# Check-in
## Review: Flex

Write a Flex rule for Fortran real literals:
- An optional sign (+ or -)
- One of the following:
  - An integer
  - One or more digits followed by a '.' followed by 0 or more digits
  - a '.' followed by one or more digits

# Administrivia
## Housekeeping

**P1 Out**

EECS 665

COMPILER CONSTRUCTION

Lecture 3
Syntactic Definition
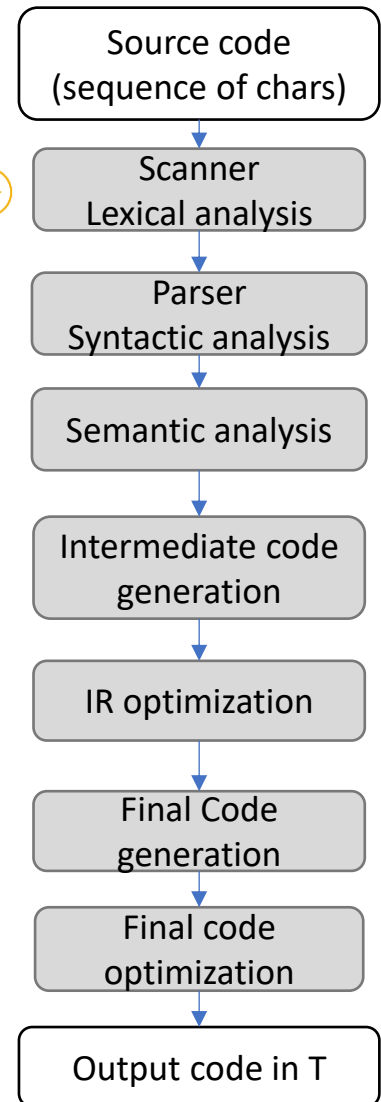
# Compiler Construction
Progress Pics

**Done:**

- Shown RegEx as a good token formalism

- Lexical specification

- Lexical recognition

In progress 🕐

```
┌─────────────────────────┐
│      Source code        │
│  (sequence of chars)    │
└─────────────────────────┘
           │
┌─────────────────────────┐
│        Scanner          │
│    Lexical analysis     │
└─────────────────────────┘
           │
┌─────────────────────────┐
│         Parser          │
│    Syntactic analysis   │
└─────────────────────────┘
           │
┌─────────────────────────┐
│    Semantic analysis    │
└─────────────────────────┘
           │
┌─────────────────────────┐
│   Intermediate code     │
│      generation         │
└─────────────────────────┘
           │
┌─────────────────────────┐
│     IR optimization     │
└─────────────────────────┘
           │
┌─────────────────────────┐
│       Final Code        │
│      generation         │
└─────────────────────────┘
           │
┌─────────────────────────┐
│       Final code        │
│      optimization       │
└─────────────────────────┘
           │
┌─────────────────────────┐
│     Output code in T    │
└─────────────────────────┘
```

# Last Time
## Lecture Review

RegEx → $\varepsilon$-NFA → $\varepsilon$-free NFA → DFA

**Thompson's Construction Algorithm**

**$\varepsilon$-elimination**

**Rabin–Scott Powerset Construction**

Replace sub-RegExes with sub-FSMs "bottom-up" in the expression tree

Use the $\varepsilon$-closure to "bypass middleman" states and transitions

Create a DFA that tracks all possible states the NFA could be in

**Good news**

DFAs have a natural implementation (use a 2-D array)

**Bad news**

DFAs don't exactly do tokenization

# Today's Lecture
## Outline

**Finish lexical analysis**

- How to build a scanner

**Begin discussing syntax**

- How to specify the syntax of a programming language

# DFA ≠ Tokenizer
## Limitations

- Finite automata only check for language membership of a string (recognition)

- The Scanner needs to
  - Break the input into many different tokens
  - Know what characters comprise the token

**input string**

int a; b = 3 + 2;

**DFA**



**recognition**

👍 or 👎

# DFA ≠ Tokenizer

Limitations

- Finite automata only check for language membership of a string (recognition)

- The Scanner needs to
  - Break the input into many different tokens
  - Know what characters comprise the token

We need to go… *beyond recognition*

# Lecture Outline
## Syntactic Definition

**From DFAs to Tokenizer**

- Algorithms

- Implementation details

**How Language Syntax is Formally Defined: CFGs**

- Why we need context-free grammars

- How we use context-free grammars

# From RegExes to Tokenizer
## Algorithms

| RegEx | Thompson's Construction → | $\varepsilon$-NFA | $\varepsilon$-elimination → | $\varepsilon$-free NFA | Powerset Construction → | DFA |

**Language**

true   { Token(**true**) }
'+'   { Token(**cross**) }
'-'   { Token(**dash**) }

**true**

$T_0 \xrightarrow{\text{'t'}} T_1 \xrightarrow{\text{'r'}} T_2 \xrightarrow{\text{'u'}} T_3 \xrightarrow{\text{'e'}} T_4$

**cross**

$C_0 \xrightarrow{\text{'+'}} C_1$

**dash**

$D_0 \xrightarrow{\text{'-'}} D_1$

# From RegExes to Tokenizer
## Algorithms

| RegEx | →(Thompson's Construction)→ | $\varepsilon$-NFA | →($\varepsilon$-elimination)→ | $\varepsilon$-free NFA | →(Powerset Construction)→ | DFA | →(???)→ | Tokenizer |

---
## 1st Idea (flawed)

Consume char stream to **accept** state: return accepted token, restart DFAs with next char

**Language**

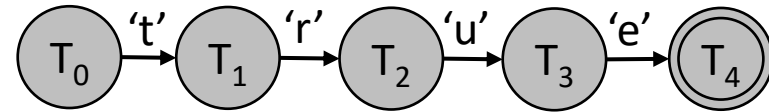| | |
|---|---|
| true | { Token(**true**) } |
| '+' | { Token(**cross**) } |
| '-' | { Token(**dash**) } |

**true**

$T_0 \xrightarrow{\text{'t'}} T_1 \xrightarrow{\text{'r'}} T_2 \xrightarrow{\text{'u'}} T_3 \xrightarrow{\text{'e'}} T_4$

**cross**

$C_0 \xrightarrow{\text{'+'}} C_1$

**dash**

$D_0 \xrightarrow{\text{'-'}} D_1$
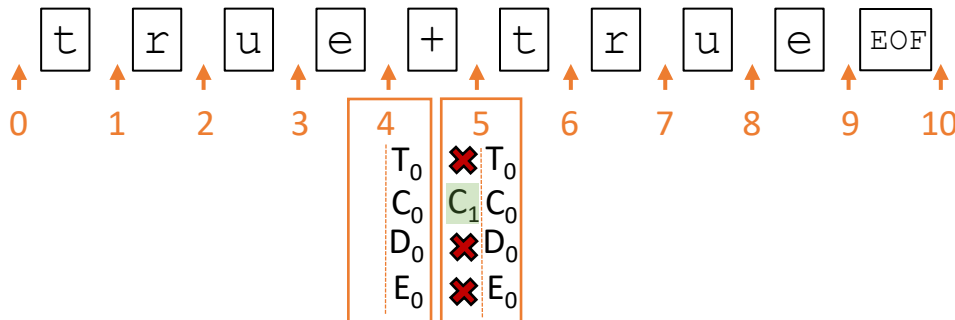
# From RegExes to Tokenizer

Algorithms

---

1st Idea (flawed)

---

Consume char stream to **accept** state: return accepted token, restart DFAs with next char
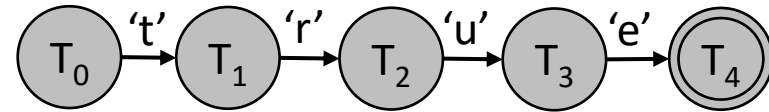
**Language**

true     { Token(**true**) }
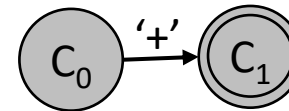'+'      { Token(**cross**) }
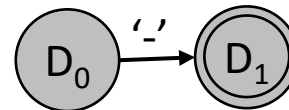'-'      { Token(**dash**) }

**Char Stream**

| t | r | u | e | + | t | r | u | e | EOF |

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |

| $T_0$ | $T_1$ | $T_2$ | $T_3$ | $T_4$ $T_0$ |
| $C_0$ | ✖ | ✖ | ✖ | ✖ $C_0$ |
| $D_0$ | ✖ | ✖ | ✖ | ✖ $D_0$ |
| $E_0$ | ✖ | ✖ | ✖ | ✖ $E_0$ |

**Token Stream**

**true**
[0,4)

**true**

$T_0$ —'t'→ $T_1$ —'r'→ $T_2$ —'u'→ $T_3$ —'e'→ $T_4$

**cross**

$C_0$ —'+'→ $C_1$

**dash**

$D_0$ —'-'→ $D_1$

**eof**

$E_0$ —EOF→ $E_1$
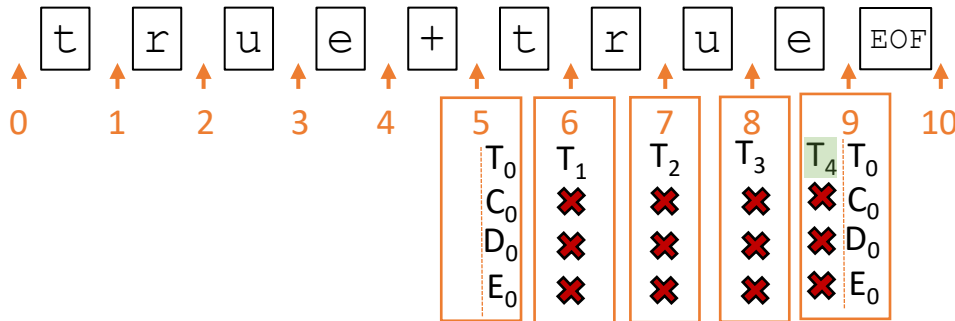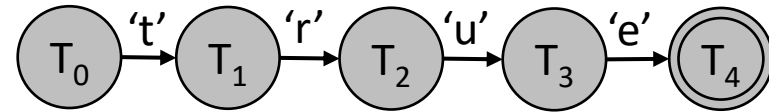
**Special State:**
Return EOF
Accept stream

# From RegExes to Tokenizer
## Algorithms

Consume char stream to **accept** state: return accepted token, restart DFAs with next char

## Language

true    { Token(**true**) }
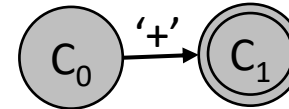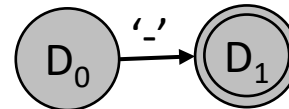'+'    { Token(**cross**) }
'-'    { Token(**dash**) }

**true**

$T_0$ —'t'→ $T_1$ —'r'→ $T_2$ —'u'→ $T_3$ —'e'→ $T_4$

## Char Stream

| t | r | u | e | + | t | r | u | e | EOF |

0   1   2   3   4   5   6   7   8   9   10

4: $T_0$ $C_0$ $D_0$ $E_0$
5: ✖$T_0$ $C_1$✖$C_0$ ✖$D_0$ ✖$E_0$

**cross**

$C_0$ —'+'→ $C_1$

**dash**

$D_0$ —'-'→ $D_1$

## Token Stream

**true** [0,4)    **cross** [4,5)

**eof**

$E_0$ —EOF→ $E_1$
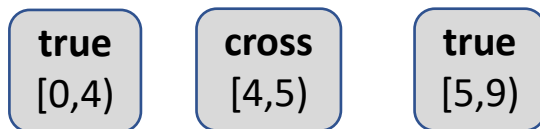
**Special State:**
Return EOF
Accept stream

# From RegExes to Tokenizer
## Algorithms

---
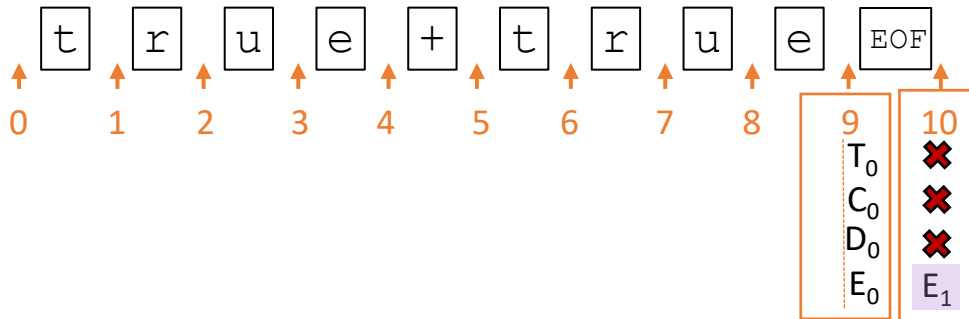
### 1st Idea (flawed)

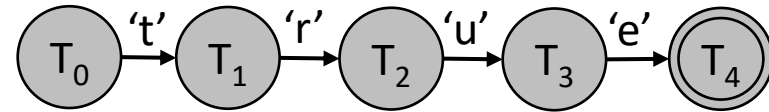Consume char stream to **accept** state: return accepted token, restart DFAs with next char

**Language**

| | |
|---|---|
| true | { Token(**true**) } |
| '+' | { Token(**cross**) } |
| '-' | { Token(**dash**) } |

**Char Stream**



**Token Stream**

| | | |
|---|---|---|
| **true** [0,4) | **cross** [4,5) | **true** [5,9) |

**true**



**cross**



**dash**



**eof**



**Special State:** Return EOF Accept stream

# From RegExes to Tokenizer
Algorithms

---

Consume char stream to **accept** state: return accepted token, restart DFAs with next char

## Language

true    { Token(**true**) }
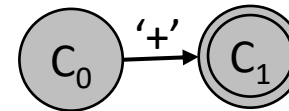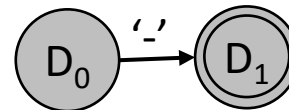'+'    { Token(**cross**) }
'-'    { Token(**dash**) }

**true**

$T_0$ --'t'--> $T_1$ --'r'--> $T_2$ --'u'--> $T_3$ --'e'--> $T_4$

## Char Stream

| t | r | u | e | + | t | r | u | e | EOF |

0  1  2  3  4  5  6  7  8  9  10

|  | 9 | 10 |
|---|---|---|
| $T_0$ | | ✖ |
| $C_0$ | | ✖ |
| $D_0$ | | ✖ |
| $E_0$ | | $E_1$ |

**cross**

$C_0$ --'+'--> $C_1$

**dash**

$D_0$ --'-'--> $D_1$

## Token Stream

| **true** [0,4) | **cross** [4,5) | **true** [5,9) | **eof** [9,10) |

**eof**

$E_0$ --EOF--> $E_1$

**Special State:**
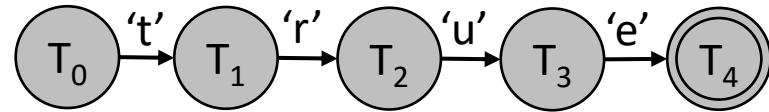Return EOF
Accept stream

# From RegExes to Tokenizer
## Algorithms

Consume char stream to **accept** state: return accepted token, restart DFAs with next char

## Language

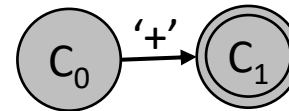true    { Token(**true**) }
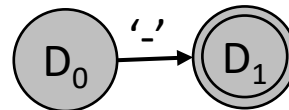
'+'    { Token(**cross**) }

'-'    { Token(**dash**) }
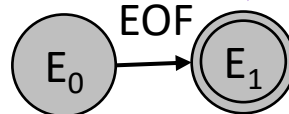
**true**

$T_0$ —'t'→ $T_1$ —'r'→ $T_2$ —'u'→ $T_3$ —'e'→ $T_4$

## Char Stream

| t | r | u | e | + | t | r | u | e | EOF |
|---|---|---|---|---|---|---|---|---|-----|

0   1   2   3   4   5   6   7   8   9   10

**cross**

$C_0$ —'+'→ $C_1$

**dash**

$D_0$ —'-'→ $D_1$

**Accept!**

**eof**

$E_0$ —EOF→ $E_1$

**Special State:**
Return EOF
Accept stream

## Token Stream

| **true** [0,4) | **cross** [4,5) | **true** [5,9) | **eof** [9,10) |
|---|---|---|---|

What happens when token languages overlap / prefix each other?

# From RegExes to Tokenizer
### Algorithms

---
## 1st Idea (flawed)

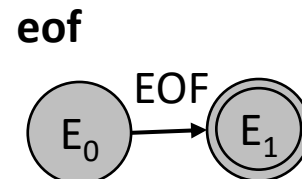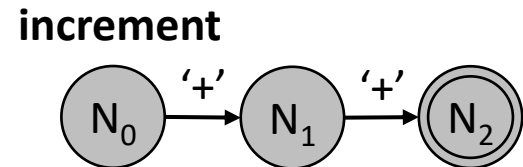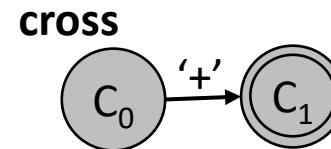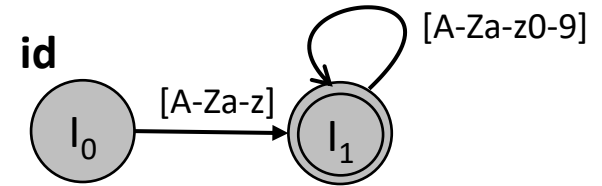Consume char stream to **accept** state: return accepted token, restart DFAs with next char
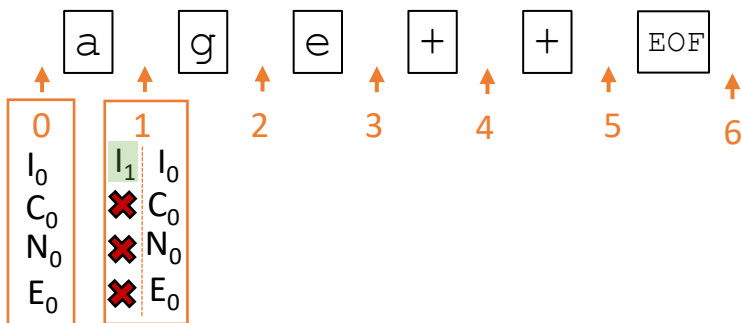
---
## Problem

What happens when token languages overlap / prefix each other?

### Language

[A-Za-z][A-Za-z0-9]*    { Token(**id**) }

+    { Token(**cross**) }

++    { Token(**increment**) }

### Char Stream

| a | g | e | + | + | EOF |
|---|---|---|---|---|-----|

0   1   2   3   4   5   6

| 0 | 1 | | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|---|
| $I_0$ | $I_1$ | $I_0$ | | | | | |
| $C_0$ | ✖ | $C_0$ | | | | | |
| $N_0$ | ✖ | $N_0$ | | | | | |
| $E_0$ | ✖ | $E_0$ | | | | | |

### Token Stream

**id: a**
[0,1)

**id**

$I_0$ —[A-Za-z]→ $I_1$ , $I_1$ —[A-Za-z0-9]→ $I_1$ (self loop)

**cross**

$C_0$ —'+'→ $C_1$

**increment**

$N_0$ —'+'→ $N_1$ —'+'→ $N_2$

**eof**

$E_0$ —EOF→ $E_1$

# From RegExes to Tokenizer

Algorithms

---
1st Idea (flawed)
---

Consume char stream to **accept** state: return accepted token, restart DFAs with next char

---
Problem
---

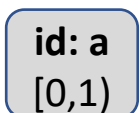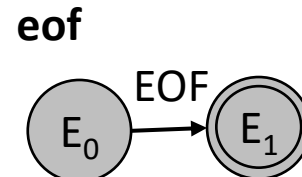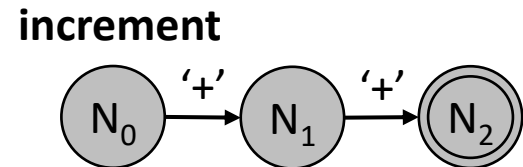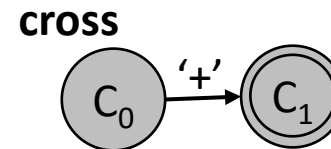What happens when token languages overlap / prefix each other?

## Language

[A-Za-z][A-Za-z0-9]*    { Token(**id**) }

$+$    { Token(**cross**) }

$++$    { Token(**increment**) }

**id**

$I_0$ — [A-Za-z] → $I_1$ with self-loop [A-Za-z0-9]

## Char Stream

| a | g | e | + | + | EOF |

0    1    2    3    4    5    6

Position 1: $I_0$, $C_0$, $N_0$, $E_0$

Position 2: $I_1$ ✖, $I_0$ ✖, $C_0$, $N_0$, $E_0$

**cross**

$C_0$ — '+' → $C_1$

**increment**

$N_0$ — '+' → $N_1$ — '+' → $N_2$

## Token Stream

**id: a**
[0,1)

**id: g**
[1,2)

**eof**

$E_0$ — EOF → $E_1$

# From RegExes to Tokenizer
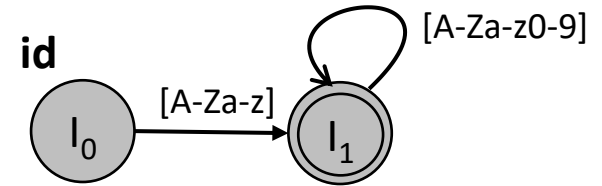
Algorithms

---
1st Idea (flawed)
---

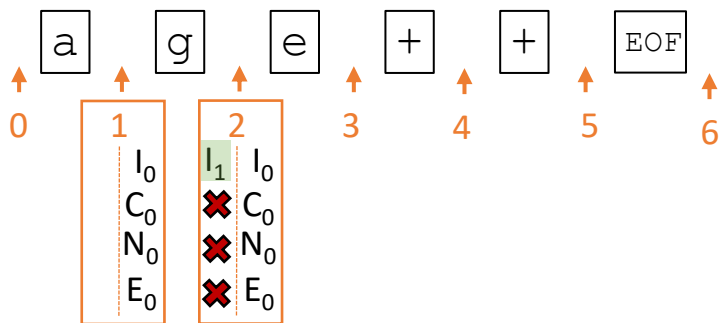Consume char stream to **accept** state: return accepted token, restart DFAs with next char

---
Problem
---

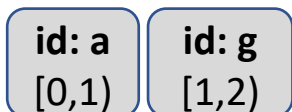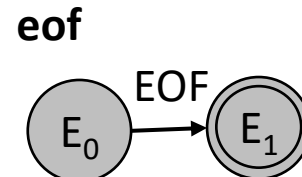What happens when token languages overlap / prefix each other?

## Language

[A-Za-z][A-Za-z0-9]*  { Token(**id**) }

+  { Token(**cross**) }

++  { Token(**increment**) }

**id**

$I_0$ →[A-Za-z]→ $I_1$  [A-Za-z0-9]

**cross**

$C_0$ →'+'→ $C_1$

**increment**

$N_0$ →'+'→ $N_1$ →'+'→ $N_2$

**eof**

$E_0$ →EOF→ $E_1$

## Char Stream

| a | g | e | + | + | EOF |

0  1  2  3  4  5  6

2:
$I_0$
$C_0$
$N_0$
$E_0$

3:
$I_1$ ✖ ✖ ✖ | $I_0$ $C_0$ $N_0$ $E_0$

## Token Stream

| **id: a** [0,1) | **id: g** [1,2) | **id: e** [2,3) |

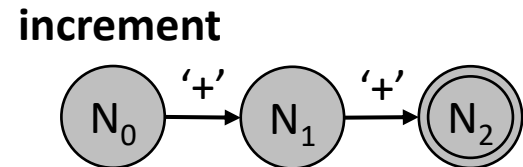# From RegExes to Tokenizer
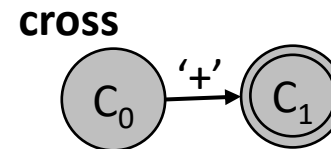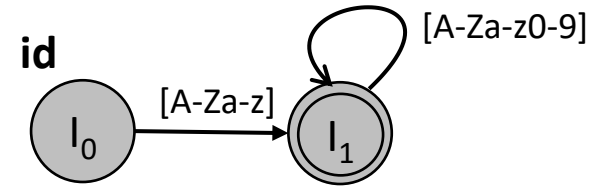## Algorithms

---
### 1st Idea (flawed)
---

Consume char stream to **accept** state: return accepted token, restart DFAs with next char
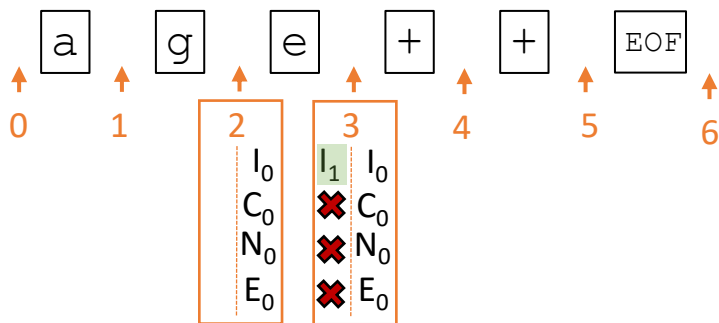
---
### Problem
---

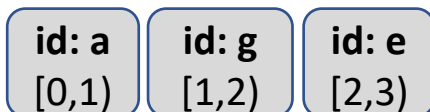What happens when token languages overlap / prefix each other?

## Language

[A-Za-z][A-Za-z0-9]*   { Token(**id**) }

$+$   { Token(**cross**) }

$++$   { Token(**increment**) }

**id**

$I_0 \xrightarrow{[A-Za-z]} I_1 \circlearrowright [A-Za-z0-9]$

## Char Stream

| a | g | e | + | + | EOF |

0   1   2   3   4   5   6

3:
$I_0$
$C_0$
$N_0$
$E_0$

4:
✖ $I_0$
$C_1$ $C_0$
$N_1$ $N_0$
✖ $E_0$

**cross**

$C_0 \xrightarrow{'+'} C_1$

**increment**

$N_0 \xrightarrow{'+'} N_1 \xrightarrow{'+'} N_2$

**eof**

$E_0 \xrightarrow{EOF} E_1$

## Token Stream

| **id: a** [0,1) | **id: g** [1,2) | **id: e** [2,3) | **cross** [3,4) |

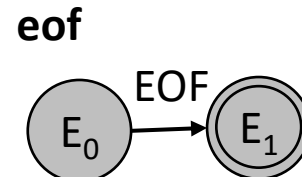# From RegExes to Tokenizer
## Algorithms

---
## 1st Idea (flawed)
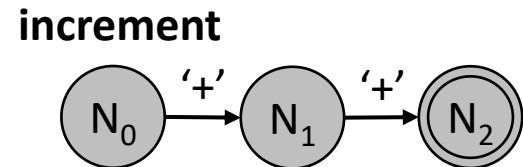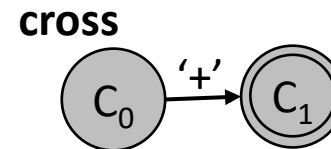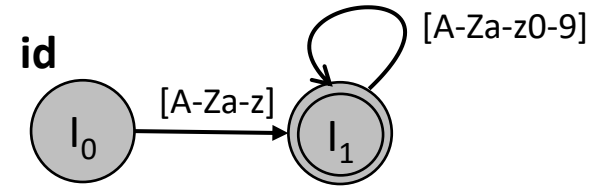
Consume char stream to **accept** state: return accepted token, restart DFAs with next char
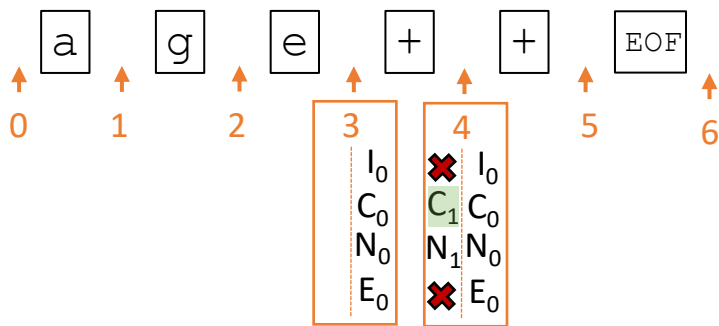
---
## Problem

What happens when token languages overlap / prefix each other?

### Language

[A-Za-z][A-Za-z0-9]*    { Token(**id**) }

+    { Token(**cross**) }

++    { Token(**increment**) }

**id**

$I_0$ →[A-Za-z]→ $I_1$ ⟲ [A-Za-z0-9]

**cross**

$C_0$ →'+'→ $C_1$

**increment**

$N_0$ →'+'→ $N_1$ →'+'→ $N_2$

**eof**

$E_0$ →EOF→ $E_1$

### Char Stream

| a | g | e | + | + | EOF |

0    1    2    3    4    5    6

Position 4:
$I_0$
$C_0$
$N_0$
$E_0$

Position 5:
❌ $I_0$
$C_1$ $C_0$
$N_1$ $N_0$
❌ $E_0$

### Token Stream

| **id: a** [0,1) | **id: g** [1,2) | **id: e** [2,3) | **cross** [3,4) | **cross** [4,5) |

# From RegExes to Tokenizer
## Algorithms

---
1st Idea (flawed)
---

Consume char stream to **accept** state: return accepted token, restart DFAs with next char

---
Problem
---

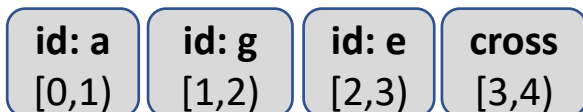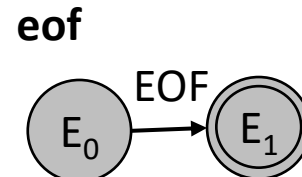What happens when token languages overlap / prefix each other?

## Language

[A-Za-z][A-Za-z0-9]*  { Token(**id**) }

+  { Token(**cross**) }

++  { Token(**increment**) }

**id**



## Char Stream



**cross**



**increment**



## Token Stream

| id: a | id: g | id: e | cross | cross | eof |
|-------|-------|-------|-------|-------|-----|
| [0,1) | [1,2) | [2,3) | [3,4) | [4,5) | [5,6) |

**eof**

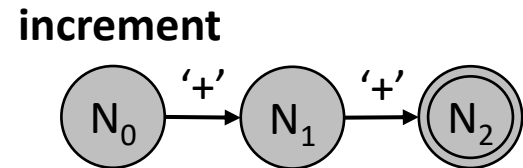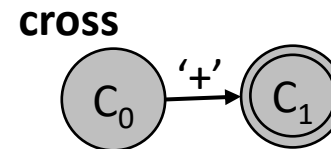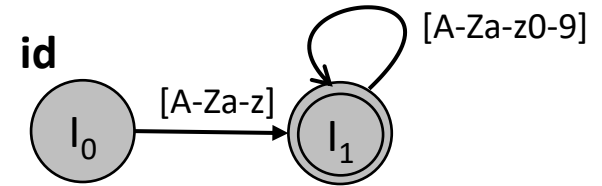# From RegExes to Tokenizer
## Algorithms

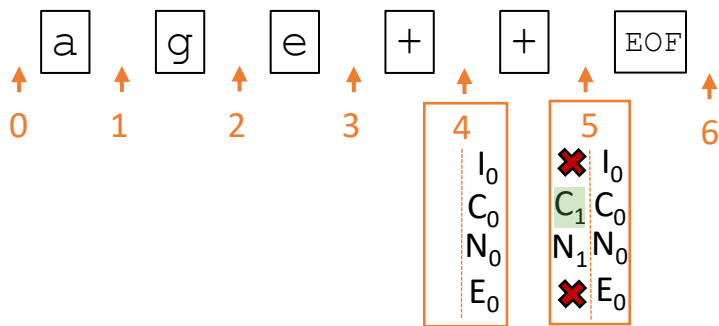Consume char stream to **accept** state: return accepted token, restart DFAs with next char

Problem

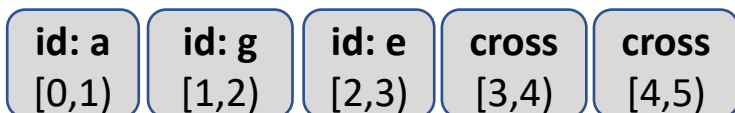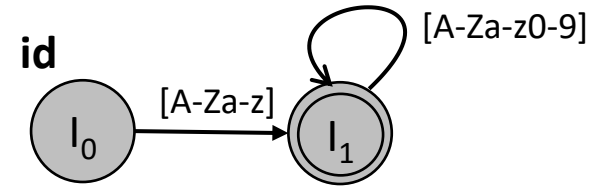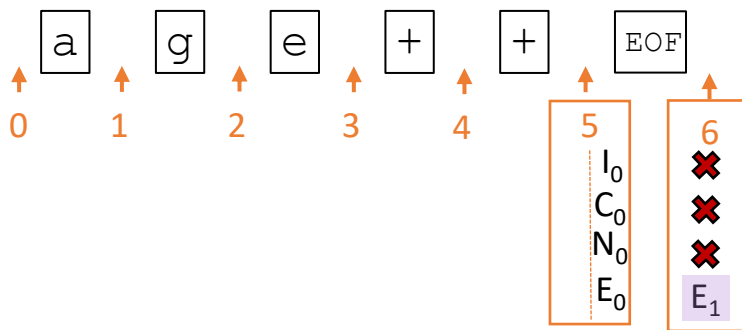What happens when token languages overlap / prefix each other?

**Language**

[A-Za-z][A-Za-z0-9]*    { Token(**id**) }

+    { Token(**cross**) }

++    { Token(**increment**) }

**id**

[A-Za-z0-9]

$I_0$ → [A-Za-z] → $I_1$

**cross**

$C_0$ → '+' → $C_1$

**increment**

$N_0$ → '+' → $N_1$ → '+' → $N_2$

**Char Stream**

| a | g | e | + | + | EOF |

**Accept, but…**    **shortest match!**

**eof**

$E_0$ → EOF → $E_1$

**Token Stream**

| id: a [0,1) | id: g [1,2) | id: e [2,3) | cross [3,4) | cross [4,5) | eof [5,6) |

# From RegExes to Tokenizer
## Algorithms

─── 1st Idea (flawed) ───

Consume char stream to **accept** state: return accepted token, restart DFAs with next char

─── Problem ───

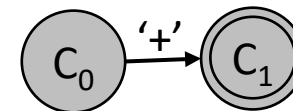What happens when token languages overlap / prefix each other?

**Language**

[A-Za-z][A-Za-z0-9]*   { Token(**id**) }

[A-Za-z0-9]

**id**                          [A-Za-z]  $I_1$

**Char Stream**

| a | | g |

We want the *immediate longest* (AKA "greedy") match

**Accept**
**but…**

$N_2$

**Token Stream**

| **id: a** [0,1) | **id: g** [1,2) | [2,3) | [3,4) | [4,5) | [5,6) |

$E_0$ → $E_1$

# From RegExes to Tokenizer
## Algorithms

Consume char stream to **accept** state: return accepted token, restart DFAs with next char

Consume char stream to **reject** states: return **last accepted** token, restart DFAs with that char

**Language**

[A-Za-z][A-Za-z0-9]*     { Token(**id**) }

**id**     [A-Za-z0-9]

[A-Za-z]

$I_1$

**Char Stream**

| a | g |
|---|---|

We want the
*immediate longest*
(AKA "greedy")
match

**Accept**
**but…**

Greed
is
Good

Let's modify the
Implementation to
support this!

$N_2$

**Token Stream**

| id: a | id: g |       |       |       |       |
|-------|-------|-------|-------|-------|-------|
| [0,1) | [1,2) | [2,3) | [3,4) | [4,5) | [5,6) |

$E_0$  $E_1$

# From RegExes to Tokenizer
## Algorithms

---1st Idea (flawed)---

Consume char stream to **accept** state: return accepted token, restart DFAs with next char

---NEW Idea (good)---

Consume char stream to **reject** states: return **last accepted** token, restart DFAs with that index

## Language

[A-Za-z][A-Za-z0-9]*    { Token(**id**) }

        +      { Token(**cross**) }
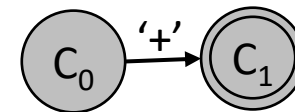
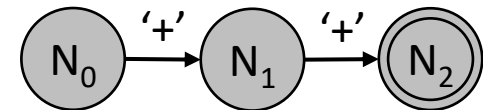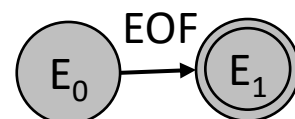       ++    { Token(**increment**) }

## Char Stream

| a | g | e | + | + | | EOF |

| 0 | 1 | 2 | 3 | 4 | 5 | 6 |

| $I_0$ $C_0$ $N_0$ $E_0$ | $I_1$ ✗ ✗ ✗ | $I_1$ ✗ ✗ ✗ | $I_1$ ✗ ✗ ✗ $\mid$ $I_0$ $C_0$ $N_0$ $E_0$ | ✗ ✗ ✗ ✗ | | |

**Last accept**

## Token Stream

**id: age**
[0,3)

**id**

$I_0$ →[A-Za-z]→ $I_1$ ⟲ [A-Za-z0-9]

**cross**

$C_0$ →'+'→ $C_1$

**increment**

$N_0$ →'+'→ $N_1$ →'+'→ $N_2$

**eof**

$E_0$ →EOF→ $E_1$

# From RegExes to Tokenizer

Algorithms

---
1st Idea (flawed)
---

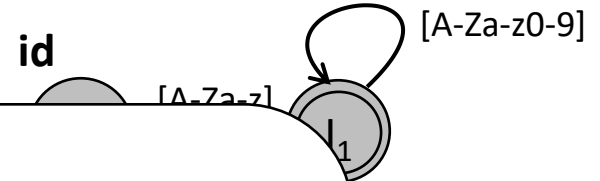Consume char stream to **accept** state: return accepted token, restart DFAs with next char

---
NEW Idea (good)
---

Consume char stream to **reject** states: return **last accepted** token, restart DFAs with that index

## Language

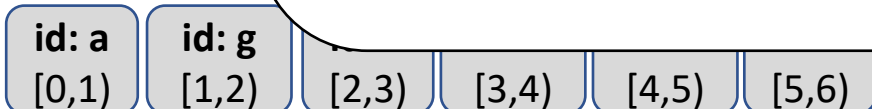[A-Za-z][A-Za-z0-9]*     { Token(**id**) }

$+$     { Token(**cross**) }

$++$     { Token(**increment**) }

**id**



**cross**



**increment**



## Char Stream



**Last accept**

**eof**



## Token Stream

| id: age [0,3) | increment [3,5) |

# From RegExes to Tokenizer
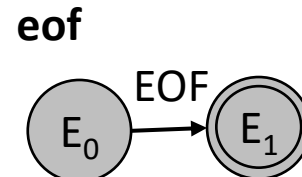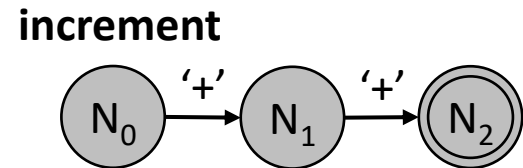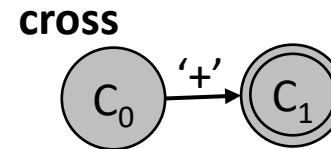Algorithms

---
1st Idea (flawed)
---

Consume char stream to **accept** state: return accepted token, restart DFAs with next char

---
NEW Idea (good)
---

Consume char stream to **reject** states: return **last accepted** token, restart DFAs with that index

## Language

[A-Za-z][A-Za-z0-9]*     { Token(**id**) }

$+$     { Token(**cross**) }

$++$     { Token(**increment**) }

## Char Stream



## Token Stream

| id: age | increment | eof |
|---------|-----------|-----|
| [0,3)   | [3,5)     | [5,6) |

# From RegExes to Tokenizer
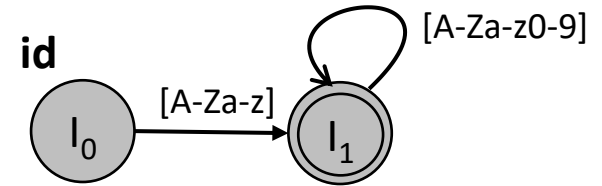## Algorithms

---

## 1st Idea (flawed)

Consume char stream to **accept** state: return accepted token, restart DFAs with next char
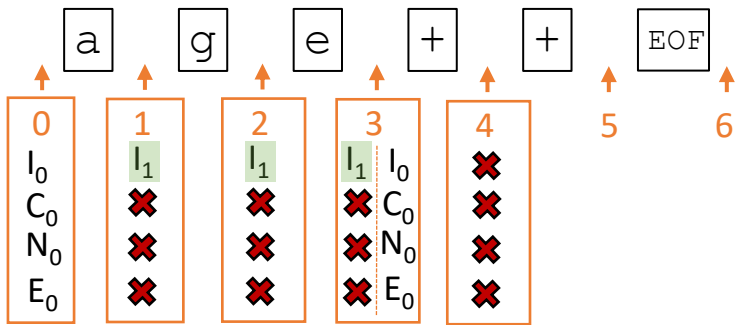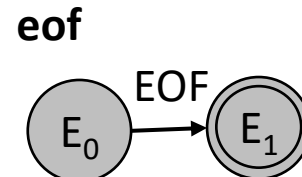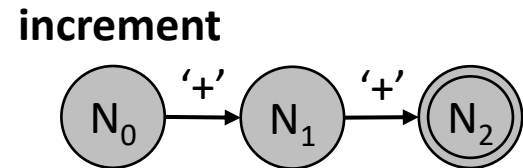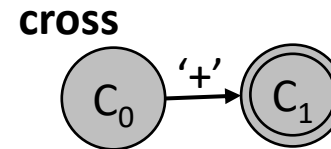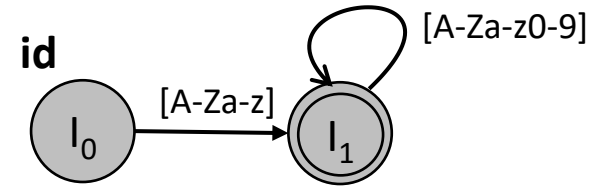
---

## NEW Idea (good)

Consume char stream to **reject** states: return **last accepted** token, restart DFAs with that index

### Language

[A-Za-z][A-Za-z0-9]*     { Token(**id**) }

       +     { Token(**cross**) }

     ++     { Token(**increment**) }

### Char Stream

| a | g | e | + | + | EOF |

Accept longest match!

### Token Stream

| id: age [0,3) | increment [3,5) | eof [5,6) |

**id**

$I_0$ $\xrightarrow{[A-Za-z]}$ $I_1$ with self-loop $[A-Za-z0-9]$

**cross**

$C_0$ $\xrightarrow{'+'}$ $C_1$

**increment**

$N_0$ $\xrightarrow{'+'}$ $N_1$ $\xrightarrow{'+'}$ $N_2$

**eof**

$E_0$ $\xrightarrow{EOF}$ $E_1$
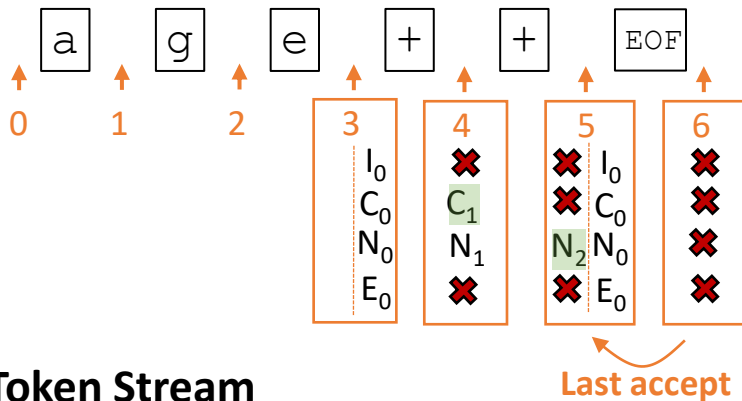
# Tokenizer Action Tables

Implementation

Consume char stream to **reject** states: return **last accepted** token, restart DFAs with that index

**Language**

[A-Za-z][A-Za-z0-9]*  { Token(**id**) }
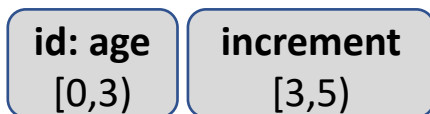
       +  { Token(**cross**) }

      ++  { Token(**increment**) }

| | + | letter | digit | EOF |
|---|---|---|---|---|
| $I_0$ | | $I_1$ | | |
| $I_1$ | | $I_1$ | $I_1$ | |
| $C_0$ | $C_1$ | | | |
| $C_1$ | | | | |
| $N_0$ | $N_1$ | | | |
| $N_1$ | $N_2$ | | | |
| $N_2$ | | | | |
| $E_0$ | | | | $E_1$ |
| $E_1$ | | | | |

| | Token |
|---|---|
| $I_0$ | |
| $I_1$ | **id** |
| $C_0$ | |
| $C_1$ | **cross** |
| $N_0$ | |
| $N_1$ | |
| $N_2$ | **increment** |
| $E_0$ | |
| $E_1$ | **eof** + accept |

**id**



**cross**



**increment**



**eof**

# Tokenizer Action Tables
## Implementation

| | + | letter | digit | EOF | | Token |
|---|---|---|---|---|---|---|
| $I_0$ | | $I_1$ | | | $I_0$ | |
| $I_1$ | | $I_1$ | $I_1$ | | $I_1$ | **id** |
| $C_0$ | $C_1$ | | | | $C_0$ | |
| $C_1$ | | | | | $C_1$ | **cross** |
| $N_0$ | $N_1$ | | | | $N_0$ | |
| $N_1$ | $N_2$ | | | | $N_1$ | |
| $N_2$ | | | | | $N_2$ | **increment** |
| $E_0$ | | | | $E_1$ | $E_0$ | |
| $E_1$ | | | | | $E_1$ | **eof** + accept |

**This basic machinery lets us implement a scanner from a RegEx spec!**

# Lexical Analysis Done
## Lecture 3 Preview

Source code in S
(sequence of chars)

|  | Lexical Definition | Lexical Recognition | Tokenization |
|---|---|---|---|
| Scanner<br>Lexical analysis | ✔ | ✔ | ✔ |

|  | Syntactic Definition | Syntactic Recognition | Parsing |
|---|---|---|---|
| Parser<br>Syntactic analysis | **Now** | **(Soon)** | **(Soon)** |

Semantic analysis

Intermediate code generation

IR optimization

Final Code generation

Final code optimization

Output code in T

**Naïve Approach**

"Just use a RegEx for all valid syntax"

# Regular Languages Lack Strength
## How Languages are Defined: CFGs

- Our RegEx-based scanner can emit a stream of tokens:

**Char Stream**

| X | \t | Y | = | Z | + | EOF |
|---|----|---|---|---|---|-----|



*Cute, but weak*

**Token Stream**

| ID: X | ID: Y | ASSIGN | ID: Z | PLUS | EOF |
|-------|-------|--------|-------|------|-----|

Observation: scanner ignores token order

**Audience Q: Could you enforce construct order another RegEx?**

Answer: Nope!   RegEx simply can't capture all PL structures (e.g. parentheses nesting)

# Defining Languages with Grammars
Syntactic Definition: How we use CFGs

- A set of (recursive) rewriting rules to rewrite sequence of symbols

- Any "completed" sequence represents a string in the language



I Has A Grammar

# Defining Languages with Grammars

Syntactic Definition: How we use CFGs

- A set of (recursive) rewriting rules to rewrite sequence of symbols

- Any "completed" sequence represents a string in the language

CFG = $(N, \Sigma, P, S)$ where:
- $N$: set of nonterminal symbols
- $\Sigma$: set of terminal symbols
- $P$: set of productions
- $S$: start nonterminal in $N$

*Rules where*
*LHS: a single nonterminal symbol*
*RHS: a sequence of any symbols*

# Defining Languages with Grammars

Syntactic Definition: How we use CFGs

Example:

$N = \{ A \}$

$\Sigma = \{ \textbf{(,)} \}$

$S = A$

$P = \left\{ \begin{array}{l} \text{P1: } A \rightarrow ( \, A \, ) \\ \text{P2: } A \rightarrow \boldsymbol{\varepsilon} \end{array} \right\}$

---

CFG = $(N,\Sigma,P,S)$ where:
- $N$: set of nonterminal symbols
- $\Sigma$: set of terminal symbols
- $P$: set of productions
- $S$: start nonterminal in $N$

# Defining Languages with Grammars
## Syntactic Definition: How we use CFGs

## Example:

$N = \{ A \}$

$\Sigma = \{ \textbf{(} , \textbf{)} \}$

$S = A$

$$P = \left\{ \begin{array}{l} \text{P1: } A \rightarrow ( \, A \, ) \\ \text{P2: } A \rightarrow \boldsymbol{\varepsilon} \end{array} \right\}$$

## Producing a string

$A$  — Begin with start symbol

$A \rightarrow ( \, A \, )$  — Apply production P1 (a *derivation step, denoted* ⇒)

$( \, A \, )$  — Get a new symbol string

$A \rightarrow ( \, A \, )$  — Apply production P1 again

$( \, ( \, A \, ) \, )$  — Get a new symbol string

$A \rightarrow \boldsymbol{\varepsilon}$  — Apply another production in P

$( \, ( \, ) \, )$  — Get a new symbol string

All terminals, this string is in language (a *sentence*)

# Simplifying Notation: Shorthand

Syntactic Definition: How we use CFGs

Example:

$N = \{ A \}$

$\Sigma = \{ \textbf{(,)} \}$

$S = A$

$$P = \left\{ \begin{array}{l} P1: A \rightarrow ( A ) \\ P2: A \rightarrow \boldsymbol{\varepsilon} \end{array} \right\}$$

Say N and Σ Implicit: Whatever symbols appears in productions

Say S Implicit: LHS of top production

Collapse rules with the same LHS using bar



TIME FOR A LITTLE SMACKEREL

of context-free grammar notation

# Simplifying Notation: Shorthand
## Syntactic Definition: How we use CFGs

Example:

$N = \{ A \}$

$\Sigma = \{ \textbf{(}, \textbf{)} \}$

Say N and $\Sigma$ Implicit: Whatever symbols appears in productions

$S = A$

Say S Implicit: LHS of top production

$P = \left\{ \begin{array}{l} \text{P1: } A \rightarrow \textbf{(} A \textbf{)} \\ \text{P2: } A \rightarrow \boldsymbol{\varepsilon} \end{array} \right\}$

Collapse rules with the same LHS using bar

Denote grammar as

$$A ::= \textbf{(} A \textbf{)}$$
$$A ::= \boldsymbol{\varepsilon}$$

or equivalently as

$$A ::= \textbf{(} A \textbf{)}$$
$$| \quad \boldsymbol{\varepsilon}$$

or even

$$A ::= \textbf{(} A \textbf{)} \mid \boldsymbol{\varepsilon}$$

# Simplifying Notation: Shorthand

Syntactic Definition: How we use CFGs

"BNF"

Denote grammar as

$$A ::= ( A )$$

$$A ::= \varepsilon$$

or equivalently as

$$A ::= ( A )$$
$$| \quad \varepsilon$$

or even

$$A ::= ( A ) | \varepsilon$$

# Some languages denoted in BNF
Syntactic Definition: How we use CFGs

$A ::= (A)$
$\quad | \quad \varepsilon$

---

$F ::= \mathbf{b}\; G\; \mathbf{y}\; \mathbf{e}$     $F \Rightarrow \mathbf{b}\; G\; \mathbf{y}\; \mathbf{e} \Rightarrow \mathbf{b}\; \mathbf{y}\; \mathbf{e}$

$\quad | \; \mathbf{see\; ya}$     $F \Rightarrow \mathbf{b}\; G\; \mathbf{y}\; \mathbf{e} \Rightarrow \mathbf{b}\; G\; \mathbf{y}\; \mathbf{y}\; \mathbf{e} \Rightarrow \mathbf{b}\; \mathbf{y}\; \mathbf{y}\; \mathbf{e}$

$G ::= G\mathbf{y}$     $F \Rightarrow \mathbf{see\; ya}$

$\quad | \quad \varepsilon$

---

$Y ::= \mathbf{a}\, Y$

$Z ::= \mathbf{w\; t\; f}$     $Y \Rightarrow a\; Y \Rightarrow a\; a\; Y \Rightarrow a\; a\; a\; Y \Rightarrow \ldots$

Accepts no strings
(not even the empty string)

# Parse Trees
Syntactic Definition: How we use CFGs

$F \Rightarrow$ **b** *G* **y e** $\Rightarrow$ **b** *G* **y y e** $\Rightarrow$ **b y y e**

## Represent Derivations

- Nodes are symbols in a tree

- Rooted at start symbol

- Children are derivation step

- Leaves are final string (if all nonterminals)

# CFG use in the Compiler
Syntactic Definition: How we use CFGs

# CFG use in the Compiler
Syntactic Definition: How we use CFGs

## CFG for PL Syntactic Structure

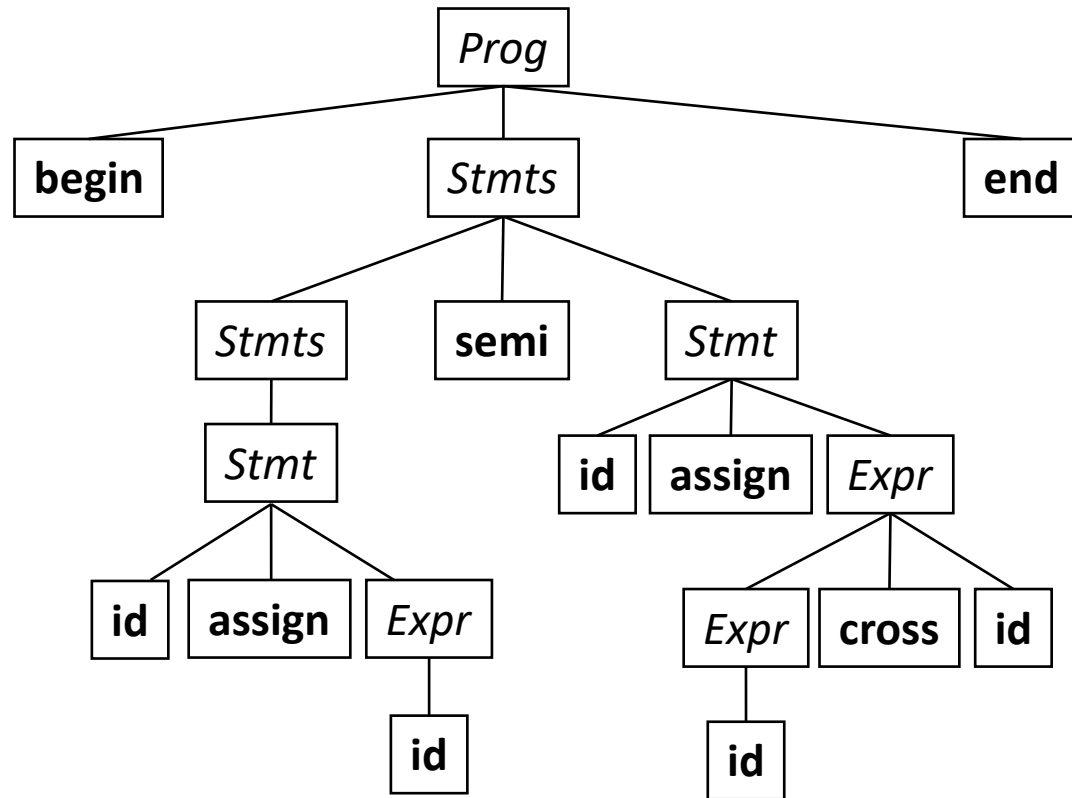Productions specify valid programs

- Let the terminals be the tokens in the language

- Let the nonterminals be the groupings that form language constructs
    - (loops, statements, functions, calls, etc)

- The grammar will recognize (or reject) the stream of tokens from the Lexer

Let's see an example with this grammar

### Productions
1. *Prog* ::= **begin** *Stmts* **end**
2. Stmts ::= *Stmts* **semi** *Stmt*
3.          | *Stmt*
4. *Stmt* ::= **id assign** *Expr*
5. *Expr* ::= **id**
6.          | *Expr* **cross id**

# Parse Tree



# Derivation Sequence

*Prog*
  **Prod. 1**

⇒ **begin** Stmts **end**
  **Prod. 2**

⇒ **begin** Stmts **semi** Stmt **end**
  **Prod. 3**

⇒ **begin** Stmt **semi** Stmt **end**
  **Prod. 4**

⇒ **begin id assign** *Expr* **semi** *Stmt* **end**
  **Prod. 4**

⇒ **begin** id assign *Expr* **semi** id assign *Expr* **end**
  **Prod. 5**

⇒ **begin** id assign **id** semi id assign *Expr* **end**
  **Prod. 6**

⇒ **begin** id assign id semi id assign *Expr* **cross id** end
  **Prod. 5**

⇒ **begin** id assign id semi id assign **id** cross id end

## Productions
1.  *Prog*  ::= **begin** *Stmts* **end**
2.  *Stmts* ::= *Stmts* **semi** *Stmt*
3.         | *Stmt*
4.  *Stmt*  ::= **id assign** *Expr*
5.  *Expr*  ::= **id**
6.         | *Expr* **cross id**

# End of Lecture
## Syntactic Definition

## Next Time

Parsing - Beyond specification for CFGs

- Extracting the *correct* tree from a token stream

## Time Permitting

- Proof sketch: why RegExs can't match PL constructs