# 6.110 Re-lecture 1

Regular expressions, automata, grammars, parse trees
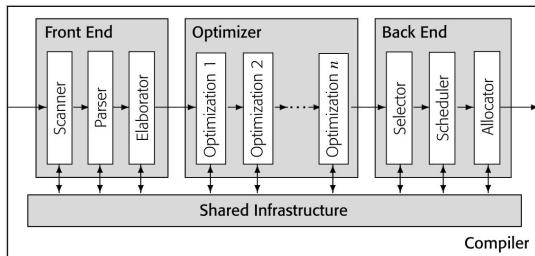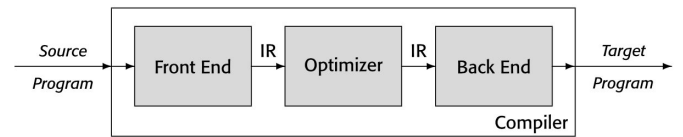


4 lectures in 1

Focus on theory

Ideal: 90 minutes

Stop me

## Plan

- Overview of compiler components and optimizations
- Regular expressions
- Finite automata
- Duality and constructions
- Tokenization
- Context-free grammar
- Ambiguity
- Precedence





# Front end

```
# Comment 1
 # Comment 2

# Factorial:

def fact(  x\
):

  if x == -1:
    return 1.j

  elif x ==0:

    return 1
  else:

      return x* fact(x

- 1)

s = "foo\
\\ \n\'\""
```

```
# Comment 1
 # Comment 2

# Factorial:

def fact(  x\
):

  if x == -1:
    return 1.j

  elif x ==0:

    return 1
  else:

      return x* fact(x

- 1)

s = "foo\
\\ \n\'\""
```

```
(KEYWORD def)          (LIT 1)
(ID "fact")            (NEWLINE)
(PUNCT "(")            (DEDENT)
(ID "x")               (KEYWORD else)
(PUNCT ")")            (PUNCT ":")
(PUNCT ":")            (NEWLINE)
(NEWLINE)              (INDENT)
(INDENT)               (KEYWORD return)
(KEYWORD if)           (ID "x")
(ID "x")               (PUNCT "*")
(PUNCT "==")           (ID "fact")
(PUNCT "-")            (PUNCT "(")
(LIT 1)                (ID "x")
(PUNCT ":")            (PUNCT "-")
(NEWLINE)              (LIT 1)
(KEYWORD return)       (PUNCT ")")
(LIT +1.i)             (NEWLINE)
(NEWLINE)              (DEDENT)
(DEDENT)               (DEDENT)
(KEYWORD elif)         (ID "s")
(ID "x")               (PUNCT "=")
(PUNCT "==")           (LIT "foo\\ \n\'\"")
(LIT 0)                (NEWLINE)
(PUNCT ":")            (ID "fact")
(NEWLINE)              (PUNCT "(")
(INDENT)               (LIT 20)
(KEYWORD return)       (PUNCT ")")
                       (NEWLINE)
                       (ENDMARKER)
```

```
# Comment 1
 # Comment 2

# Factorial:

def fact(  x\
):

  if x == -1:
    return 1.j

  elif x ==0:

    return 1
  else:

      return x* fact(x

- 1)

s = "foo\
\\ \n\'\""
```

```
Module(
  body=[
    FunctionDef(
      name='fact',
      args=arguments(
        posonlyargs=[],
        args=[
          arg(arg='x')],
        kwonlyargs=[],
        kw_defaults=[],
        defaults=[]),
      body=[
        If(
          test=Compare(
            left=Name(id='x', ctx=Load()),
            ops=[
              Eq()],
            comparators=[
              UnaryOp(
                op=USub(),
                operand=Constant(value=1))]),
          body=[
            Return(
              value=Constant(value=1j))],
          orelse=[
            _
```

# Intermediate representation

$$a \leftarrow a \times 2 \times b \times c \times d$$

$$a \leftarrow a \times 2 \times b \times c \times d$$

$$t_0 \leftarrow a \times 2$$
$$t_1 \leftarrow t_0 \times b$$
$$t_2 \leftarrow t_1 \times c$$
$$t_3 \leftarrow t_2 \times d$$
$$a \;\; \leftarrow t_3$$

$$a \leftarrow a \times 2 \times b \times c \times d$$

```
loadAI   r_arp, @a ⇒ r_a        // load 'a'
loadI    2        ⇒ r_2         // constant 2 into r_2
loadAI   r_arp,@b ⇒ r_b         // load 'b'
loadAI   r_arp,@c ⇒ r_c         // load 'c'
loadAI   r_arp,@d ⇒ r_d         // load 'd'
mult     r_a,r_2  ⇒ r_a         // r_a ← a × 2
mult     r_a,r_b  ⇒ r_a         // r_a ← (a × 2) × b
mult     r_a,r_c  ⇒ r_a         // r_a ← (a × 2 × b) × c
mult     r_a,r_d  ⇒ r_a         // r_a ← (a × 2 × b × c) × d
storeAI  r_a      ⇒ r_arp,@a    // write r_a back to 'a'
```
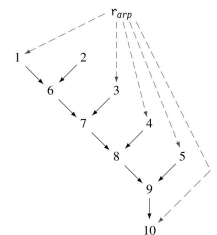
$$a \leftarrow a \times 2 \times b \times c \times d$$

```
 1   loadAI   r_arp,@a ⇒ r_a
 2   loadI    2        ⇒ r_2
 3   loadAI   r_arp,@b ⇒ r_b
 4   loadAI   r_arp,@c ⇒ r_c
 5   loadAI   r_arp,@d ⇒ r_d
 6   mult     r_a,r_2  ⇒ r_a
 7   mult     r_a,r_b  ⇒ r_a
 8   mult     r_a,r_c  ⇒ r_a
 9   mult     r_a,r_d  ⇒ r_a
10   storeAI  r_a      ⇒ r_arp,@a
```
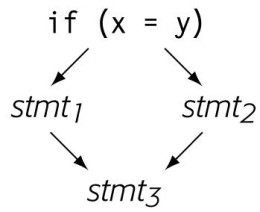


(a) Example Code from Chapter 1    (b) Dependence Graph for the Example

```
if (x = y)
    then stmt₁
    else stmt₂
stmt₃
```

```
        if (x = y)
        ↙        ↘
   stmt₁          stmt₂
        ↘        ↙
        stmt₃
```

Back end

$$a \leftarrow a \times 2 \times b \times c \times d$$

```
loadAI   r_arp, @a ⇒ r_a    // load 'a'
loadI    2         ⇒ r_2    // constant 2 into r_2
loadAI   r_arp,@b ⇒ r_b     // load 'b'
loadAI   r_arp,@c ⇒ r_c     // load 'c'
loadAI   r_arp,@d ⇒ r_d     // load 'd'
mult     r_a,r_2  ⇒ r_a     // r_a ← a × 2
mult     r_a,r_b  ⇒ r_a     // r_a ← (a × 2) × b
mult     r_a,r_c  ⇒ r_a     // r_a ← (a × 2 × b) × c
mult     r_a,r_d  ⇒ r_a     // r_a ← (a × 2 × b × c) × d
storeAI  r_a      ⇒ r_arp,@a // write r_a back to 'a'
```

```
addi   sp,sp,-32
sw     ra,28(sp)
sw     s0,24(sp)
addi   s0,sp,32
sw     a0,-20(s0)
sw     a1,-24(s0)
sw     a2,-28(s0)
sw     a3,-32(s0)
lw     a4,-20(s0)
lw     a5,-24(s0)
mul    a4,a4,a5
lw     a5,-28(s0)
mul    a4,a4,a5
lw     a5,-32(s0)
mul    a5,a4,a5
slli   a5,a5,1
mv     a0,a5
lw     ra,28(sp)
lw     s0,24(sp)
addi   sp,sp,32
jr     ra
```

**Formal languages**

Alphabet

$\Sigma = \{a, b, c, ..., z\}$

Σ = {0, 1}

Σ = {false, true}

Σ = English words

String

abcdababab

11100011001

ε

00000...                                    'i' 'like' 'six' 'oh' 'three' 'five'

Language                L = {1, 01, 10, 001, 010, 100, 0001, 0010, 0100, 1000,
                              00001, 00010, 00100, 01000, 10000, 000001, ...}

                                    (assuming Σ = {0, 1})

L = set of binary strings that contain exactly one 1

(assuming Σ = {0, 1})


L(s) = whether s contains exactly one 1 (yes or no)

(assuming Σ = {0, 1})


L = set of decimal numbers that are divisible by 3

(assuming Σ = {0, 1, 2, ..., 9})


L = set of valid hexadecimal numbers

(assuming Σ = ASCII characters)


L = set of syntactically valid Python programs

(assuming Σ = ASCII characters)


L = set of syntactically valid Python programs

(assuming Σ = Python tokens)

L = set of Python source interpretable without error

**Regular languages**

Regular expression

`(617|857)-253-(0|1|..|9)(0|1|..|9)(0|1|..|9)(0|1|..|9)`

0  Empty string

0  Empty string
ε

1 A letter from Σ
0

2 Concatenation
$a \cdot b$

2 Concatenation
$ab$

2 Concatenation
234324

3 Alternation
$a \,|\, b$

3 Alternation/Union
$a \cup b$

3  Alternation
`0|1`

3  Alternation
`000|001|100|101`

3  Alternation
`(0|1)0(0|1)`

3  Alternation
`(617|857)-253-(0|1|..|9)(0|1|..|9)(0|1|..|9)(0|1|..|9)`

4  Kleene star
`a*`

4  Kleene star
`(0|1)*`
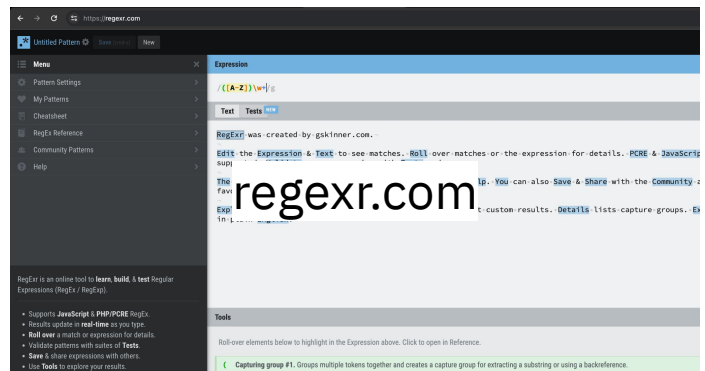
4 Kleene star
0*10*

"regex"

[0-9A-Fa-f]

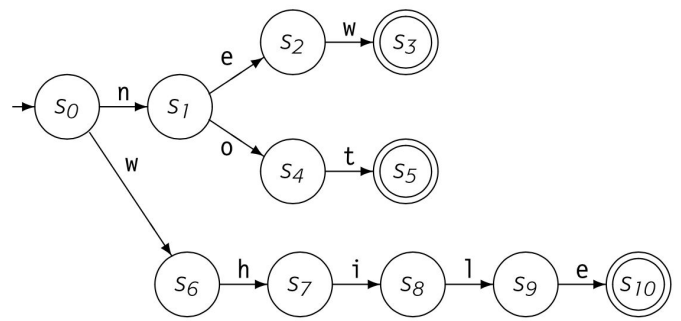[01]{8}

0|(1[01]*)

0x[0-9A-Fa-f]+

0x[0-9A-Fa-f][0-9A-Fa-f]*
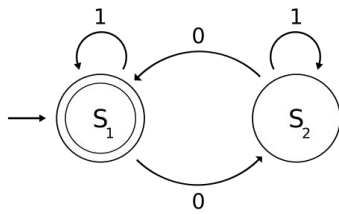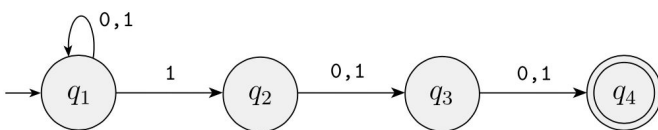
617-253-\d{4}

.*



Deterministic Finite Automata

$M = (Q, \ \Sigma, \delta, q_0, F)$

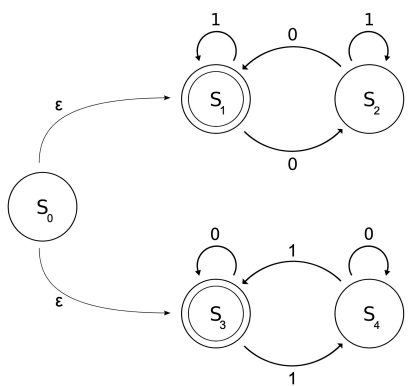**Non-deterministic Finite Automata**

States: $S_1$ (start, accepting), $S_2$

- $S_1 \xrightarrow{1} S_1$ (self-loop)
- $S_1 \xrightarrow{0} S_2$
- $S_2 \xrightarrow{1} S_2$ (self-loop)
- $S_2 \xrightarrow{0} S_1$

---

$S_0 \xrightarrow{n} S_1$

$S_1 \xrightarrow{e} S_2 \xrightarrow{w} S_3$

$S_1 \xrightarrow{o} S_4 \xrightarrow{t} S_5$

$S_0 \xrightarrow{w} S_6 \xrightarrow{h} S_7 \xrightarrow{i} S_8 \xrightarrow{l} S_9 \xrightarrow{e} S_{10}$

**L = set of binary strings containing a 1 in the third position from the end**

---

$q_1 \xrightarrow{0,1} q_1$ (self-loop)

$q_1 \xrightarrow{1} q_2 \xrightarrow{0,1} q_3 \xrightarrow{0,1} q_4$

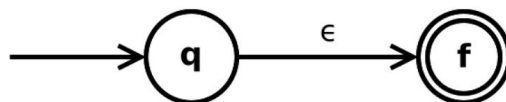**L = set of binary strings with even number of 0s or even number of 1s**

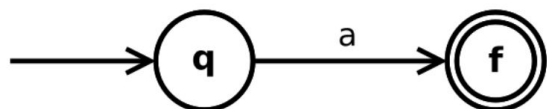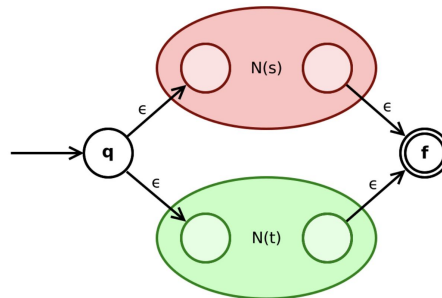Regex → NFA

Thompson's construction

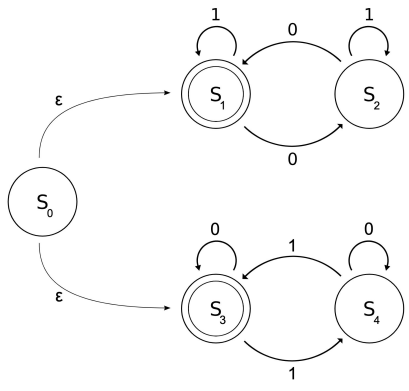The **empty-expression** ε is converted to



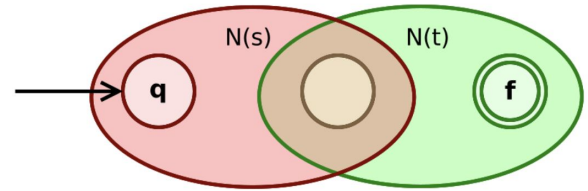A **symbol** *a* of the input alphabet is converted to



The **union expression** *s|t* is converted to

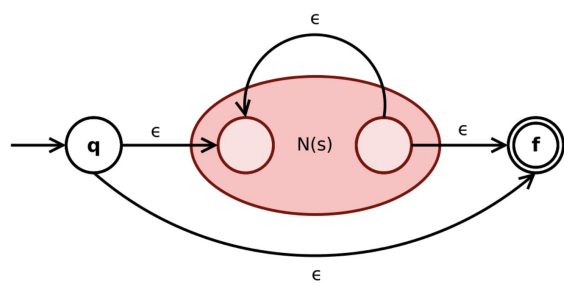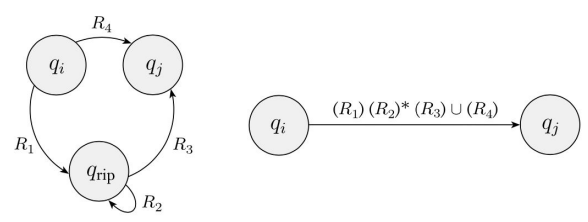The **concatenation expression** *st* is converted to

N(s)  N(t)

q  f

The **Kleene star** expression $s^*$ is converted to



DFA/NFA → Regex

Generalized NFA



$R_4$

$q_i$  $q_j$

$R_1$  $R_3$

$q_{rip}$

$R_2$

$q_i$  $(R_1)(R_2)^*(R_3) \cup (R_4)$  $q_j$

$a$   $b$   $2$

$1$   $a$   $3$   $b$   $a$

(a)

$b$

$1$   $a$   $2$   $\varepsilon$   $a$

$s$   $\varepsilon$   $b$   $3$   $a$   $\varepsilon$

(b)

$aa \cup b$

$2$

$s$   $a$   $ab$   $ba \cup a$   $a$   $\varepsilon$

$b$   $3$   $\varepsilon$

$bb$

(c)

$a(aa \cup b)^*$

$s$   $a$

$a(aa \cup b)^* ab \cup b$   $(ba \cup a)(aa \cup b)^* \cup \varepsilon$

$3$

$(ba \cup a)(aa \cup b)^* ab \cup bb$

(d)

$s$   $a$

$(a(aa \cup b)^* ab \cup b)((ba \cup a)(aa \cup b)^* ab \cup bb)^*((ba \cup a)(aa \cup b)^* \cup \varepsilon) \cup a(aa \cup b)^*$

NFA → DFA

$0,1$

$q_1$   $1$   $q_2$   $0,1$   $q_3$   $0,1$   $q_4$

DFA → NFA

$0$

$q_{000}$   $0$   $q_{100}$   $0$   $q_{010}$   $0$   $q_{110}$

$1$   $1$   $0$   $1$   $1$   $0$   $0$

$q_{001}$   $q_{101}$   $1$   $q_{011}$   $1$   $q_{111}$   $1$

$1$

Trivial

DFA Minimization                                    Hard

Why?



Greedy

L = set of binary strings that start with 0s, followed by an equal number of 1s

Pumping lemma

Take 6.045 or 6.840.

For the quiz, you should know how to:

- Simulate regex/DFA/NFA
- Design a regex/DFA/NFA
- Convert regex to NFA
- Convert NFA to DFA

How to practice: Do textbook exercises!

**Context-free grammar**

$$\langle\text{SENTENCE}\rangle \rightarrow \langle\text{NOUN-PHRASE}\rangle\langle\text{VERB-PHRASE}\rangle$$
$$\langle\text{NOUN-PHRASE}\rangle \rightarrow \langle\text{CMPLX-NOUN}\rangle \mid \langle\text{CMPLX-NOUN}\rangle\langle\text{PREP-PHRASE}\rangle$$
$$\langle\text{VERB-PHRASE}\rangle \rightarrow \langle\text{CMPLX-VERB}\rangle \mid \langle\text{CMPLX-VERB}\rangle\langle\text{PREP-PHRASE}\rangle$$
$$\langle\text{PREP-PHRASE}\rangle \rightarrow \langle\text{PREP}\rangle\langle\text{CMPLX-NOUN}\rangle$$
$$\langle\text{CMPLX-NOUN}\rangle \rightarrow \langle\text{ARTICLE}\rangle\langle\text{NOUN}\rangle$$
$$\langle\text{CMPLX-VERB}\rangle \rightarrow \langle\text{VERB}\rangle \mid \langle\text{VERB}\rangle\langle\text{NOUN-PHRASE}\rangle$$
$$\langle\text{ARTICLE}\rangle \rightarrow \texttt{a} \mid \texttt{the}$$
$$\langle\text{NOUN}\rangle \rightarrow \texttt{boy} \mid \texttt{girl} \mid \texttt{flower}$$
$$\langle\text{VERB}\rangle \rightarrow \texttt{touches} \mid \texttt{likes} \mid \texttt{sees}$$
$$\langle\text{PREP}\rangle \rightarrow \texttt{with}$$

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 1 | *Expr* | $\rightarrow$ | ( *Expr* ) | 4 | *Op* | $\rightarrow$ | + |
| 2 | | \| | *Expr Op Expr* | 5 | | \| | - |
| 3 | | \| | name | 6 | | \| | × |
| | | | | 7 | | \| | ÷ |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 1 | *Expr* | $\rightarrow$ | ( *Expr* ) | 4 | *Op* | $\rightarrow$ | + |
| 2 | | \| | *Expr Op Expr* | 5 | | \| | - |
| 3 | | \| | name | 6 | | \| | × |
| | | | | 7 | | \| | ÷ |

(a + b) × c

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 1 | *Expr* | $\rightarrow$ | ( *Expr* ) | 4 | *Op* | $\rightarrow$ | + |
| 2 | | \| | *Expr Op Expr* | 5 | | \| | - |
| 3 | | \| | name | 6 | | \| | × |
| | | | | 7 | | \| | ÷ |



| | | | |
|---|---|---|---|
| 1 | *Stmt* | $\rightarrow$ | if *Expr* then *Stmt* |
| 2 | | \| | if *Expr* then *Stmt* else *Stmt* |
| 3 | | \| | *Other* |

Ambiguity

## Left factoring

| | | | |
|---|---|---|---|
| 1 | *Stmt* | → | `if` *Expr* `then` *Stmt* |
| 2 | | \| | `if` *Expr* `then` *Stmt* `else` *Stmt* |
| 3 | | \| | *Other* |

| | | | |
|---|---|---|---|
| 1 | *Stmt* | → | `if` *Expr* `then` *Stmt* |
| 2 | | \| | `if` *Expr* `then` *WithElse* `else` *Stmt* |
| 3 | | \| | *Other* |
| 4 | *WithElse* | → | `if` *Expr* `then` *WithElse* `else` *WithElse* |
| 5 | | \| | *Other* |

## Precedence climbing

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 1 | *Expr* | → | ( *Expr* ) | | 4 | *Op* | → | + |
| 2 | | \| | *Expr* *Op* *Expr* | | 5 | | \| | - |
| 3 | | \| | `name` | | 6 | | \| | × |
| | | | | | 7 | | \| | ÷ |

## Pushdown Automata

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 1 | *Expr* | → | ( *Expr* ) | | 4 | *Op* | → | + |
| 2 | | \| | *Expr* *Op* *Expr* | | 5 | | \| | - |
| 3 | | \| | `name` | | 6 | | \| | × |
| | | | | | 7 | | \| | ÷ |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 0 | *Goal* | → | *Expr* | | 5 | | \| | *Term* ÷ *Factor* |
| 1 | *Expr* | → | *Expr* + *Term* | | 6 | | \| | *Factor* |
| 2 | | \| | *Expr* - *Term* | | 7 | *Factor* | → | ( *Expr* ) |
| 3 | | \| | *Term* | | 8 | | \| | `num` |
| 4 | *Term* | → | *Term* × *Factor* | | 9 | | \| | `name` |

L = set of binary strings that start with 0s,
followed by an equal number of 1s

$q_1 \xrightarrow{\varepsilon,\varepsilon \to \$} q_2$

$0,\varepsilon \to 0$

$1,0 \to \varepsilon$

$1,0 \to \varepsilon$

$q_4 \xleftarrow{\varepsilon,\$ \to \varepsilon} q_3$

Remember this?

```
expression ::= resize ('|' resize)*;
resize     ::= primitive ('@' size)*;
size       ::= (number 'x' number);
primitive  ::= filename | '(' expression ')';

topToBottomOperator ::= '---' '-'*;
filename            ::= [A-Za-z0-9.][A-Za-z0-9._-]*;
number              ::= [0-9]+;
whitespace          ::= [ \\t\\r\\n]+;
```

## Extended Backus–Naur form*
(* in spirit)

$A \ ::= \ B\star \ C$

$$A \ ::= \ A' \ C$$
$$A' \ ::= \ \varepsilon \ | \ BA'$$

For the quiz, you should know how to:

- Parse a string using a given grammar (draw parse trees)
- Eliminate ambiguity
- Fix precedence issues
    - Make sure you understand the arithmetic examples.
    - Reminder: You can collaborate/ask for help on miniquiz.

How to practice: Do textbook exercises!

**Top-down parsing**

Recursive descent parser

<rant>

Use first principles

Ask TAs

</rant>

Project 1!

Left factoring (again)

$Factor \rightarrow$ name
$\quad | \quad$ name $[\ ArgList\ ]$
$\quad | \quad$ name $(\ ArgList\ )$
$ArgList \rightarrow Expr\ MoreArgs$
$MoreArgs \rightarrow$ , $Expr\ MoreArgs$
$\quad | \quad \epsilon$

$Factor \rightarrow$ name $Arguments$
$Arguments \rightarrow [\ ArgList\ ]$
$\quad | \quad (\ ArgList\ )$
$\quad | \quad \epsilon$
$ArgList \rightarrow Expr\ MoreArgs$
$MoreArgs \rightarrow$ , $Expr\ MoreArgs$
$\quad | \quad \epsilon$

Left recursion

```
Expr   ::= Expr + Term
         | Expr - Term
         | Term;
Term   ::= Term × Factor
         | Term ÷ Factor
         | Factor;
Factor ::= ( Expr )
         | num
         | name;
```

$$Fee \rightarrow Fee\,\alpha \qquad\qquad Fee \rightarrow \beta\,Fee'$$
$$\qquad\quad |\quad \beta \qquad\qquad\qquad\quad Fee' \rightarrow \alpha\,Fee'$$
$$\qquad\qquad\qquad\qquad\qquad\qquad\quad\; |\quad \epsilon$$

```
Expr   ::= Term Expr';
Expr'  ::= + Term Expr'
         | - Term Expr'
         | ε;
Term   ::= Factor Term'
Term'  ::= × Factor Term'
         | ÷ Factor Term'
         | ε;
Factor ::= ( Expr )
         | num
         | name
```

```
Expr   ::= Term ((+|-) Term)*
Term   ::= Factor ((×|÷) Factor)*
Factor ::= ( Expr )
         | num
         | name;
```

Indirect left recursion                    Constraint propagation

$$NT \to \varepsilon$$
$$\Rightarrow$$
$$NT \to^* \varepsilon$$

$$NT_0 \to NT_1 NT_2 \ldots \text{ and } NT_i \to^* \varepsilon$$
$$\Rightarrow$$
$$NT_0 \to^* \varepsilon$$

$$\mathbf{T} \in \text{First}(\mathbf{T})$$

$$x \in \text{First}(S)$$
$$\Rightarrow$$
$$x \in \text{First}(S\, S_1\, S_2\, S_3 \ldots)$$

$$x \in \text{First}(S)$$
$$\Rightarrow$$
$$x \in \text{First}(S\beta)$$

$$\text{First}(S) \subseteq \text{First}(S\beta)$$

$$x \in \text{First}(\beta) \quad \text{and} \quad NT \rightarrow^* \varepsilon$$
$$\Rightarrow$$
$$x \in \text{First}(NT\ \beta)$$

$$x \in \text{First}(S\beta) \quad \text{and} \quad (NT \rightarrow S\beta)$$
$$\Rightarrow$$
$$x \in \text{First}(NT)$$