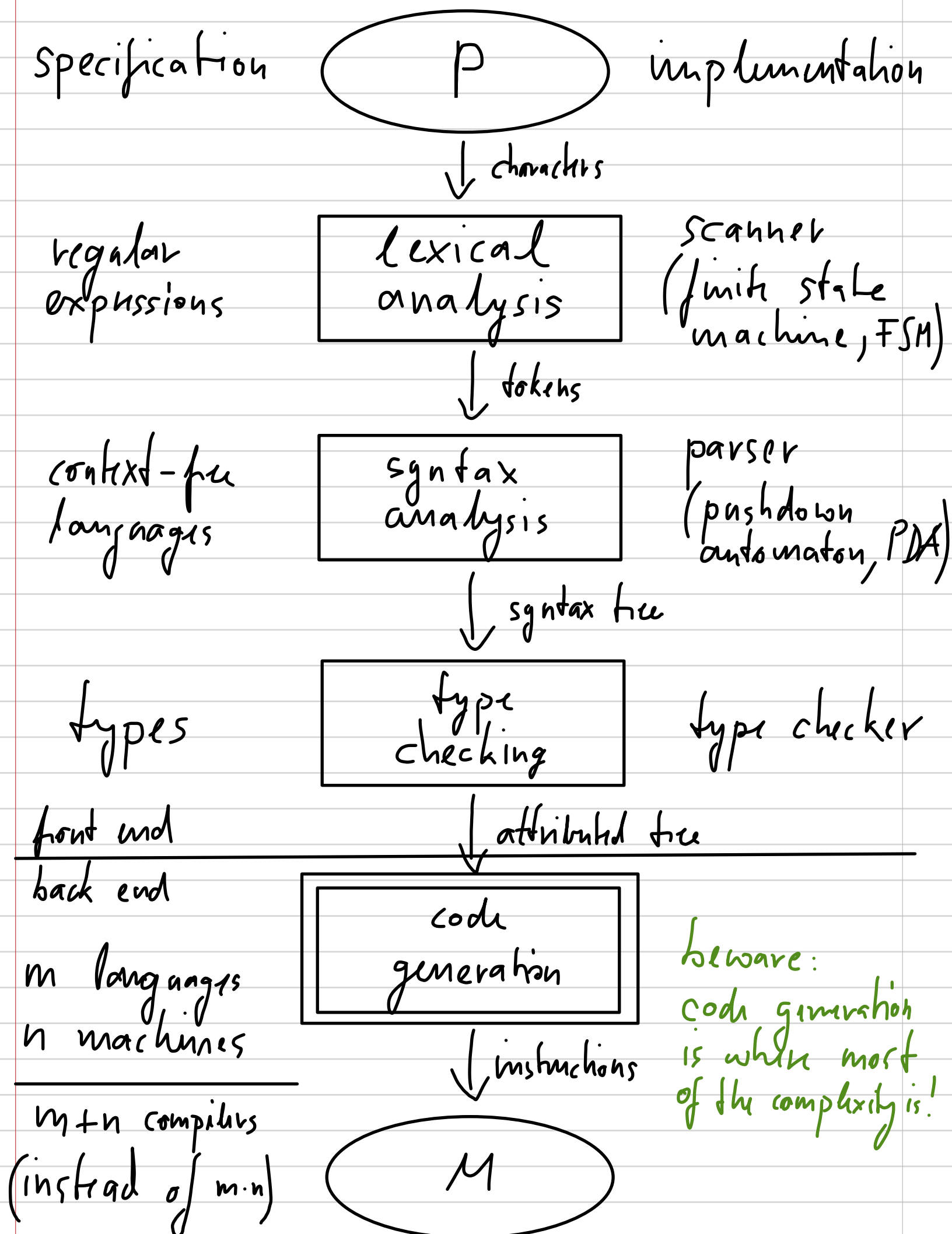


Structure



Scanning and Parsing

Sequence of symbols

①

```
main() {
    int i;
    i = 0;
    while (i < 10) {
        i = i + 1;
    }
    i = 11;
}
```

③

```
main() { |
    __int__ i; |
    __i__ = __0; |
    __while__(i < __10) { |
        ___i__ = __i__ + __1; |
    } |
    __i__ = __11; |
}
```

scanning

token for int keyword
whitespace

④

```
main() { int i; i = 0; while (i < 10) { i = i + 1; } i = 11; }
```

sequence of tokens

parsing

24-"main" 8 9 27
43 24-"i" 17
24-"i" 21 16-'0' 17
33 8 24-"i" 18 16-'10' 9 27
24-"i" 21 24-"i" 11 16-'1' 17
28
24-"i" 21 16-'11' 17
28

token for number

token for identifier

token for while keyword

⑤

```
main
__int__
__i__
__=__
__i__
__0__
__while__
__<__
__i__
__10__
__=__
__i__
__+__
__i__
__1__
__=__
__i__
__11__
```

a scanner recognizes symbols (sequences of characters) and returns tokens (integers) that represent symbols uniquely

How to Construct a Scanner!

the scanner maintains a global variable currentCharacter which is initialized to the first character of the input program by invoking a library procedure:

readCharacter(); ← implemented in runtime

which reads the next character from the input program. the scanner is invoked by the parser through a procedure:

getSymbol(); ← keep track of location in source (line, col.)

which returns the token that represents the next symbol (in another global variable). for each invocation of getSymbol() the scanner checks if currentCharacter already constitutes a valid symbol. if yes, the scanner invokes readCharacter() (to prepare for the next invocation of getSymbol()) and returns the appropriate token. if not, the scanner keeps invoking readCharacter() until it recognizes a valid symbol or returns an error.

1. define the set of valid symbols

(identifiers are sequences of letters and digits that start with a letter; numbers are sequences of digits; strings are sequences of printable characters in quotes)

2. define the set of keywords

3. define what a comment is

4. define symbol-to-token mapping

5. implement in your language

// or /* */

distinguish
"=" from
"=="

group symbols of
of the same "class"

check identifiers
if they are keywords
using a lookup table
or code

Summary:

- scanning _{words} vs. parsing _{sentences}
- focus on scanner, no parsing yet
- write first input for unit testing (very important!)
- get started today! even if you have to rewrite later

including meaningless input,
syntactically correct but also
syntactically incorrect!