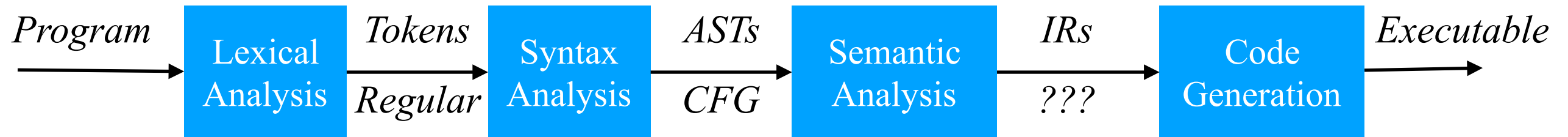


CS 160 Compilers

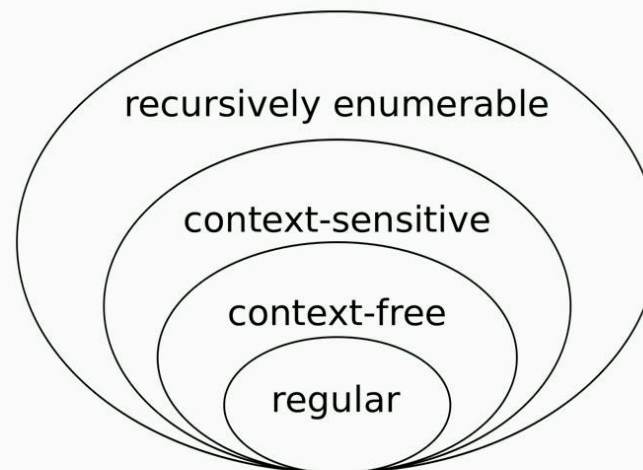
# Lecture 5: Lexical Analysis

Yu Feng  
Fall 2021

# A typical flow of a compiler

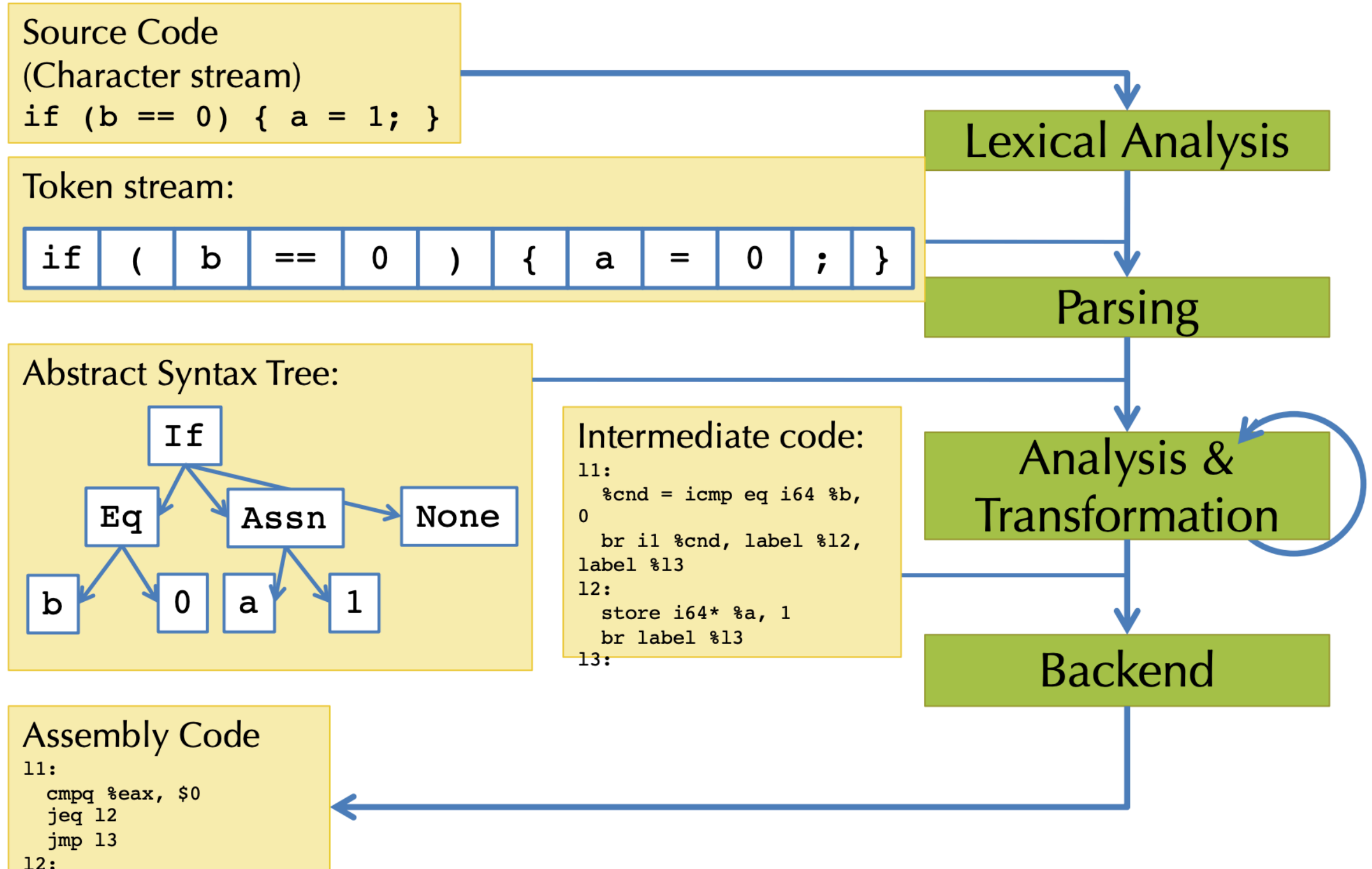


## Chomsky hierarchy



<https://en.wikipedia.org/wiki/File:Chomsky-hierarchy.svg>

# A typical flow of a compiler



# Lexical analysis

- Main Question: How to give structure to strings
- Analogy: Understanding an English sentence
  - First, we separate a string into words
  - Second, we understand sentence structure by diagramming the sentence
- Separating a string into words is called *lexing*
- Note that lexing is not necessarily trivial

# Lexical analysis

- Consider the following Patina program:

```
if x > y
```

```
then 10
```

```
else 8
```

- This program is just a string of characters

```
if x > y\nthen\t10\nelse\t8
```

- Goal: Portion the input string into substrings where the substrings are *tokens*

# What is a Token?

- Token is a syntactic category
- Example in English: noun, verbs, adjectives,...
- In a programming language: constants, identifiers, keywords, whitespaces...

# Tokens in Patina

- Tokens correspond to sets of strings
- Identifier: strings of letters, digits and '\_' starting with a letter
- Integer: a non-empty string of digits
- Keywords: “let”, “if”, ...
- Whitespace: a non-empty sequence of blanks, newlines, and tabs

# What are tokens for?

- Classify program substrings according to their role
- Output of lexical analysis is a stream of tokens...
- ...which is input to the parser
- Parser relies on token distinction
  - An identifier is treated different than a keyword



# Regular language/expressions

- We could specify tokens in many ways
- Regular Languages are the most popular
  - Simple and useful theory
  - Easy to understand
  - Efficient to implement

# Languages

- Definition: Let  $\Sigma$  be a set of characters, A **language over  $\Sigma$**  is a set of strings from characters drawn from  $\Sigma$
- Alphabet: English characters  $\Rightarrow$  Language: English sentences
- Languages are sets of strings
- Need some notation for specifying which sets we want
- The standard notation for regular languages is **regular expressions**

# Regular expressions

- Atomic Regular Expressions
  - Single character:  $c = \{“c”\}$
  - Epsilon:  $\varepsilon = \{“”\}$
- Compound Regular Expressions
  - Union:  $A+B = \{s \mid s \in A \text{ or } s \in B\}$
  - Concatenation:  $AB = \{ab \mid a \in A \text{ and } b \in B\}$
  - Iteration:  $A^* = \bigcup_{i \geq 0} A^i$  where  $A^i = A \dots i \text{ times } A$

# Regular expressions

- ▶ The **regular expressions** over  $\Sigma$  are the smallest set of expressions including
  - ▶  $\varepsilon$
  - ▶  $'c'$  where  $c \in \Sigma$
  - ▶  $A + B$  where  $A, B$  are regular expressions over  $\Sigma$
  - ▶  $AB$  where  $A, B$  are regular expressions over  $\Sigma$
  - ▶  $A^*$  where  $A$  is a regular expression over  $\Sigma$
- ▶ Regular expressions are simple, but **very useful**

# Example: Integers

- Integer: non-empty string of digits.
- $\text{digit} = '0' + '1' + '2' + '3' + '4' + '5' + '6' + \dots$
- $\text{integer} = \text{digit digit}^*$
- Abbreviation:  $A^+ = AA^*$

# Example: Identifier

- Identifier: strings of letters or digits, starting with a letter
- $\text{letter} = 'A' + \dots + 'Z' + 'a' + \dots + 'z' + '_'$
- $\text{identifier} = \text{letter} (\text{letter} + \text{digit})^*$
- How about  $(\text{letter}^* + \text{digit}^*)$ ?

# Example: Whitespace

- Whitespace: a non-empty sequence of blanks, newlines and tabs
- Whitespace = (' ' + '\n' + '\t')<sup>+</sup>

# Last example: email

- Consider UCSB cs emails: anyone@cs.ucsb.edu format
- $\Sigma = \text{letters} \cup \{., @\}$
- $\text{name} = \text{letter}^+$
- $\text{address} = \text{name } '@' \text{name } '.' \text{name } '.' \text{name}$



# TODOs by next lecture

- Come to the discussion session or office hour if you have questions
- Continue with your good work on HW1