

John Backus

"How do you specify a language?"

- Syntax

"3 < > 4"

is badly formed C.

Colourless green ideas sleep furiously.

"What is the shape of well-formed code?"

- Semantics

- Semantics

"What does a program mean?"

Backus-Naur form

Context-free grammars

A Python program means what the CPython interpreter does.

```
void foo() {
    int x = 3;
}
```

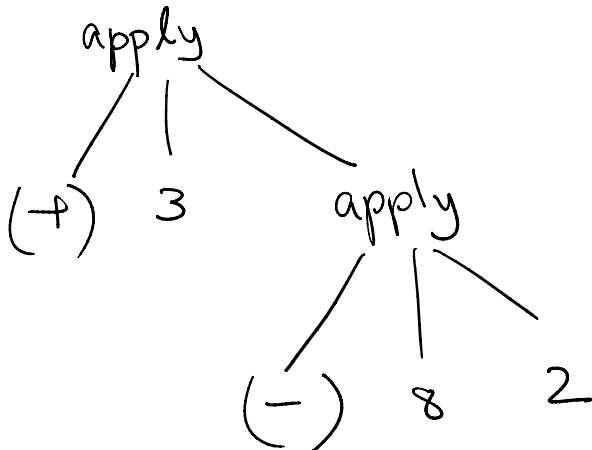
```
void foo() {
    int int = 5;
}
```

```
void >> {
    x
    int = 8
}
```

Well-formed program

Ill-formed programs

3 + _____ (8 _____ -2)



* + 3 -



There is no syntactic structure

[0-9]⁺ ~ Sequence of one or more ASCII digits

$[0-9]^+$ \rightsquigarrow Sequence of one or more ASCII digits
Expr ::= Integer Literal \rightsquigarrow 38 | 42 | 96 | ...

| Expr₁ '+' Expr₂ \rightsquigarrow Two expressions,
with a plus
sign in between

| Expr₁ '-' Expr₂ \rightsquigarrow Two expressions,
with a minus
sign in the middle

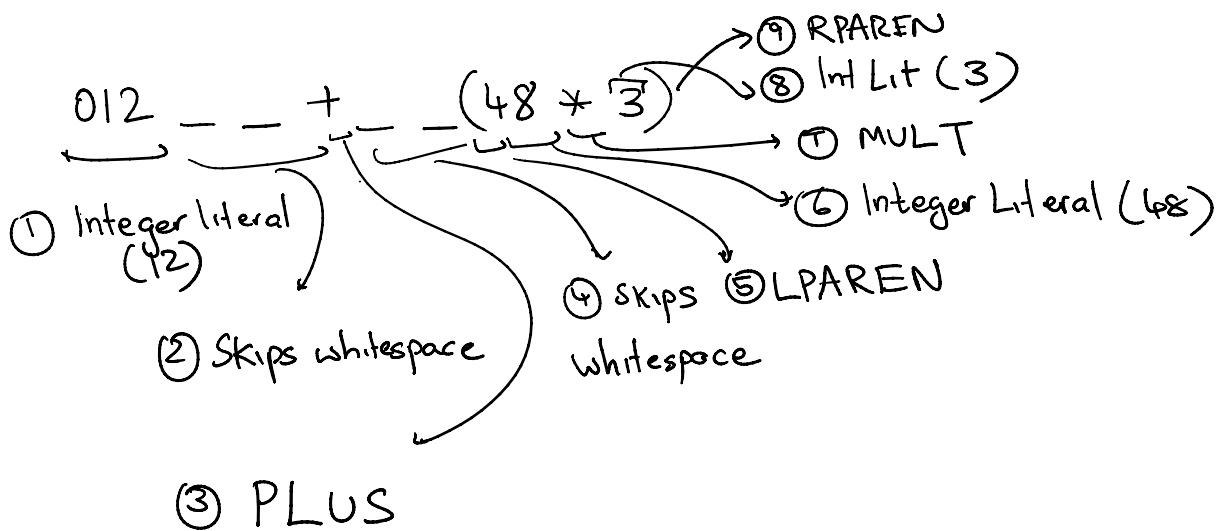
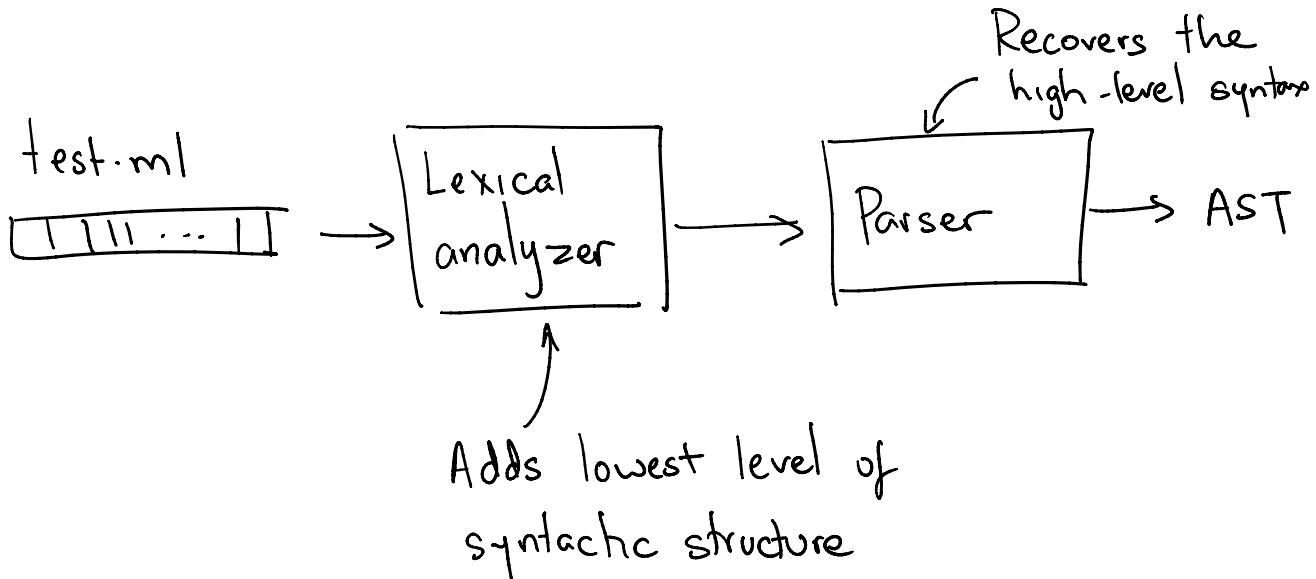
| (' Expr ') \rightsquigarrow An expression, within
parentheses.

Question : How does one go from
sequences of characters
to a parse tree?

Lots of remaining issues?

- Is "for" an identifier or a keyword?
- What about whitespace?
- What about comments?
- What about preprocessor macros?

- What about preprocessor macros?



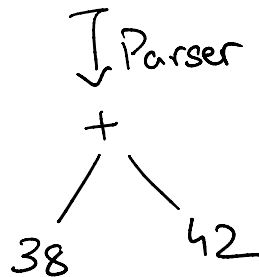
Set of tokens:

- ① PLUS ② MINUS ③ MULT
 - ④ LPAREN ⑤ RPAREN
- } Simple tokens

⑥ Integer Literal (C: int) Token with payload.

"38" \mapsto Integer Literal (38)

"38 + 42" $\xrightarrow{\text{Lexer}}$ [Integer Literal (38);
PLUS ;
Int Literal (42)]



"96 95" $\not\rightarrow$ Int (9695)

"9695" \rightarrow [Int (9); Int (6); Int (9); Int (5)]

"9695" \rightarrow [Int (9); Int (695)]

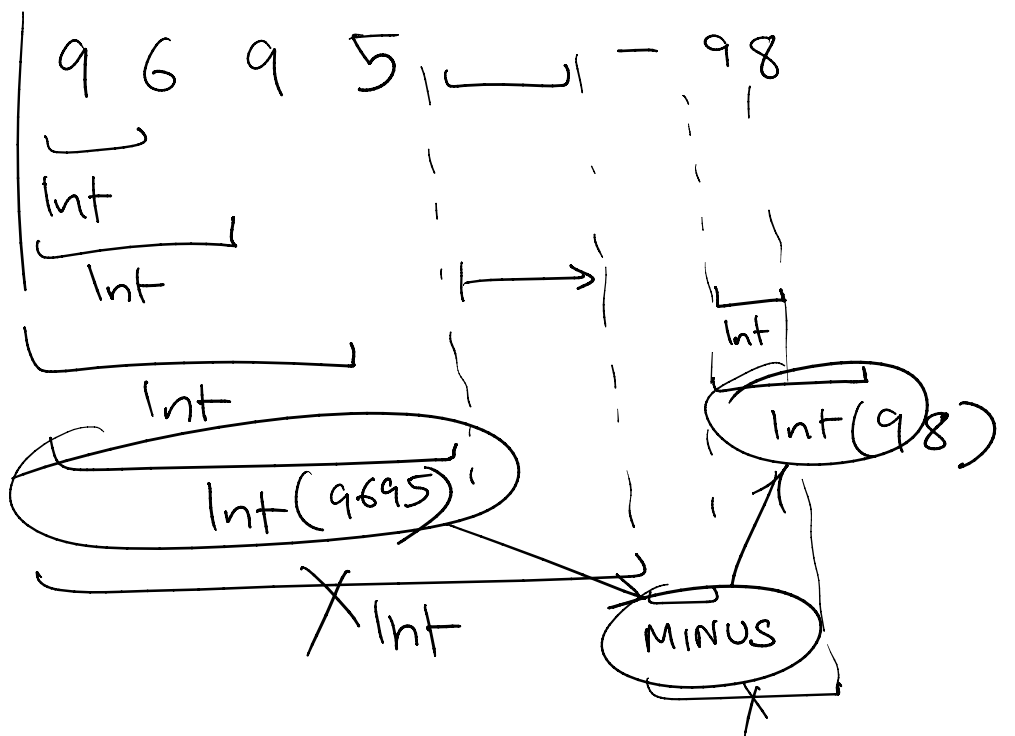
"9695" \rightarrow [Int (96); Int (95)]

⋮

⌈ . ⌋

[Int(9695)]

"Longest match principle".



For practical regex matching,

look for how it is done in Perl.

Simple "academic" regular expressions.

Digit = '0', '1', '2', ..., '9'

= ['0' - '9'] ← A character or a set of characters is a regular expr.

INT = "int"
FLOAT = "float"] A sequence of characters is a regular expr
↓
Regular exprs.

Int Lit = DIGIT[⊕] | '_' DIGIT[⊕]
An option between two reg exprs is itself a reg expr.

Ident = Alpha [AlphaNum]^{*}
'a' - 'z' 'a' - 'z'
'A' - 'Z' 'A' - 'Z'
 '0' - '9'
A sequence of zero or more occurrences of Alpha Num.

Alpha Num.

If r is a regular expr.

then r^* is also a reg. expr.
 r^+