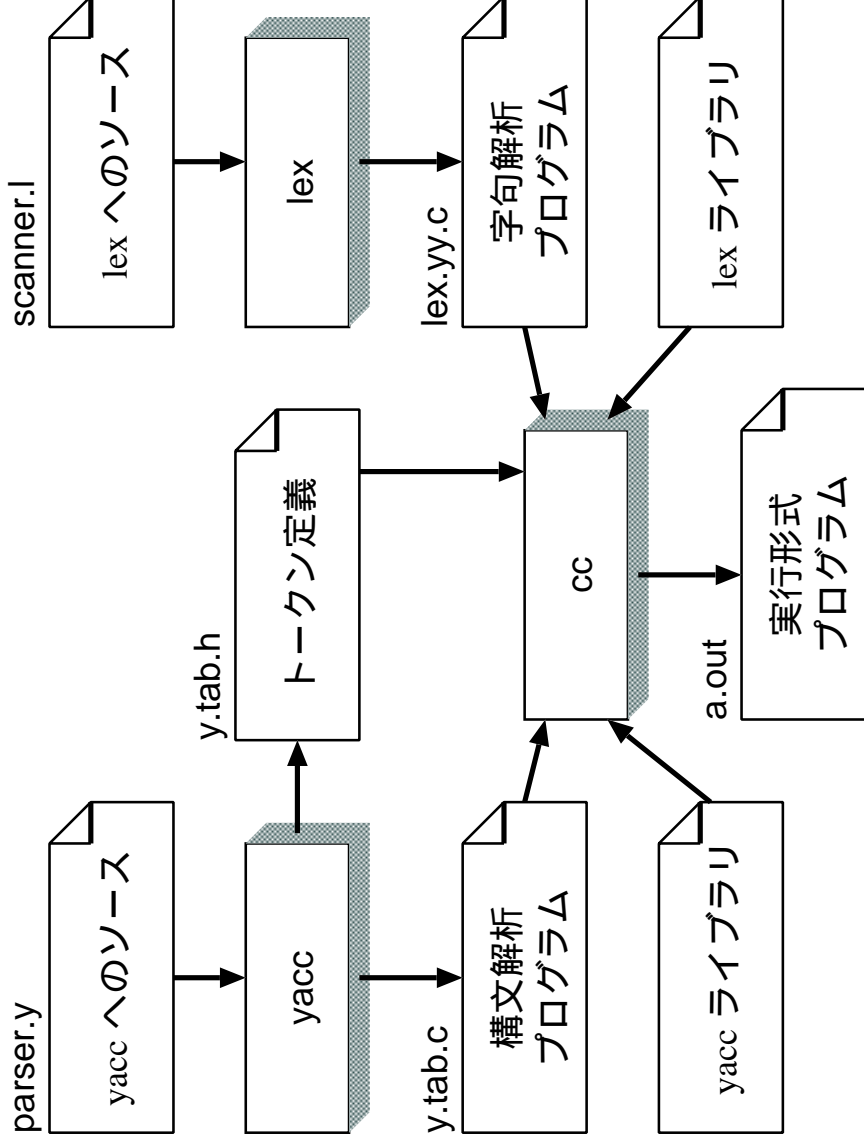


## 構文解析の自動化 (yacc – Yet Another Compiler Compiler)

```
% cat parser.y
%token PLUS TIMES LPAR RPAR IDENTIFIER
%%
E: E PLUS T {$$=make_tuple("+", $1, $3);}
  | T       {$$=$1;}
T: T TIMES F {$$=make_tuple("*", $1, $3);}
  | F       {$$=$1;}
F: LPAR E RPAR {$$=$2;}
  | IDENTIFIER {$$=$1;}
%%
...   ここにmake_tuple等の定義 ...
% yacc -d parser.y
% ls -l
...   565 Aug  2 23:36 parser.y
...  11838 Aug  2 23:36 y.tab.c
...    92 Aug  2 23:36 y.tab.h
% cat y.tab.h
#define PLUS 257
#define TIMES 258
...
%
```

yacc用にlexを使う例：

```
% cat scanner.l
%{
#include "y.tab.h"
extern int yylval;
}%
%%
[\\t ]+ {}
"+" {return PLUS;}
"*" {return TIMES;}
"(" {return LPAR;}
")" {return RPAR;}
[a-z][a-z0-9]* {yylval = copy_lexeme(); return IDENTIFIER;}
\\n {return 0;}
. {printf("lexical error: '%c'\\n", yytext[0]); return 0;}
}%
char StringTable[1000]; char *STp = StringTable;
int copy_lexeme() {
    int p = (int) STp; strcpy(STp, yytext);
    STp += strlen(yytext)+1; return p;
}
%
```



make\_tupleに構文木を表示させる：

```

int prev_tuple = 0;
int make_tuple(char *op, int arg1, int arg2) {
    printf("T%d: (%s, ", ++prev_tuple, op);
    if (arg1 > 0) printf("%s", arg1); else printf("T%d, ", -arg1);
    if (arg2 > 0) printf("%s", arg2); else printf("T%d", -arg2);
    printf("\n"); return -prev_tuple;
}
  
```

実行例：

```
% lex scanner.l
% cc y.tab.c lex.yy.c -ly -ll
% a.out
x*y+z*w
T1: (*, x, y)
T2: (*, z, w)
T3: (+, T1, T2)
% a.out
x+y+
T1: (+, x, y)
syntax error
%
```

yaccを使って，動作表と行先表の情報を調べる：

```
% yacc -v parser.y
% cat y.output
state 0
    $accept : . E $end (0)
    IDENTIFIER shift 5
    LPAR shift 4
    . error
    E goto 1
    T goto 2
    F goto 3
state 1
    ... 中略 ...
state 11
    F : LPAR E RPAR . (5)
    . reduce 5
7 terminals, 4 nonterminals
7 grammar rules, 12 states
%
```

## あいまいな文法への対処

文法 Gif =  $\langle \text{Pif}, S \rangle$  を使用する .

$$\text{Pif} = \{ S \rightarrow \text{if } ( E ) S \mid \text{if } ( E ) S \text{ else } S \mid i = E ; \\ E \rightarrow i \}$$

## LL(1)構文解析の場合

くくり出しの結果  $\text{Gif}' = \langle \text{Pif}', S \rangle$

$$\text{Pif}' = \{ S \rightarrow \text{if } ( E ) S A \mid i = E ;$$
$$A \rightarrow \text{else } S \mid \varepsilon$$
$$E \rightarrow i \}$$
$$\text{First}(E) = \{i\}$$
$$\text{First}(A) = \{\text{else}\} \cup \text{Follow}(A)$$
$$= \{\text{else}\} \cup \text{Follow}(S) = \{\text{else}, \$\}$$
$$\text{First}(S) = \{\text{if}, i\}$$
$$\text{Follow}(S) = \{\$, \} \cup \text{First}(A) = \{\$, \text{else}\}$$
$$\text{Follow}(A) = \text{Follow}(S) = \{\$, \text{else}\}$$
$$\text{Follow}(E) = \{ \}, ; \}$$
$$\text{Director}(S, \text{if } ( E ) S A) = \{\text{if}\}$$
$$\text{Director}(S, i = E ; ) = \{i\}$$
$$\text{Director}(A, \text{else } S) = \{\text{else}\}$$
$$\text{Director}(A, \varepsilon) = \text{Follow}(A) = \{\$, \text{else}\}$$
$$\text{Director}(E, i) = \{i\}$$

$\text{Director}(A, \text{else } S) \cap \text{Director}(A, \varepsilon) = \{\text{else}\} \neq \phi$

なので、 $\text{Gif}'$  は、LL(1)文法でない。

次が else なら,  $A \rightarrow \text{else } S$  を採用することにする.

```
tuple parse_A() {
    if (is_else(*p)) {
        p++;
        return parse_S();
    } else if (is_eof(*p))
        return EMPTY;
    else
        syntax_error();
}

tuple parse_S() {
    tuple x, y, z;
    if (is_if(*p)) {
        p++;
        if (!is_lpar(*p++)) syntax_error();
        x = parse_E();
        if (!is_rpar(*p++)) syntax_error();
        y = parse_S(); z = parse_A();
        return make_tuple("if", x, y, z);
    } else if (is_i(*p)) {
        x = *p++;
        if (!is_eq(*p++)) syntax_error();
        y = parse_E();
        if (!is_semicolon(*p++)) syntax_error();
        return make_tuple("=", x, y);
    } else
        syntax_error();
}
}
```



# LR構文解析の場合

Gifの正準集合(抜粋):

$$\begin{aligned} I_0 &= \text{Closure}(\{S' \rightarrow \cdot S\}) \\ &= \{S' \rightarrow \cdot S \$, S \rightarrow \cdot \text{if}(E) S, S \rightarrow \cdot \text{if}(E) S \text{ else } S, S \rightarrow \cdot i = E; \} \\ I_1 &= \text{Goto}(I_0, \text{if}) = \{S \rightarrow \text{if} \cdot (E) S, S \rightarrow \text{if} \cdot (E) S \text{ else } S\} \\ I_2 &= \text{Goto}(I_1, () ) = \text{Closure}(\{S \rightarrow \text{if}(\cdot E) S, S \rightarrow \text{if}(\cdot E) S \text{ else } S\}) \\ &= \{S \rightarrow \text{if}(\cdot E) S, S \rightarrow \text{if}(\cdot E) S \text{ else } S, E \rightarrow \cdot i\} \\ I_3 &= \text{Goto}(I_2, E) = \{S \rightarrow \text{if}(E \cdot) S, S \rightarrow \text{if}(E \cdot) S \text{ else } S\} \\ I_4 &= \text{Goto}(I_3, () ) = \text{Closure}(\{S \rightarrow \text{if}(E) \cdot S, S \rightarrow \text{if}(E) \cdot S \text{ else } S\}) \\ &= \{S \rightarrow \text{if}(E) \cdot S, S \rightarrow \text{if}(E) \cdot S \text{ else } S, \\ &\quad S \rightarrow \cdot \text{if}(E) S, S \rightarrow \cdot \text{if}(E) S \text{ else } S, S \rightarrow \cdot i = E; \} \\ I_5 &= \text{Goto}(I_4, S) = \{S \rightarrow \text{if}(E) S \cdot, S \rightarrow \text{if}(E) S \cdot \text{else } S\} \end{aligned}$$

不都合な状態  $I_5$  について

$$\text{Follow}(S) = \{\$, \text{else}\}$$

なので, GifはSLR(1)文法ではない.

シフトすると決めてしまえばSLR(1)構文解析は可能.

Action( $I_5, \text{else}$ ) = s6として動作表を作る.

GifはLR(1)文法でもLALR(1)文法でもない.

同様に対処すれば, LR(1)・LALR(1)構文解析も可能.

## yaccの場合

### yaccのルール

1. 還元よりもシフトを優先する .
2. 先に与えた生成規則による還元を優先する .

```
% cat Gif.y
... 中略 ...

S:  IF LPAR E RPAR S {$$=make_tuple("if", $3, $5, EMPTY);}
    | IF LPAR E RPAR S ELSE S {$$=make_tuple("if", $3, $5, $7);}
    | IDENTIFIER EQ E SEMICOLON {$$=make_tuple("=", $1, $3);}

E:  IDENTIFIER {$$=$1;}
... 中略 ...

% yacc -d Gif.y
conflicts: 1 shift/reduce
% cc y.tab.c lex.yy.c -ly -ll
% a.out
if (x) if (y) a=b; else b=a;
T1: (=, a, b)
T2: (=, b, a)
T3: (if, y, T1, T2)
T4: (if, x, T3, -)
%
```

## エラーリカバリ

1回のコンパイルで、できるだけ多くのエラーを検出  
エラー検出で構文解析をただちに終了しないで、  
なんらかの処置を行って実行を継続。

### LL(1)構文解析のエラーリカバリ

parse\_Aがエラーを検出するのは：

1. 次の入力 $\alpha$ が、どのDirector( $A, \alpha$ )にも属さない。
2.  $A \rightarrow \alpha$ の解析中に、次の入力 $\alpha$ が、期待した終端記号と異なる。

例：Gif'の parse\_S

Aから生成されたであろうトークン列を予測し、それらを読み飛ばす。

1の場合：

$A \xrightarrow{*} a$ である終端記号 $a$ まで

または、Follow( $A$ )の要素の直前まで読み飛ばす。

2の場合：

期待する終端記号の自動挿入（読み飛ばしなし）

または、期待する終端記号まで読み飛ばす。

# LR構文解析のエラーリカバリ

動作表の値が「未定義」のときにエラー検出

非終端記号ごとのエラーリカバリは困難

→ 大きな単位(文, 宣言, 関数定義など)について

非終端記号  $A$  についてのエラーリカバリ:

$$(q_0 a_1 q_1 \cdots a_n q_n, x_k x_{k+1} \cdots x_m \$)$$

の配置で,  $Action(q_n, x_k)$  が未定義

エラーがなく  $A$  への還元が成功していれば,

$$(q_0 a_1 q_1 \cdots q_i A q, x_j x_{j+1} \cdots x_m \$) \quad 0 \leq i < n, k \leq j \leq m, q = Goto(q_i, A)$$

スタックの巻き戻し:  $Goto(q_i, A)$  が定義されている  $q_i$  が見つかるまで

入力トークンの読み飛ばし:  $Follow(A)$  の要素  $x_j$  が見つかるまで

Gifの  $S$  についてのエラーリカバリ:

```
while ((state = lookup_goto(top_state(), S)) < 0) {
    pop_state(); pop_value();
}
pop_value(); push_value(ERROR);
while (!is_eof(*p) && !is_else(*p)) p++;
```

## yaccのエラーリカバリ: エラー規則を使用

```
S: IF LPAR E RPAR S { ... }
  | IF LPAR E RPAR S ELSE S { ... }
  | IDENTIFIER EQ E SEMICOLON { ... }
  | error EQ E SEMICOLON {
    printf("Error: bad lvar.\n"); $$=$$3;}
  | IDENTIFIER error E SEMICOLON {
    printf("Error: '=' expected.\n");
    $$=make_tuple("=", $1, $3);};
```

「 $S \rightarrow i = E ;$ 」に対応するエラー規則

```
S → error = E ;
S → i error E ;
```

```
% a.out
if (x) =b; else b=a;
Error: bad lvar.
T1: (=, b, a)
T2: (if, x, b, T1)
% a.out
if (x) a b; else b=a;
Error: '=' expected.
T1: (=, a, b)
T2: (=, b, a)
T3: (if, x, T1, T2)
%
```

「 $A \rightarrow \alpha_1 \text{ error } \alpha_2$ 」の error の位置で構文エラー：

「 $A \rightarrow \alpha_1 \cdot \text{error } \alpha_2$ 」を含む状態で検出

1. エラートークンを生成して、

$$A \rightarrow \alpha_1 \text{ error } \cdot \alpha_2$$

にシフト。

2.  $\alpha_2$  に還元できるトークン列を探し、 $A \leftarrow$ 還元。  
(トークンを読み飛ばしながら、通常の解析)
3. 通常の構文解析を再開。

## 第5章

# 意味解析

# Cコンパイラの主要な意味解析処理

1. 名前とオブジェクトとの対応づけ
2. 型チェック

オブジェクト:名前をつけられる構成要素  
変数, 関数, ラベル, 構造体, 共用体,  
構造体/共用体のメンバ, enum 定数, 型

例:

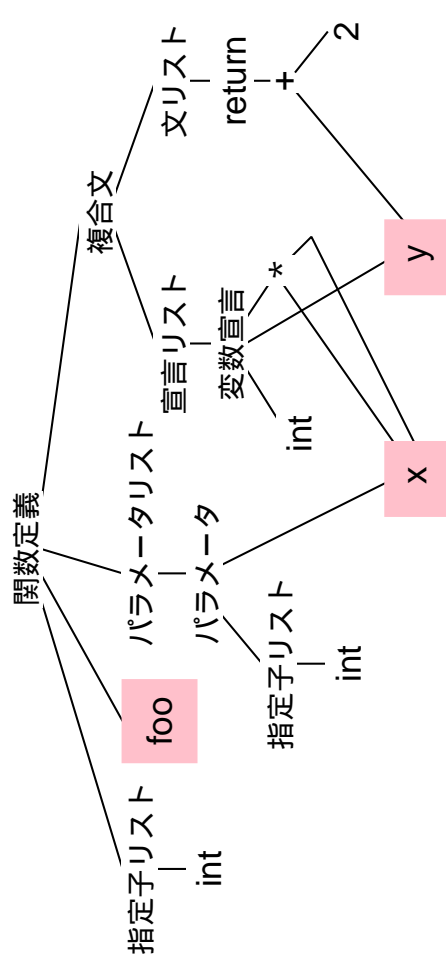
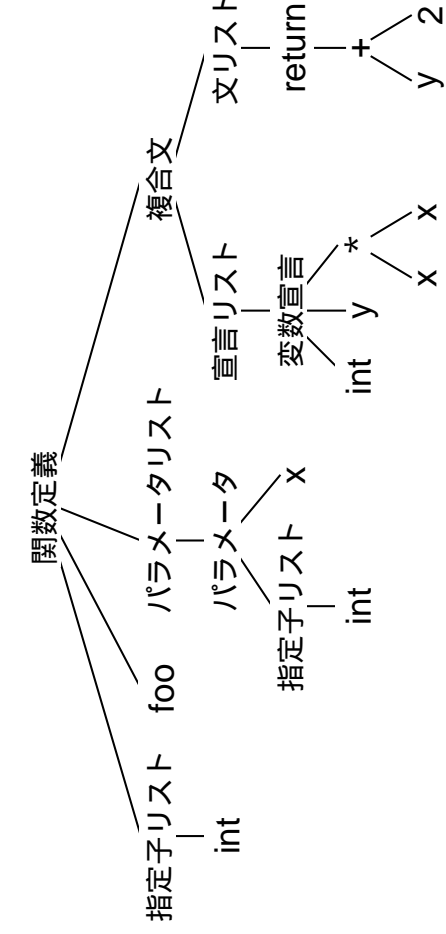
```
int foo(int x) {  
    int y = x*x;  
    return y+2;  
}  
  
double bar(int x) {  
    double y = x*x;  
    return y+2;  
}
```



# オブジェクト構造体

オブジェクトを表現する構造体

- オブジェクトの種類 (変数, 関数, ラベルなど)
- オブジェクトの名前
- 型情報
- 属性 (auto, register, static など)
- 最適化のための情報
- オブジェクトを割り当てる場所



# 名前空間とスコープ

名前空間：名前が表すオブジェクトの検索空間

C言語の主要な名前空間：

- ・変数，関数，typedef名，enum定数の名前空間
- ・ラベルの名前空間
- ・構造体と共用体の名前空間

例：

```
struct tarao {  
    int tarao;  
} tarao;
```

```
struct tarao x = tarao;
```

スコープ：オブジェクトを参照できる範囲

```
int foo(int x) {  
    int y = x*x;  
    return y+2;  
}
```

```
int foo(int x) {  
    int y = x*x;  
    return y+2;  
}
```

## スコープ間の関係

- ・ 二つのスコープには共通部分がない．
- ・ 二つのスコープはまったく一致する．
- ・ 片方のスコープが，もう一方に含まれる．

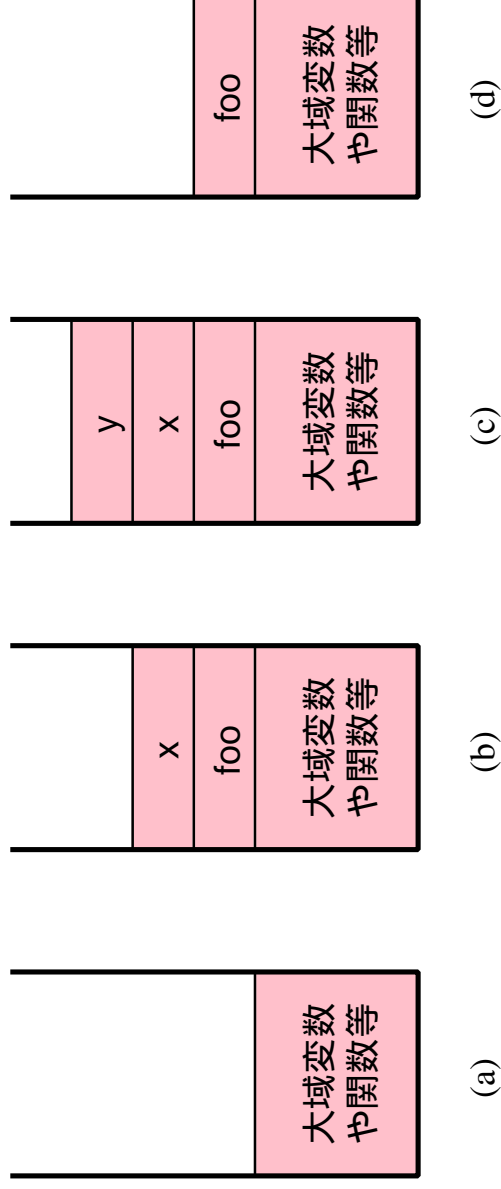
同じ名前空間に属する二つのオブジェクトのスコープが一致するときは，異なる名前をつけなければならない．

例：C言語の関数パラメータ

## 名前とオブジェクトの対応づけ

1. まず、 $x$ の現れているプログラム中の位置によって、 $x$ が属する名前空間を特定する。
2. その名前空間に登録されているオブジェクトのうち、 $x$ の現れている位置をスコープに含むものを、内側のスコープを持つものから順に調べる。
3.  $x$ を名前とするオブジェクトが見つかれば、最初に見つかったオブジェクトを、 $x$ は表す。そのようなオブジェクトがなければ、 $x$ の表すオブジェクトは未定義である。

## スタックを使った実装



大域変数、関数、typedef名などは、ハッシュ表等に登録するほうが検索効率が良い。

## 前方参照

```
int sigma(int n) {
    int i = 0, s = 0;
    goto test;
loop:
    s += i;
test:
    if (++i <= n) goto loop;
    return s;
}
```

1. “goto test;” で、未定義ラベル test をラベル表に登録。  
{<test, 未定義, 参照済>}
2. “loop:” で、ラベル loop を登録。  
{<test, 未定義, 参照済>, <loop, 定義済, 未参照>}
3. “test:” で、ラベル test が「定義済」に。  
{<test, 定義済, 参照済>, <loop, 定義済, 未参照>}
4. “goto loop;” で、ラベル loop が「参照済」に。  
{<test, 定義済, 参照済>, <loop, 定義済, 参照済>}
5. 関数の意味解析終了時に、ラベル表を調べ、  
「未定義」のラベルがあればエラー。「未参照」のラベルがあれば警告。

## 型子エックと型変換

型表現：型を表すデータ

```
<char, (signed|unsigned|-)>  
<int, (signed|unsigned|-, (long|short|-))>  
<float>  
<double, (long|-)>  
<enum, ...>
```

多くのプログラミング言語（C言語を含む）では，構文木を葉から根の方向に向かって調べてれば，式の型が決定できる．

例： $e_1 * e_2$ の型子エック  
   $A : e_1$ の型（型表現）  
   $B : e_2$ の型（型表現）

```
if (A==Tdouble) {  
    check_double(B); return A;  
} else if (B==Tdouble) {  
    check_double(A); return B;  
} else if (A==Tfloat) {  
    check_float(B); return A;  
} else if (B==Tfloat) {  
    check_float(A); return B;  
} else {  
    check_int(A); check_int(B); return  
Tint; }  
126
```

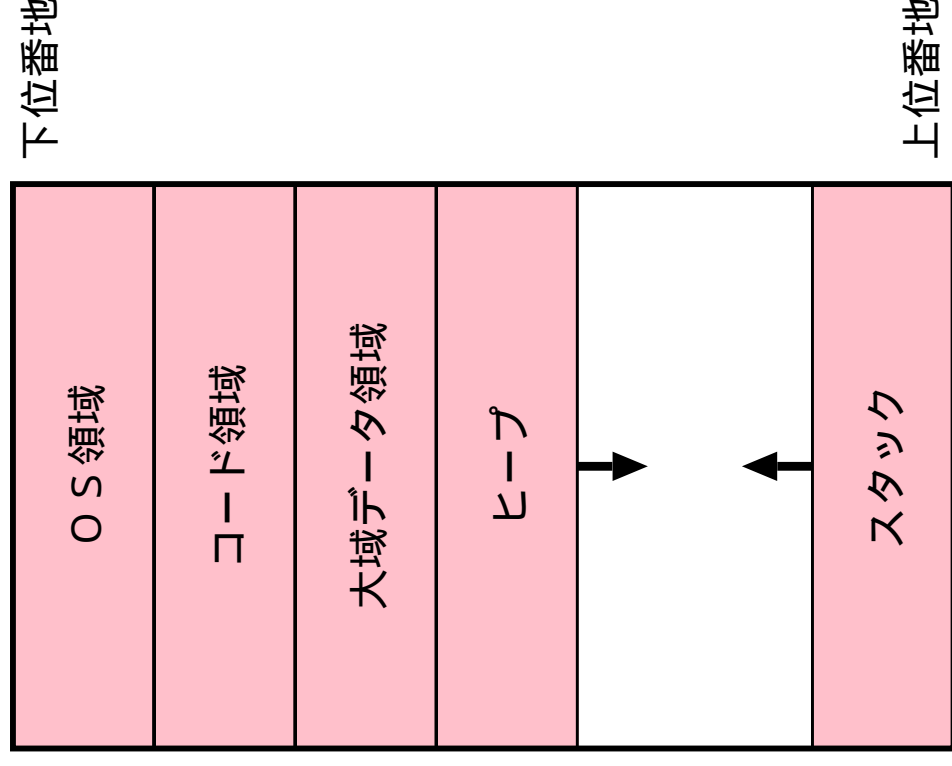
## 第6章

# コード生成

# 実行環境のモデル

## メモリ領域

- OS領域（場所・サイズは固定）
- ユーザ領域
  - － コード領域（ロード時に確定）
  - － データ領域
    - \* 大域データ領域（ロード時に確定）
    - \* スタック（実行時に拡大・縮小）
    - \* ヒープ（実行時に拡大）





# CPUとレジスタ

中央処理装置 (Central Processing Unit, CPU):

コード領域から機械語命令を一つずつ取り出し実行

算術論理装置 (Arithmetic Logic Unit, ALU):

算術演算や比較演算を実行

## Pentium のレジスタ

1. 汎用レジスタ (general-purpose register):

オペランドとして使用できる32ビット長レジスタ

ベースポインタ `ebp` , スタックポインタ `esp` , `eax` , `ebx` , `ecx` , `edx` など

2. 条件フラグ (condition flag):

比較命令の結果を格納する1ビット長レジスタ

ゼロフラグ `zf` , 符号フラグ `sf` など

3. 命令ポインタ (プログラムカウンタ, instruction pointer):

次の命令の場所を指す32ビット長レジスタ `eip`

CPUは次の動作を繰り返す:

1. `eip`の指す位置から命令を一つ取り出す
2. 次の命令を指すように`eip`を更新する
3. 取り出した命令を実行する

## アセンブリ命令

### ラベル宣言

〈ラベル〉：

### アセンブリ命令

〈命令名〉 〈オペランド<sub>1</sub>〉 , … , 〈オペランド<sub>*n*</sub>〉

まとめて

〈ラベル〉： 〈命令名〉 〈オペランド<sub>1</sub>〉 , …

オペランド：

- 汎用レジスタ
- メモリ番地
  - ラベル
  - 相対番地:  $n[R]$
- 整数定数

## 算術命令:

```
add  eax,ebx    ; eaxにebxの値を足す
sub  esp,4      ; espから4を引く
imul ebx,_x     ; ebxにxの値を掛ける
```

## 移動命令:

```
mov  eax,4      ; eaxに整数4を格納する
mov  ebp,esp    ; espの値をebpに格納する
```

## 無条件ジャンプ命令:

```
jmp  ラベル
```

## 比較命令:

```
cmp  x,y
```

## 条件付きジャンプ命令:

```
j... ラベル
```

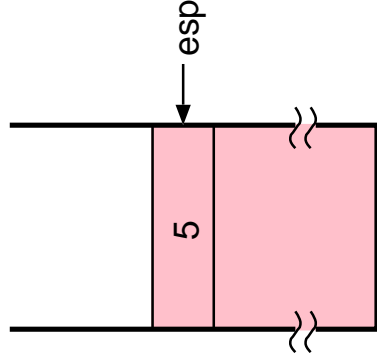
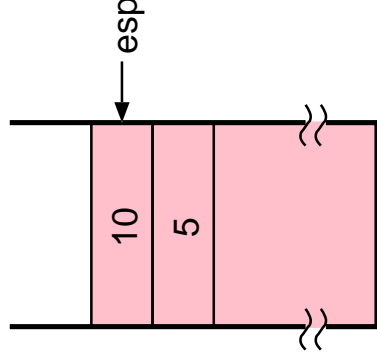
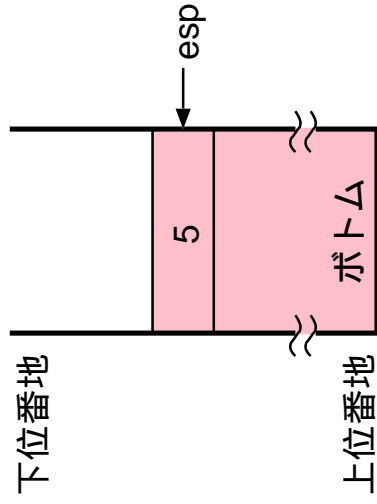
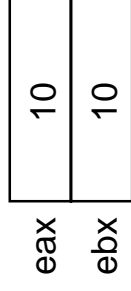
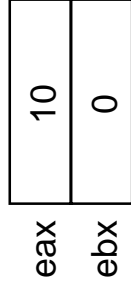
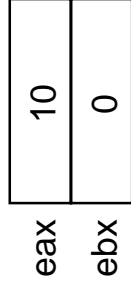
---

### 命令 条件 意味

jg	$x > y$	Jump if Greater
jge	$x \geq y$	Jump if Greater or Equal
je	$x = y$	Jump if Equal
jne	$x \neq y$	Jump if Not Equal
j1	$x < y$	Jump if Less
jle	$x \leq y$	Jump if Less or Equal

# スタック操作命令:

```
push eax  
pop ebx
```



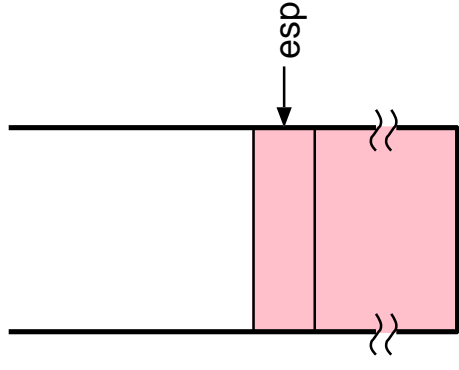
(a) 実行前

(b) push eax 実行後

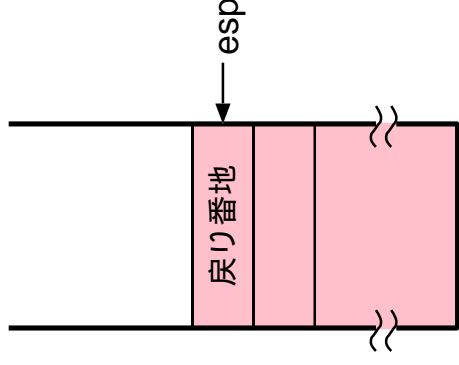
(c) pop ebx 実行後

関数呼出し命令: call ラベル

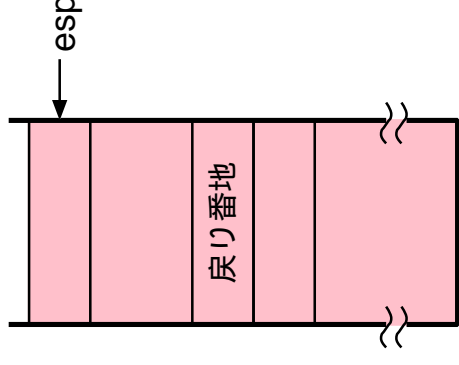
リターン命令: ret



(a) 呼出し直前  
(e) リターン直後



(b) 呼出し直後  
(d) リターン直前

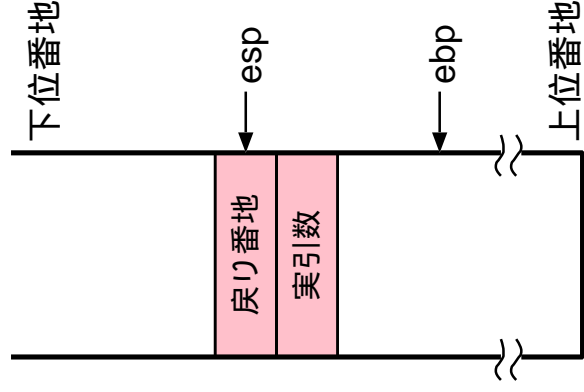


(c) 関数実行中

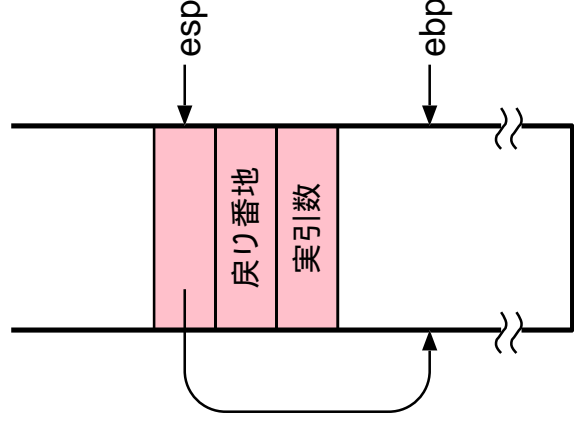
# 関数呼出し

```
int foo(int x) {  
    int y = x*x;  
    return y+2;  
}
```

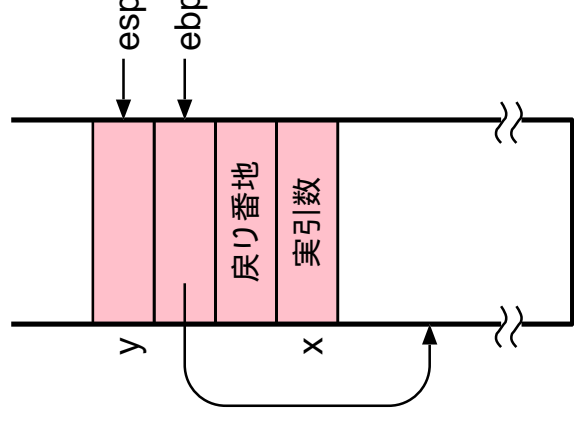
```
1  _foo:  push  ebp  
2      mov  ebp, esp  
3      sub  esp, 4  
4      ...  
9      mov  esp, ebp  
10     pop  ebp  
11     ret
```



(a) 呼出し直後



(b) ebp の値の保存



(c) 本体実行中

# fooの呼出し

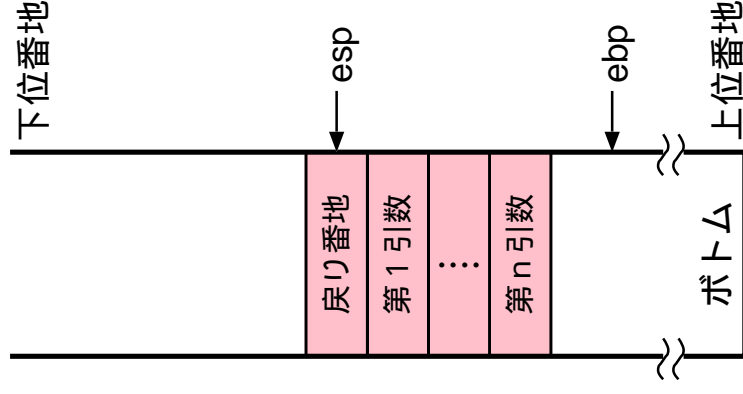
```
push 5  
call _foo  
add esp, 4
```

# 一般の関数 f のコード

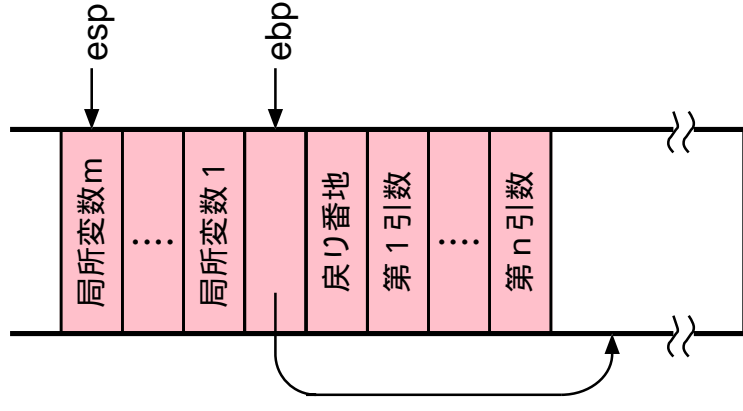
```
-f: push ebp  
mov ebp, esp  
sub esp, Nlocal
```

本体の実行

```
Lret: mov esp, ebp  
pop ebp  
ret
```



(a) 呼出し直後



(b) 本体実行中

## 関数呼出し式 $f(e_1, e_2, \dots, e_n)$ のコード

$e_n$  を計算し, 結果をプッシュする

...

$e_2$  を計算し, 結果をプッシュする

$e_1$  を計算し, 結果をプッシュする

call  $-f$

add esp,  $n \times 4$

例:  $f(1, g(2, 3), 4)$  のコード

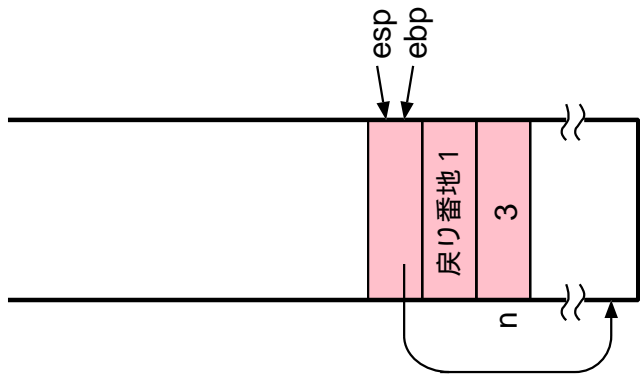
```
push 4      ;  $f$  への第3引数
push 3      ;  $g$  への第2引数
push 2      ;  $g$  への第1引数
call  $-g$ 
add esp, 8
push eax    ;  $f$  への第2引数
push 1      ;  $f$  への第1引数
call  $-f$ 
add esp, 12
```



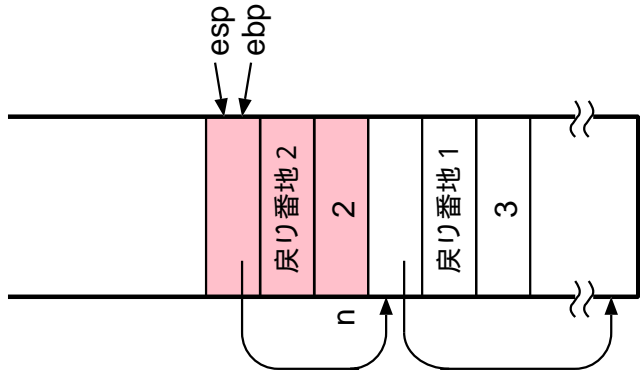
例: 階乗を求める関数 fact のコードとその実行

```
int fact(int n) {  
    if (n == 1) return 1;  
    else return n * fact(n-1);  
}
```

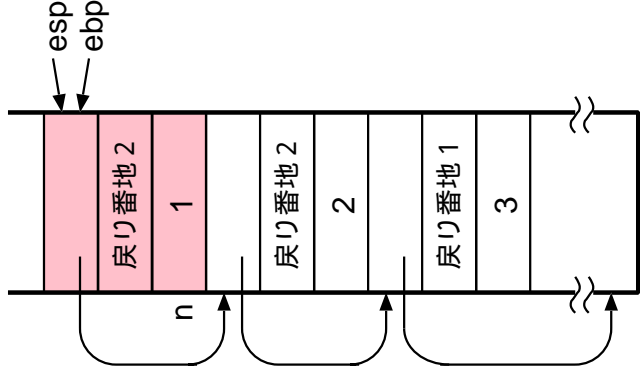
```
_fact: push ebp  
       mov  ebp,esp  
       cmp  8[ebp],1      ; n = 1かどうか  
       jne  L2           ; n ≠ 1なら L2へ  
       mov  eax,1        ; n = 1なら戻り値は1  
       jmp  L1  
L2:   mov  eax,8[ebp]    ; n - 1の計算  
       sub  eax,1  
       push eax  
       call _fact       ; 再帰呼出しで (n - 1)! の計算  
       add  esp,4  
       imul eax,8[ebp]  ; n × (n - 1)! の計算  
L1:   pop  ebp  
       ret
```



(a) fact(3) 実行中



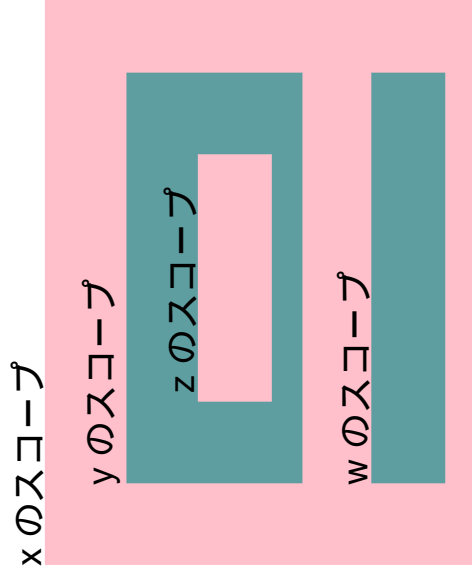
(b) fact(2) 実行中



(c) fact(1) 実行中

# 局所変数等の割当て

1. 局所変数領域のサイズ  $N_{local}$  をできるだけ小さく
  - 同じ位置に複数の局所変数を割り当てる
2. スコープが重なる局所変数は、同じ位置に割り当てない



順序	動作	相対番地	last_alloc	max_alloc
0	初期化		0	0
1	xの割当て	-4	1	1
2	yの割当て	-8	2	2
3	zの割当て	-12	3	3
4	zの解放		2	3
5	yの解放		1	3
6	wの割当て	-8	2	3
7	wの解放		1	3
8	xの解放		0	3

$$N_{local} = 3 \times 4 = 12 \text{ バイト使用}$$

# レジスタの退避

Pentiumの汎用レジスタは、  
espとebpの他に6個だけ

1. 呼出し後保存レジスタ (callee-saved register)  
呼び出す側が使用中かもしれない汎用レジスタ。  
呼び出される関数は、使用前に退避
2. 呼出し前保存レジスタ (caller-saved register)  
上記以外の汎用レジスタ。  
呼び出される関数は、退避せずに使用してよい  
別の関数を呼び出すときは、退避しておく

少なくともeaxレジスタは呼出し前保存  
→ 全レジスタが呼出し前保存と仮定

# 文のコード生成

## 複合文のコード

$\{d_1 \dots d_n \ s_1 \dots s_m\}$

$d_1$  の初期化

...

$d_n$  の初期化

$s_1$  の実行

...

$s_m$  の実行

宣言 `int x = 10;` の初期化

`mov n[ebp], 10`

# if文のコード

if (e)  $s_1$  else  $s_2$

$e$ を計算し，結果が偽ならば  $L_1$ へジャンプ

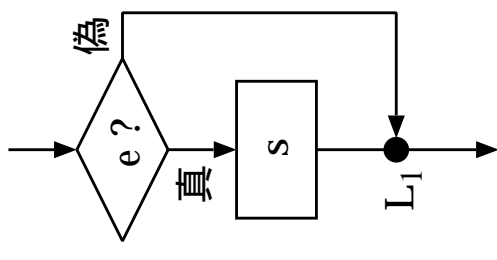
$s_1$ の実行

jmp  $L_2$

$s_2$ の実行

$L_1$ :

$L_2$ :

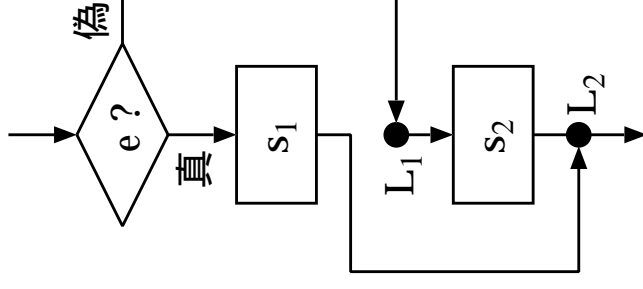


if (e)  $s$

$e$ を計算し，結果が偽ならば  $L_1$ へジャンプ

$s$ の実行

$L_1$ :



## while 文のコード

while (e) s

$L_1$ :  $e$ を計算し, 結果が偽ならば  $L_2$ へジャンプ

$s$ の実行

jmp  $L_1$

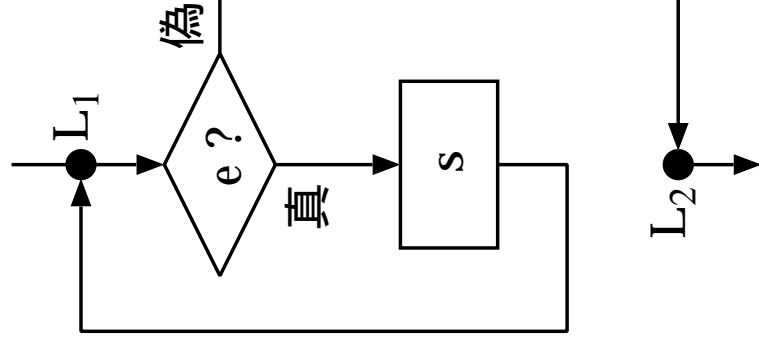
$L_2$ :

## break 文のコード

jmp  $L_2$

## continue 文のコード

jmp  $L_1$



# 式文と代入式のコード

式文  $e$ ;

$e$  の計算 (結果は破棄)

変数  $v$  への代入式  $v = e'$

$e'$  の計算 ( $R$  に結果)

mov loc( $v$ ),  $R$

例:  $x = y = z + 1$ ;

mov  $R$ , loc( $z$ )

add  $R$ , 1

mov loc( $y$ ),  $R$

mov loc( $x$ ),  $R$



# ラベルとジャンプ命令の抑制

```
if (e1) {  
    if (e2) s1 else s2  
}
```

$e_1$  を計算し，結果が偽なら  $L_1$  へジャンプ

$e_2$  を計算し，結果が偽なら  $L_2$  へジャンプ

$s_1$  の実行

jmp  $L_3$

$L_2$ :  $s_2$  の実行

$L_3$ :

$L_1$ :

$L_1$  と  $L_3$  は融合可能

## 抑制の方法

1. まったく参照されないラベルを抑制する
2. 直後に来るラベル（もしあれば）を再利用する
3. 直後が無条件ジャンプなら，そのラベルへジャンプする

# 算術式のコード生成

演算式  $e_1 \circ e_2$ :

1.  $e_1$  の値をあるレジスタ  $R$  にロードし,
2.  $R$  を第1オペランド,  $e_2$  の値を第2オペランドとして,  
‘ $\circ$ ’ に対応する命令 *inst* を実行する.

例:  $x+y$

```
mov  R, loc(x)
add  R, loc(y)
```

例:  $a*b+y$

```
mov  R, loc(a)
imul R, loc(b)
add  R, loc(y)
```

例:  $a*b+x*y$

```
mov  R1, loc(a)
imul R1, loc(b)
mov  R2, loc(x)
imul R2, loc(y)
add  R1, R2
```

例:  $R_1$  だけを使って  $a*b+x*y$  を計算

```
mov  R1, loc(a)
imul R1, loc(b)
mov  temp, R1
mov  R1, loc(x)
imul R1, loc(y)
add  R1, temp
```

## 可換演算と非可換演算

例:  $R_1$  だけを使って  $a*b-x*y$  を計算

```
mov  R1, loc(a)
imul R1, loc(b)
mov  temp, R1
mov  R1, loc(x)
imul R1, loc(y)
sub  R1, temp
```

ではなく

```
mov  R1, loc(x)
imul R1, loc(y)
mov  temp, R1
mov  R1, loc(a)
imul R1, loc(b)
sub  R1, temp
```

利用可能なレジスタ数が  $N$  ,  
 $e_1$  と  $e_2$  の計算に必要なレジスタ数も  $N$  のとき ,  
 $e_1 - e_2$  は

$e_1$  の計算 (  $R_1$  に結果 )

mov  $temp, R_1$

$e_2$  の計算 (  $R_2$  に結果 )

mov  $R_3, temp$

sub  $R_3, R_2$

よりも

$e_2$  の計算 (  $R_2$  に結果 )

mov  $temp, R_2$

$e_1$  の計算 (  $R_1$  に結果 )

sub  $R_1, temp$

右辺の計算を先に !

## 算術式のコード生成アルゴリズム

$N$ : 式の計算に利用できるレジスタ数 ( $N \geq 2$ )

$\rho(e)$ :  $e$ の計算に必要なレジスタ数

$$N \geq \rho(e) \geq 0$$

$\rho(e) = 0 \leftrightarrow e$ は変数か定数

アルゴリズム:

算術式  $e_1 \circ e_2$  に対して,

まず  $\rho(e_1)$  と  $\rho(e_2)$  を求め

各  $e_i$  のコード生成開始時点で,

$\rho(e_i)$  個以上のレジスタを空ける

## RSL (Right-Save-Left) 型

$\rho(e_1) = \rho(e_2) = N$  のとき

$e_2$  の計算 (  $R_2$  に結果 )

mov temp,  $R_2$

$e_1$  の計算 (  $R_1$  に結果 )

inst  $R_1, temp$

## RL (Right-Left) 型

$\rho(e_1) < N$  のとき

$e_2$  の計算 (  $R_2$  に結果 )

$e_1$  の計算 (  $R_1$  に結果 )

inst  $R_1, R_2$

## R (Right) 型

$e_1$  が変数/定数  $v$  で, 演算が可換のとき

$e_2$  の計算 (  $R_2$  に結果 )

inst  $R_2, loc(v)$

## LR (Left-Right) 型

$\rho(e_2) < N$  のとき

$e_1$  の計算 (  $R_1$  に結果 )

$e_2$  の計算 (  $R_2$  に結果 )

inst  $R_1, R_2$

## L (Left) 型

$e_2$  が変数/定数  $v$  のとき

$e_1$  の計算 (  $R_1$  に結果 )

inst  $R_1, loc(v)$

	$\rho(e_2) = N$	その他	$\rho(e_2) = 0$
$\rho(e_1) = N$	RSL	LR	L
$N > \rho(e_1) > 0$	RL	RL/LR	L
$\rho(e_1) = 0$ (非可換)	RL	RL	L
$\rho(e_1) = 0$ (可換)	R	R	L

## RSL型コード生成ルーチン

```
reg emit_RSL_code(char *inst, tree e1, tree e2) {
    reg R1, R2 = emit_expr(e2);
    loc temp = allocate_temp();
    emit("mov", temp, R2);
    release_register(R2);
    R1 = emit_expr(e1);
    emit(inst, R1, temp);
    release_temp(temp);
    return R1;
}
```

## RL型コード生成ルーチン

```
reg emit_RL_code(char *inst, tree e1, tree e2) {
    reg R2 = emit_expr(e2);
    reg R1 = emit_expr(e1);
    emit(inst, R1, R2);
    release_register(R2);
    return R1;
}
```

# 使用レジスタ数の計算

算術式  $e_1 \circ e_2$  の使用レジスタ数

$\rho(e_2)$ の値	$N$	その他	0
$\rho(e_1) = N$	$N$	$N$	$N$
その他	$N$	RL型: $\max(\rho(e_1) + 1, \rho(e_2))$ LR型: $\max(\rho(e_1), \rho(e_2) + 1)$	$\rho(e_1)$
$\rho(e_1) = 0$ (非可換)	$N$	$\max(2, \rho(e_2))$	1
$\rho(e_1) = 0$ (可換)	$N$	$\rho(e_2)$	1

$\rho(e_1) = N$  または  $\rho(e_2) = N$  なら  $\rho(e_1 \circ e_2) = N$

RL型なら  $\rho(e_1 \circ e_2) = \max(\rho(e_1) + 1, \rho(e_2))$

$e_2$  の計算 ( $R_2$  に結果)

$\rho(e_2)$  個使用

$e_1$  の計算 ( $R_1$  に結果)

$\rho(e_1) + 1$  個使用

*inst*  $R_1, R_2$

mov  $R, \text{loc}(x)$

imul  $R, \text{loc}(y)$

mov  $R', \text{loc}(a)$

imul  $R', \text{loc}(b)$

sub  $R', R$

例:  $a*b-x*y$

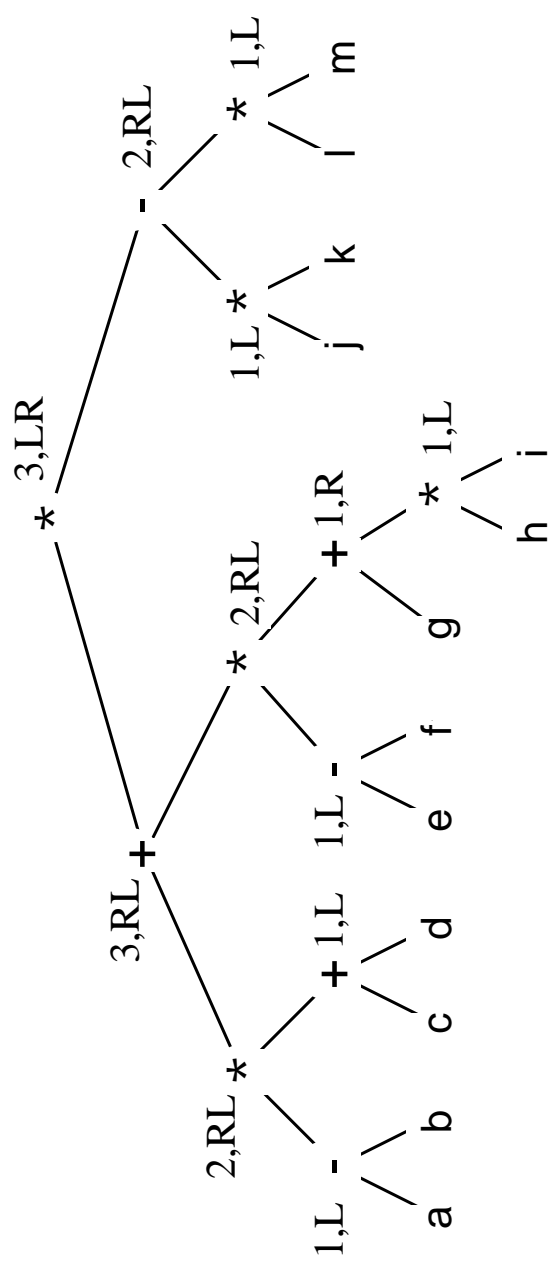
$\rho(a) = \rho(b) = 0$  より  $\rho(a*b) = 1$ , 同様に  $\rho(x*y) = 1$

したがって,  $a*b-x*y$  は RL 型



例:  $N = 3$  のとき,

$$((a-b)*(c+d)+(e-f)*(g+h*i))*(j*k-l*m)$$



```

mov  eax,h      mov  eax,c      mov  eax,l
imul eax,i      add   eax,d      imul eax,m
add   eax,g      mov  ecx,a      mov  ebx,j
mov  ebx,e      sub  ecx,b      imul ebx,k
sub  ebx,f      imul ecx,eax    sub  ebx,eax
imul ebx,eax    add  ecx,ebx    imul ecx,ebx

```

## 関数呼出し式の使用レジスタ数

$f(e_1, \dots, e_n)$

各  $e_i$  の値をスタックにプッシュ

$e_i$  の計算 ( $R$  に結果)

push  $R$

$\rho(e_i)$  個のレジスタを使用

すべての実引数をスタックにプッシュ

使用レジスタ数は  $\max(m, \rho(e_1), \dots, \rho(e_n))$

$f$  の使用レジスタ数は, 分からない.  $\rightarrow$  最大の  $N$  個と仮定

$$\rho(f(e_1, \dots, e_n)) = \max(N, \rho(e_1), \dots, \rho(e_n)) = N$$

例:  $f() - x*y$

$\rho(f()) = N$ ,  $\rho(x*y) = 1$  だから LR 型

```
call _f
mov R, loc(x)
imul R, loc(y)
sub eax, R
```

## 条件ジャンプのコード生成

「 $e$ を計算し，結果が偽なら $L$ へジャンプ」

$e$ の計算 ( $R$ に結果)

```
cmp  R,0    ; 0と比較  
je   L      ; 等しければジャンプ
```

「真なら」の場合は， $je$ を $jne$ で置き換える．

例: if (f()) x = 10; のコード

```
call _f  
cmp  eax,0  
je   L  
mov  loc(x),10
```

$L$ :

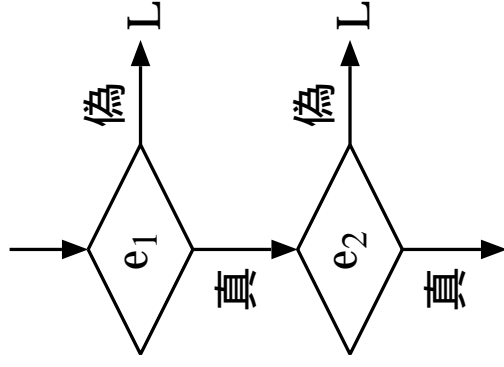
「 $e_1 >= e_2$ が真なら $L$ へ」( $e_1 >= e_2$ がRL型するとき)

$e_2$ の計算 ( $R_2$ に結果)

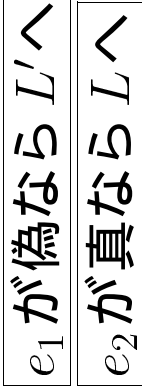
$e_1$ の計算 ( $R_1$ に結果)

```
cmp  R1,R2  
jge  L
```

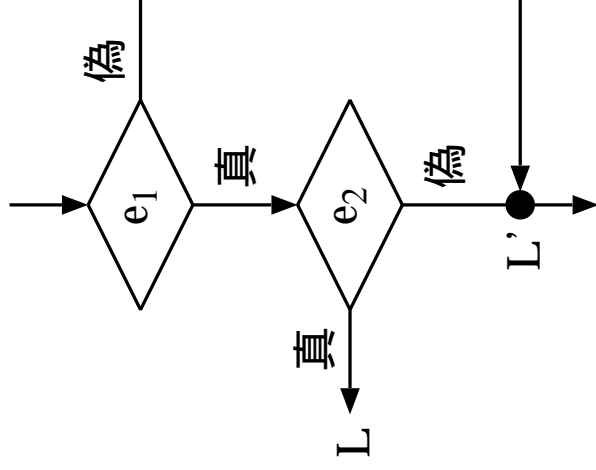
「 $e_1 \&\& e_2$ が偽なら  $L \wedge$ 」



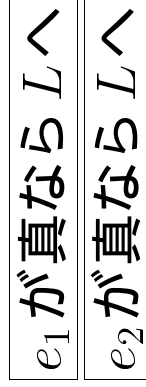
「 $e_1 \&\& e_2$ が真なら  $L \wedge$ 」



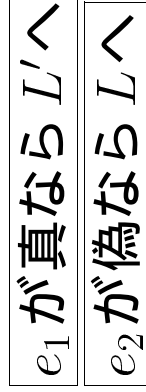
$L'$ :



「 $e_1 \vee \vee e_2$ が真なら  $L \wedge$ 」

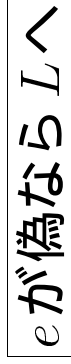


「 $e_1 \vee \vee e_2$ が偽なら  $L \wedge$ 」

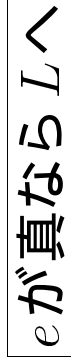


$L'$ :

「! $e$ が真なら  $L \wedge$ 」



「! $e$ が偽なら  $L \wedge$ 」



例: 「(a && b) || !(c || d) が真なら Lへ」

(a && b) || !(c || d) が真なら Lへ

a && b が真なら Lへ

a が偽なら L<sub>1</sub>へ

cmp loc(a), 0

je L<sub>1</sub>

b が真なら Lへ

cmp loc(b), 0

jne L

L<sub>1</sub>:

!(c || d) が真なら Lへ

c || d が偽なら Lへ

c が真なら L<sub>2</sub>へ

cmp loc(c), 0

jne L<sub>2</sub>

d が偽なら Lへ

cmp loc(d), 0

je L

L<sub>2</sub>:

## 戻り値の計算コード

return  $e$ ;

$e$ の計算 ( $R$ に結果)

mov eax,  $R$

jmp  $Lret$

モード指定で移動命令を削除

eaxモード: 結果をできるだけeaxに

no-eaxモード: 結果をできるだけeax以外に

freeモード: どのレジスタでもよい

例: return  $a*b-x*y$ ; のコード

$a*b-x*y$ 全体はeaxモード,  $a*b-x*y$ はRL型

$x*y$ の計算 ( $R_2$ に結果)

no-eaxモード

$a*b$ の計算 ( $R_1$ に結果)

eaxモード

*inst*  $R_1, R_2$

生成結果 (移動命令削除)

mov ebx, loc(x)

imul ebx, loc(y)

mov eax, loc(a)

imul eax, loc(b)

sub eax, ebx

jmp  $Lret$

## 関数コードの生成例

```
_fact: push ebp
      mov  ebp,esp
      本体 { if ... } の実行
      if 文 if(n==1) ... else ... の実行
      n==1が偽なら L2へ
      cmp  8[ebp],1
      jne  L2
      return 1; の実行
      mov  eax,1
      jmp  L1
L2:  return n*fact(n-1); の実行
      n*fact(n-1) の計算
      fact(n-1) の計算
      n-1を計算し, 結果をプッシュ
      mov  eax,8[ebp]
      sub  eax,1
      push eax
      call _fact
      add  esp,4
      imul eax,8[ebp]
L1:  pop  ebp
      ret
```

# その他のトピック

## 1. 局所関数

- 定義を包含する他の局所関数・最上位関数のフレーム参照
  - 静的リンクまたはディスプレイを使用
2. 3オペランド命令と1オペランド命令
    - 2オペランド命令と同様に考察
  3. 呼出し後保存レジスタ
    - 呼出し前保存レジスタを優先的に割り当てる
    - 関数呼出しを含む式は，呼出し後保存レジスタを使って高速化