# **Compiler Structure**

Martin Kellogg

## Today's Agenda

- Social contract of a compiler
- Compiler frontends
  - Lexing
  - Parsing
    - our first intermediate representation: abstract syntax trees

"Primum non nocere"

"Primum non nocere" - "First, do no harm."

"Primum non nocere" - "First, do no harm."

- one of the key tenets of bioethics
  - e.g., the Hippocratic Oath that doctors take

"Primum non nocere" - "First, do no harm."

- one of the key tenets of bioethics
  - e.g., the Hippocratic Oath that doctors take
- also a key tenet in compilers:

First <u>Rule of Compilers</u>: don't change semantics

"Primum non nocere" - "First, do no harm."

- one of the key tenets of bioethics
  - e.g., the Hippocratic Oath that doctors take
- also a key tenet in compilers:

First <u>Rule of Compilers</u>: don't change semantics

• i.e., don't change the meaning of the program

"Primum non nocere" - "First, do no harm."

- one of the key tenets of bioethics
  - e.g., the Hippocratic Oath that doctors take
- also a key tenet in compilers:

First <u>Rule of Compilers</u>: don't change semantics

- i.e., don't change the meaning of the program
- this is the **contract** that a compiler promises to its users

**Definition**: the *syntax* of a programming language is how programs are written down

• roughly "spelling and grammar"

**Definition**: the *syntax* of a programming language is how programs are written down

• roughly "spelling and grammar"

**Definition**: the *semantics* of a programming language explains what a given program **means** 

• Two relevant kinds of semantics (for now):

**Definition**: the *syntax* of a programming language is how programs are written down

• roughly "spelling and grammar"

**Definition**: the *semantics* of a programming language explains what a given program **means** 

- Two relevant kinds of semantics (for now):
  - *static semantics*: roughly, the state space of the program

**Definition**: the *syntax* of a programming language is how programs are written down

• roughly "spelling and grammar"

**Definition**: the *semantics* of a programming language explains what a given program **means** 

- Two relevant kinds of semantics (for now):
  - **static semantics**: roughly, the state space of the program
  - dynamic semantics: in some specific execution, what does this program actually evaluate to?

• A static semantics is useful for proving that certain behaviors do or do not occur

- A static semantics is useful for proving that certain behaviors do or do not occur
  - e.g., typechecking implies a particular static semantics

- A static semantics is useful for proving that certain behaviors do or do not occur
  - e.g., typechecking implies a particular static semantics
  - one way to think about this: the static semantics is the union of the program's dynamic semantics (for all possible inputs)

- A static semantics is useful for proving that certain behaviors do or do not occur
  - e.g., typechecking implies a particular static semantics
  - one way to think about this: the static semantics is the union of the program's dynamic semantics (for all possible inputs)
- We want some **dynamic** semantics to be impossible!

- A static semantics is useful for proving that certain behaviors do or do not occur
  - e.g., typechecking implies a particular static semantics
  - one way to think about this: the static semantics is the union of the program's dynamic semantics (for all possible inputs)
- We want some **dynamic** semantics to be impossible!

• e.g., 1 + "hello" should not have a dynamic semantics

- A static semantics is useful for proving that certain behaviors do or do not occur
  - e.g., typechecking implies a particular static semantics
  - one way to think about this: the static semantics is the union of the program's dynamic semantics (for all possible inputs)
- We want some **dynamic** semantics to be impossible!
  - e.g., 1 + "hello" should not have a dynamic semantics
    - typechecking prevents this "program" from ever being evaluated!

1. First, don't change the semantics

- **1.** First, don't change the semantics
  - from the textbook:

The compiler must preserve the meaning of the program being compiled.

- **1.** First, don't change the semantics
  - from the textbook:

This <u>R</u>ule is referring to the program's <u>static</u> semantics: we must not change the program's behavior for <u>any</u> input!

The compiler must preserve the meaning of the program being compiled.

- **1.** First, don't change the semantics
  - from the textbook:

This <u>R</u>ule is referring to the program's <u>static</u> semantics: we must not change the program's behavior for <u>any</u> input!

The compiler must preserve the meaning of the program being compiled.

2. Second, try to improve the program ("try to help")

- 1. First, don't change the semantics
  - from the textbook:

This <u>R</u>ule is referring to the program's <u>static</u> semantics: we must not change the program's behavior for <u>any</u> input!

The compiler must preserve the meaning of the program being compiled.

- 2. Second, try to improve the program ("try to help")
  - from the textbook:

The compiler must improve the input program in some discernible way.

- 1. First, don't change the semantics
  - from the textbook:

This <u>R</u>ule is referring to the program's <u>static</u> semantics: we must not change the program's behavior for <u>any</u> input!

7V.

The compiler must preserve the meaning of the program being compiled.

- 2. Second, try to improve th
  - from the textbook:

The compiler must improve the

If there is ever a **conflict** between these <u>R</u>ules, which one takes precedence?

- 1. First, don't change the semantics
  - from the textbook:

This <u>R</u>ule is referring to the program's <u>static</u> semantics: we must not change the program's behavior for <u>any</u> input!

7V.

The compiler must preserve the meaning of the program being compiled.

- 2. Second, try to improve th
  - from the textbook:

The compiler must improve the

If there is ever a **conflict** between these <u>R</u>ules, which one takes precedence? **The First <u>R</u>ule!** 

#### Traditional compiler/interpreter structure



#### Traditional compiler/interpreter structure



#### Traditional compiler/interpreter structure



## Today's Agenda

- Social contract of a compiler
- Compiler frontends
  - Lexing
  - Parsing
    - our first intermediate representation: abstract syntax trees

• A *lexical analyzer* (or *lexer*) divides program text into "tokens"

A lexical analyzer (or lexer) divides program text into "tokens"
 i.e., it is a function of type string -> token list

- A lexical analyzer (or lexer) divides program text into "tokens"
  i.e., it is a function of type string -> token list
- A token is a syntactic category

- A lexical analyzer (or lexer) divides program text into "tokens"
  i.e., it is a function of type string -> token list
- A token is a syntactic category
  - in English:
    - "noun", "verb", "adjective"

- A lexical analyzer (or lexer) divides program text into "tokens"
  - i.e., it is a function of type string -> token list
- A token is a syntactic category
  - in English:
    - "noun", "verb", "adjective"
  - in a programming language:
    - identifier, integer constant, keyword, whitespace
# Lexical analysis

- A *lexical analyzer* (or *lexer*) divides program text into "tokens"
  - i.e., it is a function of type string -> token list
- A token is a syntactic category
  - in English:
    - "noun", "verb", "adjective"
  - in a programming language:
    - identifier, integer constant, keyword, whitespace
- Parsers rely on token distinctions
  - e.g., identifiers are treated differently than keywords

• Note that tokens correspond to sets of strings

- Note that tokens correspond to sets of strings
  - e.g., an identifier token is any string of letters or digits that starts with a letter (in most languages)

- Note that tokens correspond to **sets** of strings
  - e.g., an identifier token is any string of letters or digits that starts with a letter (in most languages)
- A lexer needs to do two things:
  - recognize substrings that correspond to tokens

- Note that tokens correspond to **sets** of strings
  - e.g., an identifier token is any string of letters or digits that starts with a letter (in most languages)
- A lexer needs to do two things:
  - recognize substrings that correspond to tokens
  - return the *lexeme* of the token
    - the lexeme is the specific substring

- Note that tokens correspond to **sets** of strings
  - e.g., an identifier token is any string of letters or digits that starts with a letter (in most languages)
- A lexer needs to do two things:
  - recognize substrings that correspond to tokens
  - o return the *lexeme* of the token
    - the lexeme is the specific substring
    - e.g., in Cool, myVar is an identifier token whose lexeme is "myVar"

- Lexers usually discard "uninteresting" tokens that don't contribute to parsing
  - e.g., whitespace or comments

- Lexers usually discard "uninteresting" tokens that don't contribute to parsing
  - e.g., whitespace or comments
- Lexing a string is typically implemented by reading the string from start to end, **one token** at a time

- Lexers usually discard "uninteresting" tokens that don't contribute to parsing
  - e.g., whitespace or comments
- Lexing a string is typically implemented by reading the string from start to end, **one token** at a time
  - Lexers generally require some amount of *lookahead* to successfully partition a string
  - e.g., "i" vs "if"

- Lexers usually discard "uninteresting" tokens that don't contribute to parsing
  - e.g., whitespace or comments
- Lexing a string is typically implemented by reading the string from start to end, **one token** at a time
  - Lexers generally require some amount of *lookahead* to successfully partition a string
  - e.g., "i" vs "if"

looks like an identifier at first...

- Lexers usually discard "uninteresting" tokens that don't contribute to parsing
  - e.g., whitespace or comments
- Lexing a string is typically implemented by reading the string from start to end, **one token** at a time
  - Lexers generally require some amount of *lookahead* to successfully partition a string
  - e.g., "i" vs "if"

looks like an identifier at first...but it's actually a keyword

**Definition:** Let  $\Sigma$  ("sigma") be a set of characters. A *language over*  $\Sigma$  is a set of strings of characters drawn from  $\Sigma$ .  $\Sigma$  is called the *alphabet*.

**Definition:** Let  $\Sigma$  ("sigma") be a set of characters. A *language over*  $\Sigma$  is a set of strings of characters drawn from  $\Sigma$ .  $\Sigma$  is called the *alphabet*.

e.g., alphabet = English characters, language = valid English sentences

**Definition:** Let  $\Sigma$  ("sigma") be a set of characters. A *language over*  $\Sigma$  is a set of strings of characters drawn from  $\Sigma$ .  $\Sigma$  is called the *alphabet*.

- e.g., alphabet = English characters, language = valid English sentences
  - note that **not** every string of English characters is a valid sentence! E.g., "xagea' iwej"

**Definition:** Let  $\Sigma$  ("sigma") be a set of characters. A *language over*  $\Sigma$  is a set of strings of characters drawn from  $\Sigma$ .  $\Sigma$  is called the *alphabet*.

- e.g., alphabet = English characters, language = valid English sentences
  - note that **not** every string of English characters is a valid sentence! E.g., "xagea' iwej"
- for lexing, we care about the *regular languages*, which you might remember from your CS theory class

**Definition:** a *regular language* is a language that can be recognized by a *regular expression* or, equivalently, by a *deterministic finite automaton*.

**Definition:** a *regular language* is a language that can be recognized by a *regular expression* or, equivalently, by a *deterministic finite automaton*.

- More formally, the set of regular languages over some alphabet Σ is defined as:
  - the empty language Ø is regular
  - for each s in  $\Sigma$ , the singleton set { s } is regular
  - $\circ$  if A is a regular language, then A<sup>\*</sup> is a regular language
  - if A and B are regular languages, then A  $\cup$  B and A B are regular

**Definition:** a *regular language* is a language that can *regular expression* or, equivalently, by a *deterministic* 

- More formally, the set of regular languages ove defined as:
  - $\circ$  the empty language Ø is regular
  - for each s in  $\Sigma$ , the singleton set { s } is regular
  - if A is a regular language, then  $A^*$  is a regular language
  - if A and B are regular languages, then  $A \cup B$  and  $A \bullet B$  are regular

This Kleene Star operator indicates repetition ("zero or more times")

**Definition:** a *regular language* is a language that can be recognized by a *regular expression* or, equivalently, by a *deterministic finite automaton*.

- More formally, the set of regular languages over some alphabet Σ is defined as:
  - the empty language strings in both A
  - for each s in  $\Sigma$ , the si and  $B^{"}$ ) regular
  - if A is a regular language, then A is a regular language
  - if A and B are regular languages, then  $A \cup B$  and  $A \bullet B$  are regular

Definition: a regular language is a language that can be recognized by a *regular expression* or, equivalently, by a *deterministic finite automaton*.

- More formally, the set of regular languages over some alphabet  $\Sigma$  is defined as:
  - the empty language  $\left(\begin{array}{c} \text{Concatenation} (\text{``each string in A} \\ \text{for each s in } \Sigma, \text{ the si} \right)$ Ο
  - Ο
  - if A is a regular language, then A' is a regular language Ο
  - if A and B are regular languages, then  $A \cup B$  and  $A \bullet B$  are Ο regular

• In practice, we use a *lexical analysis generator* like ply, flex, or ocamllex to build lexers for us

- In practice, we use a *lexical analysis generator* like ply, flex, or ocamllex to build lexers for us
  - these tools take regular expressions as input...

- In practice, we use a *lexical analysis generator* like ply, flex, or ocamllex to build lexers for us
  - these tools take regular expressions as input...
  - …and then convert into finite automata, which they implement using lookup tables…

- In practice, we use a *lexical analysis generator* like ply, flex, or ocamllex to build lexers for us
  - these tools take regular expressions as input...
  - …and then convert into finite automata, which they implement using lookup tables…
  - to eventually generate performant lexing code

#### **Trivia Break: Computing History**

Around 825 CE, this Persian scientist and polymath wrote kitāb alhisāb al-hindī ("Book of Indian computation") and kitab al-jam' wa'l-tafrig al-hisāb al-hindī ("Addition and subtraction in Indian arithmetic"). Despite these works later having a major influence on computer science, he is best known for his popularizing treatise on algebra (the "Al-Jabr"). The word "algorithm" is derived from his name.

# Today's Agenda

- Social contract of a compiler
- Compiler frontends
  - Lexing
  - Parsing
    - our first intermediate representation: abstract syntax trees

- Otherwise, it produces an error
  - e.g., "parse error on line 3"

- Otherwise, it produces an error
  - e.g., "parse error on line 3"
- Comparison of lexing vs. parsing:

Phase	Input	Output
Lexer		
Darcar		
Parser		

- Otherwise, it produces an error
  - e.g., "parse error on line 3"
- Comparison of lexing vs. parsing:

Phase	Input	Output
Lexer	Sequence of characters	Sequence of tokens
Parser		

- Otherwise, it produces an error
  - e.g., "parse error on line 3"
- Comparison of lexing vs. parsing:

Phase	Input	Output
Lexer	Sequence of characters	Sequence of tokens
Parser	Sequence of tokens	Parse tree

# Parser Output: Abstract Syntax Trees (ASTs)

• Implicit in our discussion of compilation so far: programs in the compiler's source language are data, from the perspective of the compiler

# Parser Output: Abstract Syntax Trees (ASTs)

- Implicit in our discussion of compilation so far: programs in the compiler's source language are data, from the perspective of the compiler
  - this is a powerful perspective: you know how to write programs that manipulate data already

# Parser Output: Abstract Syntax Trees (ASTs)

- Implicit in our discussion of compilation so far: programs in the compiler's source language are data, from the perspective of the compiler
  - this is a powerful perspective: you know how to write programs that manipulate data already
  - fundamentally, your compiler is just another program that manipulates data - the only difference is that the data, in this case, are also programs
## Parser Output: Abstract Syntax Trees (ASTs)

- Implicit in our discussion of compilation so far: programs in the compiler's source language are data, from the perspective of the compiler
  - this is a powerful perspective: you know how to write programs that manipulate data already
  - fundamentally, your compiler is just another program that manipulates data - the only difference is that the data, in this case, are also programs
- So, how do we represent programs as data?

# Parser Output: Abstract Syntax Trees (ASTs)

- Implicit in our discussion of compilation so far: programs in the compiler's source language are data, from the perspective of the compiler
  - this is a powerful perspective: you know how to write programs that manipulate data already
  - fundamentally, your compiler is just another program that manipulates data - the only difference is that the data, in this case, are also programs
- So, how do we represent programs as data?
  - there are many ways, but today we'll discuss two...

#1: treat the program as a string

• allows us to easily decide **syntactic** properties

- allows us to easily decide syntactic properties
  - for example, checking if a program contains the text "foo"

- allows us to easily decide syntactic properties
  - for example, checking if a program contains the text "foo"
- key downside: cannot use the program's semantics

- allows us to easily decide syntactic properties
  - o for example, checking if a program contains the text "foo"
- key downside: cannot use the program's semantics
  - semantics are relevant for properties related to context that is, where the question to be decided depends on the rest of the program

#2: treat the program as a tree

#2: treat the program as a tree

**Definition**: an *abstract syntax tree* (or *AST*) is a tree-based representation of a program's syntactic structure

#2: treat the program as a tree

**Definition**: an *abstract syntax tree* (or *AST*) is a tree-based representation of a program's syntactic structure

• usually produced by a parser

#2: treat the program as a tree

**Definition**: an *abstract syntax tree* (or *AST*) is a tree-based representation of a program's syntactic structure

- usually produced by a parser
- nodes in the tree represent syntactic constructs

#2: treat the program as a tree

**Definition**: an *abstract syntax tree* (or *AST*) is a tree-based representation of a program's syntactic structure

- usually produced by a parser
- nodes in the tree represent syntactic constructs
  - parent-child relationships in the AST represent compound expressions in the source code (e.g., a "plus node" might have two children: the left and right side expressions)

Example: 5 + (2 + 3)



Example: 5 + (2 + 3)







- Not all sequences of tokens are programs and correspond to an AST. E.g.:
  - then x \* / + 3 while x ; y z then

- Not all sequences of tokens are programs and correspond to an AST. E.g.:
  - then x \* / + 3 while x ; y z then
- The parser must distinguish between valid and invalid sequences of tokens

- Not all sequences of tokens are programs and correspond to an AST. E.g.:
  - $\circ$  then x \* / + 3 while x ; y z then
- The parser must distinguish between valid and invalid sequences of tokens
- We need:
  - A language or formalism to describe valid sequences of tokens

- Not all sequences of tokens are programs and correspond to an AST. E.g.:
  - $\circ$  then x \* / + 3 while x ; y z then
- The parser must distinguish between valid and invalid sequences of tokens
- We need:
  - A language or formalism to describe valid sequences of tokens
  - A method (i.e., an algorithm) for distinguishing between valid and invalid token sequences

### Aside: computability theory

**Definition**: an *algorithm* is a finite sequence of mathematically rigorous instructions, typically used to solve a class of specific problems or to perform a computation. [Wikipedia's <u>Algorithm</u>]

## Aside: computability theory

**Definition**: an *algorithm* is a finite sequence of mathematically rigorous instructions, typically used to solve a class of specific problems or to perform a computation. [Wikipedia's <u>Algorithm</u>]

• when we talk about algorithms in Computer Science, we generally are referring to the computable functions, which are the formalized analogue of the intuitive notion above

## Aside: computability theory

**Definition**: an *algorithm* is a finite sequence of mathematically rigorous instructions, typically used to solve a class of specific problems or to perform a computation. [Wikipedia's <u>Algorithm</u>]

- when we talk about algorithms in Computer Science, we generally are referring to the computable functions, which are the formalized analogue of the intuitive notion above
  - a function is *computable* if there exists an algorithm that can do the job of the function which we can execute in one of our models of computation (e.g., the lambda calculus or a Turing machine)

• So, what's the **formalism** for parsing?

- So, what's the **formalism** for parsing?
  - Regular expressions/DFAs (as we used for lexing) don't work

- So, what's the **formalism** for parsing?
  - Regular expressions/DFAs (as we used for lexing) don't work
    - because parsing requires context!

- So, what's the **formalism** for parsing?
  - Regular expressions/DFAs (as we used for lexing) don't work
    - because parsing requires context!
    - e.g., we need to check if parentheses are matched

- So, what's the **formalism** for parsing?
  - Regular expressions/DFAs (as we used for lexing) don't work
    - because parsing requires context!
    - e.g., we need to check if parentheses are matched
- Perhaps you recall from your CS theory class what kind of formalism is appropriate for checking the language of balanced parentheses?

- So, what's the **formalism** for parsing?
  - Regular expressions/DFAs (as we used for lexing) don't work
    - because parsing requires context!
    - e.g., we need to check if parentheses are matched
- Perhaps you recall from your CS theory class what kind of formalism is appropriate for checking the language of balanced parentheses?
  - the context-free grammars or pushdown automata

How do we know? Programming languages have recursive structure

- How do we know? Programming languages have recursive structure
- E.g., consider the language of arithmetic expressions with integers, +, \*, and ( )

- How do we know? Programming languages have recursive structure
- E.g., consider the language of arithmetic expressions with integers, +, \*, and ( )
- An expression in this language is either:
  - $\circ$  an integer
  - an **expression** followed by "+" followed by an **expression**
  - an expression followed by "\*" followed by an expression
  - a '(' followed by an **expression** followed by ')'

- How do we know? Programming languages have recursive structure
- E.g., consider the language of arithmetic expressions with integers, +, \*, and ( )
- An expression in this language is either:
  - $\circ$  an integer
  - an **expression** followed by "+" followed by an **expression**
  - an **expression** followed by "\*" followed by an **expression**
  - a '(' followed by an **expression** followed by ')'
- 5, 5+3, (5+3) \*7 are expressions in this language

- An alternative notation for the language on the last slide:
  - **E** -> int
  - **E -> E + E**
  - **E -> E \* E**
  - **E** -> ( **E** )

- An alternative notation for the language on the last slide:
  - **E -> int**
  - **E -> E + E**
  - **E -> E \* E**
  - **E** -> ( **E** )
- We can view these rules as **rewrite rules**

- An alternative notation for the language on the last slide:
  - **E** -> int
  - **E -> E + E**
  - **E -> E \* E**
  - **E** -> ( **E** )
- We can view these rules as **rewrite rules** 
  - We start with **E** and replace occurrences of **E** with some right-hand side
- An alternative notation for the language on the last slide:
  - **E** -> int
  - **E -> E + E**
  - **E -> E \* E**
  - **E** -> ( **E** )
- We can view these rules as **rewrite rules** 
  - We start with **E** and replace occurrences of **E** with some right-hand side
  - Eventually, we'll derive an actual program in the language

- An alternative notation for the language on the last slide:
  - **E** -> int
  - **E -> E + E**
  - **E -> E \* E**
  - **E**->(**E**)

Note that there is no series of rewrites that allows us to obtain an **invalid** program (e.g., with mismatched parentheses). Why?

- We can view these rules as rewrite rules
  - We start with E and replace occurrences of E with some right-hand side
  - Eventually, we'll derive an actual program in the language

• The notation on the previous slide is exactly a *context-free grammar*.

- The notation on the previous slide is exactly a *context-free grammar*.
- Officially, a context-free grammar consists of:

- The notation on the previous slide is exactly a *context-free grammar*.
- Officially, a context-free grammar consists of:
  - $\circ$  A set of *non-terminals* N
    - Written in uppercase in these notes

- The notation on the previous slide is exactly a *context-free grammar*.
- Officially, a context-free grammar consists of:
  - A set of *non-terminals* N
    - Written in uppercase in these notes
  - A set of *terminals* T
    - Lowercase or punctuation in these notes

- The notation on the previous slide is exactly a *context-free grammar*.
- Officially, a context-free grammar consists of:
  - A set of *non-terminals* N
    - Written in uppercase in these notes
  - A set of *terminals* T
    - Lowercase or punctuation in these notes
  - A *start symbol* S (a non-terminal)

- The notation on the previous slide is exactly a *context-free grammar*.
- Officially, a context-free grammar consists of:
  - A set of *non-terminals* N
    - Written in uppercase in these notes
  - A set of *terminals* T
    - Lowercase or punctuation in these notes
  - A *start symbol* S (a non-terminal)
  - A set of *productions* (rewrite rules)

- Context-free grammar example:
  - **E** -> int
  - **E -> E + E**
  - **E -> E \* E**
  - **E** -> ( **E** )

- Context-free grammar example:
  - **E** -> int
  - **E -> E + E**
  - **E -> E \* E**
  - **E**->(**E**)
- **E** is the only **non-terminal**

- Context-free grammar example:
  - **E** -> int
  - **E -> E + E**
  - **E -> E \* E**
  - **E**->(**E**)
- **E** is the only **non-terminal**
- int, +, \*, (, and ) are the *terminals* 
  - called terminals because they are never replaced

- Context-free grammar example:
  - **E** -> int
  - **E -> E + E**
  - **E -> E \* E**
  - **E**->(**E**)
- **E** is the only **non-terminal**
- int, +, \*, (, and ) are the *terminals* 
  - called terminals because they are never replaced
- each bullet above is one of the rewrite rules (*productions*)

- Context-free grammar example:
  - **E** -> int
  - **E -> E + E**
  - **E -> E \* E**
  - **E**->(**E**)
- **E** is the only **non-terminal**
- int, +, \*, (, and ) are the *terminals* 
  - called terminals because they are never replaced
- each bullet above is one of the rewrite rules (*productions*)
- by convention, the first non-terminal in the first production is the *start symbol*

• A *derivation* is a sequence of productions, starting from the start symbol of the context-free grammar, that produces a string of only terminals

- A *derivation* is a sequence of productions, starting from the start symbol of the context-free grammar, that produces a string of only terminals
  - a derivation can be drawn as a tree

- A *derivation* is a sequence of productions, starting from the start symbol of the context-free grammar, that produces a string of only terminals
  - a derivation can be drawn as a tree
    - start symbol is the root

- A *derivation* is a sequence of productions, starting from the start symbol of the context-free grammar, that produces a string of only terminals
  - a derivation can be drawn as a tree
    - start symbol is the root
    - for each production, add children for each symbol on the right-hand side of the rule as children of the left-hand side non-terminal

- A *derivation* is a sequence of productions, starting from the start symbol of the context-free grammar, that produces a string of only terminals
  - a derivation can be drawn as a tree
    - start symbol is the root
    - for each production, add children for each symbol on the right-hand side of the rule as children of the left-hand side non-terminal
  - this derivation exactly corresponds to the AST!

• In the early days of computing, parsing was the **most expensive** part of a compiler

- In the early days of computing, parsing was the **most expensive** part of a compiler
  - e.g., consider original FORTRAN:
    - no semantic analysis or typechecking
    - little optimization
    - code generation is *syntax-directed*

- In the early days of computing, parsing was the **most expensive** part of a compiler
  - e.g., consider original FORTRAN:
    - no semantic analysis or typechecking
    - little optimization
    - code generation is syntax-directed
- Therefore, a lot of early compilers research was focused on improving parsing

- In the early days of computing, parsing was the **most expensive** part of a compiler
  - e.g., consider original FORTRAN:
    - no semantic analysis or typechecking
    - little optimization
    - code generation is syntax-directed
- Therefore, a lot of early compilers research was focused on improving parsing
  - e.g., LL(1) vs LR grammars, Earley parsing, etc.
    we aren't going to cover these in this class!

• These days, parsing is considered a mostly-solved problem

- These days, parsing is considered a mostly-solved problem
  - semantic analysis and optimization times dominate

- These days, parsing is considered a mostly-solved problem
  - semantic analysis and optimization times dominate
  - parsing is fast enough that it's not worth optimizing further

- These days, parsing is considered a mostly-solved problem
  - o semantic analysis and optimization times dominate
  - parsing is fast enough that it's not worth optimizing further
- Automatic *parser generators* (e.g., bison or ocamlyacc) can produce a parser implementation from a context-free grammar

- These days, parsing is considered a mostly-solved problem
  - semantic analysis and optimization times dominate
  - parsing is fast enough that it's not worth optimizing further
- Automatic parser generators (e.g., bison or ocamlyacc) can produce a parser implementation from a context-free grammar
  - production compilers, though, typically use a custom recursive descent parser

- These days, parsing is considered a mostly-solved problem
  - o semantic analysis and optimization times dominate
  - parsing is fast enough that it's not worth optimizing further
- Automatic *parser generators* (e.g., **bison** or **ocamlyacc**) can produce a parser implementation from a context-free grammar
  - production compilers, though, typically use a custom recursive descent parser
    - that is, compiler engineers prioritize simplicity over having a fast parser (the speed gains are no longer worthwhile for the technical debt)

- These days, parsing is considered a **mostly-solved problem** 
  - semantic analysis and optimization times de Takeaway for you:
  - parsing is fast enough that it's not worth op you can assume
- Automatic *parser generators* (e.g., **bison** or **oca** that it's **relatively** produce a parser implementation from a contex **easy** to get an AST
  - production compilers, though, typically use a custom recursive descent parser
    - that is, compiler engineers prioritize simplicity over having a fast parser (the speed gains are no longer worthwhile for the technical debt)

• PA1 (all four languages!) due today

- PA1 (all four languages!) due today
- My OH this week are modified:
  - I will hold OH today instead of Wednesday, 4-5pm
  - if you need to meet later in the week, send me an email and we'll arrange something

- PA1 (all four languages!) due today
- My OH this week are modified:
  - I will hold OH today instead of Wednesday, 4-5pm
  - if you need to meet later in the week, send me an email and we'll arrange something
- Don't forget: PA2c1 is due **Friday** 
  - this is a testing assignment: you'll just write Cool programs

- PA1 (all four languages!) due today
- My OH this week are modified:
  - I will hold OH today instead of Wednesday, 4-5pm
  - if you need to meet later in the week, send me an email and we'll arrange something
- Don't forget: PA2c1 is due **Friday** 
  - this is a testing assignment: you'll just write Cool programs
- The course staff has become aware of a bug in the Apple Silicon builds for Cool related to newlines
  - No fix is available. Use Ubuntu instead.