# Compiler concepts: Parsing
## COMSM1302 Overview of Computer Architecture

John Lapinskas, University of Bristol

# Describing languages
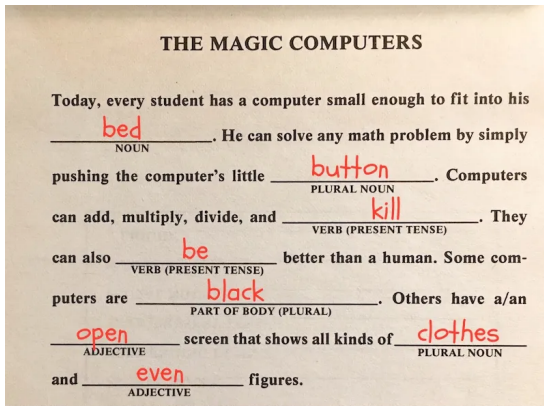
We have a description of Hack syntax in Nisan and Schocken, right?

# Describing languages

We have a description of Hack syntax in Nisan and Schocken, right?



Maybe this isn't the most convenient form possible...

# A more sophisticated approach

We'll need a highly sophisticated mathematical construction:
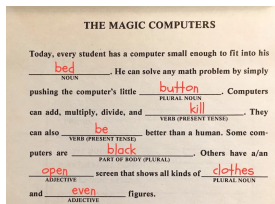
# A more sophisticated approach: Madlibs!

We'll need a highly sophisticated mathematical construction: Madlibs.



Source: Jesse Vig via Medium (here)

# A more sophisticated approach: Madlibs?

We'll need a highly sophisticated mathematical construction: Madlibs.



Source: Jesse Vig via Medium (here)

A **context-free grammar** (or just **grammar**) is a way of quickly and rigorously specifying which strings in a language have valid syntax.

There is a deep and rich mathematical theory here, which we thankfully don't need to learn! Programmers express grammars in **Backus-Naur Form** (**BNF**), and usually just understanding BNF is enough.

BNF is basically Madlibs, but recursive.

# Introduction to ~~Madlibs~~ Backus-Naur Form (BNF)

Here's a simple example:

$$\langle noun \rangle ::= \text{'lecturer'} \mid \text{'student'} \mid \text{'pizza'}$$
$$\langle presentVerb \rangle ::= \text{'eats'} \mid \text{'devours'} \mid \text{'consumes'}$$

You can read each | as "or" and each ::= as "is defined as". E.g. a $\langle noun \rangle$ is defined as one of the three strings 'lecturer', 'student', or 'pizza'.

# Introduction to ~~Madlibs~~ Backus-Naur Form (BNF)

Here's a simple example:

$$\langle noun \rangle ::= \text{'lecturer'} \mid \text{'student'} \mid \text{'pizza'}$$
$$\langle presentVerb \rangle ::= \text{'eats'} \mid \text{'devours'} \mid \text{'consumes'}$$

You can read each | as "or" and each ::= as "is defined as". E.g. a $\langle noun \rangle$ is defined as one of the three strings 'lecturer', 'student', or 'pizza'.

We can also build up definitions in terms of other definitions, e.g.

$$\langle sentence \rangle ::= \text{'The'} \; \langle noun \rangle \; \langle presentVerb \rangle \; \text{'the'} \; \langle noun \rangle$$

Here, valid $\langle sentence \rangle$s include:

# Introduction to ~~Madlibs~~ Backus-Naur Form (BNF)

Here's a simple example:

$$\langle noun \rangle ::= \text{`lecturer'} \mid \text{`student'} \mid \text{`pizza'}$$
$$\langle presentVerb \rangle ::= \text{`eats'} \mid \text{`devours'} \mid \text{`consumes'}$$

You can read each | as "or" and each ::= as "is defined as". E.g. a $\langle noun \rangle$ is defined as one of the three strings 'lecturer', 'student', or 'pizza'.

We can also build up definitions in terms of other definitions, e.g.

$$\langle sentence \rangle ::= \text{`The'} \; \langle noun \rangle \; \langle presentVerb \rangle \; \text{`the'} \; \langle noun \rangle$$

Here, valid $\langle sentence \rangle$s include:

$$\text{`The'} \; \langle noun \rangle \; \langle presentVerb \rangle \; \text{`the'} \; \langle noun \rangle$$

# Introduction to ~~Madlibs~~ Backus-Naur Form (BNF)

Here's a simple example:

$$\langle\text{noun}\rangle ::= \text{'lecturer'} \mid \text{'student'} \mid \text{'pizza'}$$
$$\langle\text{presentVerb}\rangle ::= \text{'eats'} \mid \text{'devours'} \mid \text{'consumes'}$$

You can read each | as "or" and each ::= as "is defined as". E.g. a $\langle\text{noun}\rangle$ is defined as one of the three strings 'lecturer', 'student', or 'pizza'.

We can also build up definitions in terms of other definitions, e.g.

$$\langle\text{sentence}\rangle ::= \text{'The'} \; \langle\text{noun}\rangle \; \langle\text{presentVerb}\rangle \; \text{'the'} \; \langle\text{noun}\rangle$$

Here, valid $\langle\text{sentence}\rangle$s include:

<p style="text-align:center">'The' <strong style="color:red">'lecturer' 'consumes'</strong> 'the' <strong style="color:red">'pizza'</strong></p>

# Introduction to ~~Madlibs~~ Backus-Naur Form (BNF)

Here's a simple example:

$$\langle noun \rangle ::= \text{'lecturer'} \mid \text{'student'} \mid \text{'pizza'}$$
$$\langle presentVerb \rangle ::= \text{'eats'} \mid \text{'devours'} \mid \text{'consumes'}$$

You can read each | as "or" and each ::= as "is defined as". E.g. a $\langle noun \rangle$ is defined as one of the three strings 'lecturer', 'student', or 'pizza'.

We can also build up definitions in terms of other definitions, e.g.

$$\langle sentence \rangle ::= \text{'The'} \langle noun \rangle \langle presentVerb \rangle \text{'the'} \langle noun \rangle$$

Here, valid $\langle sentence \rangle$s include:

'The' **'student' 'eats'** 'the' **'pizza'**

# Introduction to ~~Madlibs~~ Backus-Naur Form (BNF)

Here's a simple example:

$$\langle noun \rangle ::= \text{'lecturer'} \mid \text{'student'} \mid \text{'pizza'}$$
$$\langle presentVerb \rangle ::= \text{'eats'} \mid \text{'devours'} \mid \text{'consumes'}$$

You can read each | as "or" and each ::= as "is defined as". E.g. a $\langle noun \rangle$ is defined as one of the three strings 'lecturer', 'student', or 'pizza'.

We can also build up definitions in terms of other definitions, e.g.

$$\langle sentence \rangle ::= \text{'The'} \; \langle noun \rangle \; \langle presentVerb \rangle \; \text{'the'} \; \langle noun \rangle$$

Here, valid $\langle sentence \rangle$s include:

<span style="color:black">'The'</span> **'lecturer' 'devours'** <span style="color:black">'the'</span> **'student'**

Here's a simple example:

$$\langle\text{noun}\rangle ::= \text{'lecturer'} \mid \text{'student'} \mid \text{'pizza'}$$
$$\langle\text{presentVerb}\rangle ::= \text{'eats'} \mid \text{'devours'} \mid \text{'consumes'}$$

You can read each | as "or" and each ::= as "is defined as". E.g. a $\langle\text{noun}\rangle$ is defined as one of the three strings 'lecturer', 'student', or 'pizza'.

We can also build up definitions in terms of other definitions, e.g.

$$\langle\text{sentence}\rangle ::= \text{'The'} \; \langle\text{noun}\rangle \; \langle\text{presentVerb}\rangle \; \text{'the'} \; \langle\text{noun}\rangle$$

Here, valid $\langle\text{sentence}\rangle$s include:

<blockquote>'The' **'lecturer' 'devours'** 'the' **'student'**</blockquote>

Anything we define as part of the grammar must be enclosed in $\langle\rangle$s. We call these **non-terminal symbols**. Anything else (e.g. 'lecturer') is a **terminal symbol** or **token**.

## Example: Integers

The last feature of BNF is the source of its power: it allows **recursion**.

For example, suppose our tokens are '0' through '9', and we want to define a non-terminal symbol that matches precisely the non-negative whole numbers (allowing for leading zeroes). We could write:

## Example: Integers

The last feature of BNF is the source of its power: it allows **recursion**.

For example, suppose our tokens are '0' through '9', and we want to define a non-terminal symbol that matches precisely the non-negative whole numbers (allowing for leading zeroes). We could write:

$$\langle \text{digit} \rangle ::= \text{'0'} \mid \text{'1'} \mid \text{'2'} \mid \text{'3'} \mid \text{'4'} \mid \text{'5'} \mid \text{'6'} \mid \text{'7'} \mid \text{'8'} \mid \text{'9'}$$
$$\langle \text{number} \rangle ::= \langle \text{digit} \rangle \mid \langle \text{digit} \rangle \, \langle \text{number} \rangle$$

E.g. '016' is a $\langle \text{number} \rangle$ because we can expand $\langle \text{number} \rangle$'s definition as:

$$\langle \text{number} \rangle \longrightarrow \langle \text{digit} \rangle \, \langle \text{number} \rangle \longrightarrow \langle \text{digit} \rangle \, \langle \text{digit} \rangle \, \langle \text{number} \rangle$$
$$\longrightarrow \langle \text{digit} \rangle \, \langle \text{digit} \rangle \, \langle \text{digit} \rangle \longrightarrow \text{'0'} \, \text{'1'} \, \text{'6'}.$$

## Example: Integers

The last feature of BNF is the source of its power: it allows **recursion**.

For example, suppose our tokens are '0' through '9', and we want to define a non-terminal symbol that matches precisely the non-negative whole numbers (allowing for leading zeroes). We could write:

$$\langle \text{digit} \rangle ::= \text{'0'} \mid \text{'1'} \mid \text{'2'} \mid \text{'3'} \mid \text{'4'} \mid \text{'5'} \mid \text{'6'} \mid \text{'7'} \mid \text{'8'} \mid \text{'9'}$$
$$\langle \text{number} \rangle ::= \langle \text{digit} \rangle \mid \langle \text{digit} \rangle \; \langle \text{number} \rangle$$

E.g. '016' is a $\langle \text{number} \rangle$ because we can expand $\langle \text{number} \rangle$'s definition as:

$$\langle \text{number} \rangle \longrightarrow \langle \text{digit} \rangle \; \langle \text{number} \rangle \longrightarrow \langle \text{digit} \rangle \; \langle \text{digit} \rangle \; \langle \text{number} \rangle$$
$$\longrightarrow \langle \text{digit} \rangle \; \langle \text{digit} \rangle \; \langle \text{digit} \rangle \longrightarrow \text{'0'} \; \text{'1'} \; \text{'6'}.$$

That's it! That's all of BNF. It can be hard to use and hard to reason about, but the syntax is simple.

# Example: Better integers

How should we redefine ⟨number⟩ to allow negative numbers, but forbid leading zeroes? (Assume we have '0' through '9' and '−' as tokens.)

# Example: Better integers

How should we redefine ⟨number⟩ to allow negative numbers, but forbid leading zeroes? (Assume we have '0' through '9' and '−' as tokens.)

Some sanity checks for any such re-definition:

- '−' '1' '0' should be a ⟨number⟩.
- '0' should be a ⟨number⟩.
- '0' '1' shouldn't be a ⟨number⟩.
- '−' '0' shouldn't be a ⟨number⟩.

## Example: Better integers

How should we redefine ⟨number⟩ to allow negative numbers, but forbid leading zeroes? (Assume we have '0' through '9' and '−' as tokens.)

Some sanity checks for any such re-definition:

- '−' '1' '0' should be a ⟨number⟩.
- '0' should be a ⟨number⟩.
- '0' '1' shouldn't be a ⟨number⟩.
- '−' '0' shouldn't be a ⟨number⟩.

There are multiple approaches — there's no such thing as the "right" expression of a grammar in BNF. Here's one way:

⟨posDigit⟩ ::= '1' | '2' | '3' | '4' | '5' | '6' | '7' | '8' | '9'
⟨posNumber⟩ ::= ⟨posDigit⟩ | ⟨posNumber⟩ ⟨posDigit⟩ | ⟨posNumber⟩ '0'
⟨number⟩ ::= ⟨posNumber⟩ | '0' | '−' ⟨posNumber⟩

# Parse trees

⟨posDigit⟩ ::= '1' | '2' | '3' | '4' | '5' | '6' | '7' | '8' | '9'

⟨posNumber⟩ ::= ⟨posDigit⟩ | ⟨posNumber⟩ ⟨posDigit⟩ | ⟨posNumber⟩ '0'

⟨number⟩ ::= ⟨posNumber⟩ | '0' | '−' ⟨posNumber⟩

The goal of parsing is to convert a list of tokens into a **parse tree** or **concrete syntax tree** (**CST**) which gives its BNF structure. E.g. for "−886":



Each non-terminal symbol is a node. Its children are its BNF expansion, in order from left to right — so the leaves are precisely the tokens.
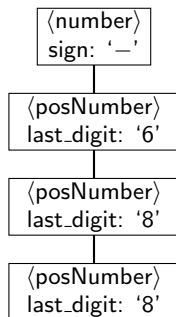
# Abstract syntax trees

$$\langle posDigit \rangle ::= \text{'1'} \mid \text{'2'} \mid \text{'3'} \mid \text{'4'} \mid \text{'5'} \mid \text{'6'} \mid \text{'7'} \mid \text{'8'} \mid \text{'9'}$$

$$\langle posNumber \rangle ::= \langle posDigit \rangle \mid \langle posNumber \rangle \langle posDigit \rangle \mid \langle posNumber \rangle \text{ '0'}$$

$$\langle number \rangle ::= \langle posNumber \rangle \mid \text{'0'} \mid \text{'−'} \langle posNumber \rangle$$

For efficiency and convenience, we may choose to process a CST into an **abstract syntax tree** (**AST**), which contains the same information in a more convenient form. E.g. we might decide we don't need the $\langle posDigit \rangle$ nodes:

Consider this grammar for simple arithmetic expressions.

$\langle$expression$\rangle$ ::= $\langle$number$\rangle$ | $\langle$expression$\rangle$ $\langle$operator$\rangle$ $\langle$expression$\rangle$ |
$\qquad\qquad\qquad$ '(' $\langle$expression$\rangle$ $\langle$operator$\rangle$ $\langle$expression$\rangle$ ')'

$\quad\langle$operator$\rangle$ ::= '+' | '−' | '∗' | '/' | '^'

# Ambiguity

Consider this grammar for simple arithmetic expressions.

$\langle$expression$\rangle ::= \langle$number$\rangle \mid \langle$expression$\rangle \langle$operator$\rangle \langle$expression$\rangle \mid$
$\qquad\qquad\qquad$ '(' $\langle$expression$\rangle \langle$operator$\rangle \langle$expression$\rangle$ ')'

$\quad\langle$operator$\rangle ::=$ '+' | '$-$' | '$*$' | '/' | '$\hat{\ }$'

Then a parser could output several valid CSTs for e.g. $(3 + 4) * (5 - 1)/3$.

This **ambiguity** can be dealt with *as long as* the semantic meaning is not ambiguous. E.g. here it is the same for all CSTs.

# Generating parse trees

Parsing is a difficult and subtle problem, but a well-understood one.

# Generating parse trees

Parsing is a difficult and subtle problem, but a well-understood one.



NOOOOOOOOOOOOOO!!!! YOU
CAN'T HAVE THE COMPUTER WRITE YOUR
PARSER FOR YOU!!!! WHAT ABOUT
NON-CONTEXT-FREE GRAMMARS AND
RECURSIVE DESCENT AND ALL THE BEAUTIFUL
SUBTLETY OF TYPE THEONO NOOOOOOOOOOOOO

haha yacc
go brrrrrrrr

Source: Generated with imgflip (here).

This means we shouldn't try to solve it again ourselves! We should instead use a **parser generator** which takes our grammar in BNF form and outputs code for a parser in a language of our choice. (E.g. yacc for C.)

# Extended Backus-Naur Form (EBNF)

Often both parser generators and language specifications add extra syntax to BNF for usability, but there's no one standard. Based loosely on ISO 14977, we'll add:

- ()s mean grouped terms, e.g. ('0' | '1') ('0' | '1') means 00, 01, 10 or 11.
- []s mean optional terms, e.g. ['0'] '1' means 01 or 1.
- {}s mean repetition, e.g. {'0' | '1'} means any number of zeroes and ones (including the empty string).
- $A - B$ means anything that matches $A$, but doesn't match $B$, e.g. ⟨number⟩ − ⟨digit⟩ means any number that's not a ⟨digit⟩.[1]

This doesn't let BNF express any grammars it couldn't before (why not?), but it does make it much nicer to read and write. For example:

$$⟨digit⟩ ::= \text{‘0’} | \text{‘1’} | \text{‘2’} | \text{‘3’} | \text{‘4’} | \text{‘5’} | \text{‘6’} | \text{‘7’} | \text{‘8’} | \text{‘9’}$$
$$⟨number⟩ ::= ([\text{‘}-\text{’}]\, (⟨digit⟩ - \text{‘0’})\, \{⟨digit⟩\}) | \text{‘0’}$$

With EBNF, we can build a readable grammar for all of C, never mind Hack! See for example here (credit Samuya Debray).

---

[1] This is only valid EBNF when $B$ expands into one of finitely many possible sequences of tokens, so e.g. ⟨number⟩ − ⟨posNumber⟩ wouldn't be valid.