

The first 3

can be understood

The structure of a compiler

by analogy to how humans comprehend a natural language e.g. if $x = y$ then $z=1$; else $z=2$;

Keyword identifier equality predicate

2. Parsing (Syntax Analysis)

Once words are understood then the next step is to understand sentence structure

Parsing = Diagramming Sentence tree

3. Semantic Analysis

e.g. scope checking and type checking

Once the structure is understood, we can try understanding the "meaning"

- but meaning is hard for compilers

- compilers use semantic analysis to catch inconsistencies (because for understanding the full meaning of a sentence I need to have background knowledge)

- but compilers have limited analysis power

- for checking the type of error we don't need background knowledge e.g. type mismatch



4. Optimization

the better (conserve resources)

1. Run faster → efficient

2. Use less memory → memory requirement

5. Code Generation (only when the code is valid)

Produces assembly code usually

A translation to another language

Programming language = Design

Compiler = implementation

Language design has big impact on compiler

- Determines what is easy and hard to compile

- many trade-offs in language design called (design and compilation trade-off)

The proportion has changed since (Fortran) → First high level programming language

Before: lexing, parsing ⇒ expensive

Now: Optimization is King

* application domains have distinctive and conflicting needs

- Science ⇒ Parallelism, floating point, arrays and ops

- business ⇒ Persistence, Report generation, Data analysis e.g. SQL

- System Languages ⇒ Embedded and Operating Systems

For such applications you need \Rightarrow very low level control over resources
Like working with registers

Lexical Analyzer

- Scans the pure HLL code line by line
- takes Lexemes as input and produces Tokens
- Removes comments and whitespaces from pure HLL code
- Helps in macro expansion in pure HLL code

* in the implementation of a lexical analyzer it works on a DFA

Errors in lexical Analysis

- Identifier that are way too long. C: 247 C++: 2048 Python: 79
- Exceeding length of numerical constants. ex: $\text{int } i = \underbrace{4567891;}_6$
- ill-formed numerical constants. ex: $\text{int } i = 47\frac{1}{8}$ size 2Byte $\sim 32,768 \rightarrow 32,767$
- illegal characters that are absent from the source code ex:

Issues with lexical analysis

- Lookahead
- Ambiguities

Char $x[0] = "Hello";$; ;

we Specify lexers using \Rightarrow RE or Examples of a RE

1950's

Output of a lexical analyzer is a stream of tokens $\xrightarrow{\text{go to}} \text{Parser (syn anal)}$

* Fortran How can a lexical analyzer be implemented:

white spaces are insignificant

1. Recognizing Substrings corresponding to tokens
2. Return the value of a lexeme of a token (lexeme is the substring)

"Lookahead" might be used to decide where one token ends and the other token begins
So we have to use lookahead

C++ do Foo<Bar<Baz>> >
 ^ ws
 can't

There are several ways to specify tokens $\xrightarrow{\text{pop}} \text{RL}$ (Simple, easy to understand, Efficient)

a language is a set of strings

Need some notation for specifying which sets we want (RL $\xrightarrow{\text{notation}}$ RE)

Atomic regex \Rightarrow Union: $A + B = \{s \mid s \in A \text{ or } s \in B\}$ Concatenation: $AB = \{s \mid s \in A \text{ and } s \in B\}$
Iteration $A^* = \bigcup_{i \geq 0} A^i$ where $A^i = A \dots i \text{ times } A$

Def. The RE over S are the smallest set of expressions ex

Syntax vs. Semantics

To be careful, we should distinguish syntax (the reg. exp.) and semantics (the langs. they denote),
Meaning function L maps syntax to semantics
L: Exp \rightarrow Sets of Strings

$$\begin{aligned} L(\epsilon) &= \{\} \\ L(c^n) &= \{c^n\} \\ L(A+B) &= L(A) \cup L(B) \\ L(AB) &= \{ab \mid a \in L(A) \text{ and } b \in L(B)\} \\ L(A^*) &= \bigcup_{i \geq 0} L(A)^i \end{aligned}$$

Identifier = letter (letter + digit)* \neq (letter* + digit*)

a letter followed by a | this means either all letters
bunch of letters or digits follows by all digits

in Summary RE describe many useful languages / RL = Language Specification

Lexical analysis (implementation)

- Union: $A \cup B \equiv A + B$
- Option: $A + \epsilon \equiv A ?$
- Range: $'a' + 'b' + \dots + 'z' \equiv [a-z]$
- excluded Range: complement of $[a-z] \equiv [\wedge a-z]$

RE \Rightarrow Lexical Spec Steps

1. write a regexp for the lexemes of each token ex. Number = digit+...
Keyword = "if" + "else" ...
2. Construct R, matching all lexemes for all tokens $R = \text{Keyword} + \text{Identifier} + \text{Number}$
3. Let input be $x_1 \dots x_n$ for $1 \leq i \leq n$ check $x_i \dots x_n \in L(R)$
 $= R_1 + R_2 + R_3$
4. If Success, then we know that $x_1 \dots x_j \in L(R_j)$ for some j
5. Remove $x_1 \dots x_j$ from input and go to (3)

Ambiguities

1. How much input is used? Like $x_1 \dots x_i \in L(R)$ Rule: Pick the longest possible string
- it's called the maximal match

2. Which Token is used ex. $x_1 \dots x_j \in L(R_k)$ which one do we choose
Rule: use rule listed first (Priority) (j if $j > k$)

Error Handling = Lowest priority

Summary

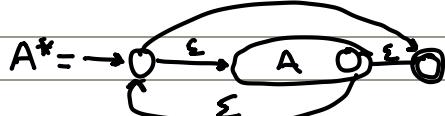
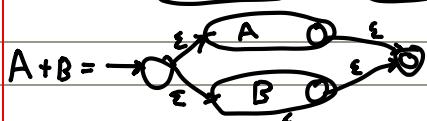
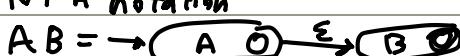
- RE provide a concise notation for string patterns
- Use in lexical analysis requires small extensions
- require only single pass over the input
- Good algorithm known
- few operations per character (table lookup)

Finite Automata

- RE = spec & FA = impl
- A FA is made of
 - 1. input alphabet Σ
 - 2. set of states S
 - 3. start state n
 - 4. set of accepting states $F \subseteq S$
 - 5. set of transitions state $\xrightarrow{\text{input}} \text{state}$

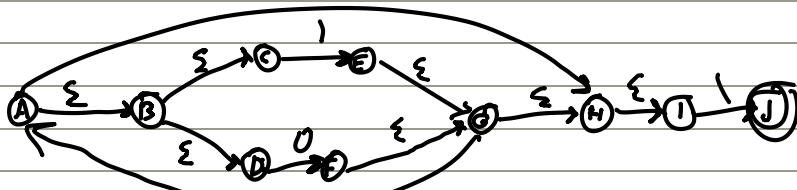
Lexical Specification \rightarrow RE \rightarrow NFA \rightarrow DFA \rightarrow Table-driven implementation of DFA

NFA notation

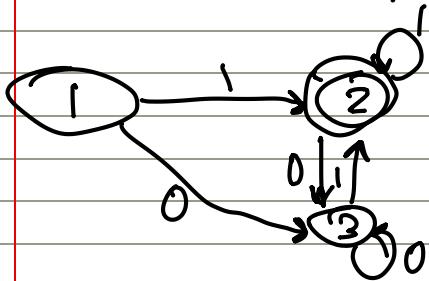


$(I + \emptyset)^*$

(ab^*ab^*) even
 $a(ab^*ab^*)$ odd



Σ -enclosure(A) $\{\Sigma, A, B, C, D, H, I\} \rightarrow 1$ | $\Sigma - \{A\}$ | I | O
 Σ -enclosure(E, J) $\{\Sigma, E, G, H, I, A, B, C, D\} \rightarrow 2$
 Σ -enclosure(F) $\{G, H, I, A, B, C, D\} \rightarrow 3$



Languages are developed around the von neumann architecture

Influences on Language design (Computer architecture / Programming methodologies)

why is the von neumann architecture dominant

1. Data and programs stored in memory
 2. memory is separate from the CPU
 3. instruction and data are piped from the CPU
 4. basis for imperative languages
- variables model memory cells
 - assignments statements model piping
 - iteration is efficient
 - selection and branching statements

algorithm

```

init Program Counter
rep loop
    get instruction pointed by the counter
    counter += 1
    decode the instruction
    execute the // /
end rep

```

* a connection between a CPU and RAM = speed of computer

Load module : the user and system load together

Linking and Loading : the process of collecting system program units and linking user programs

RAM bottleneck : instruction execution speed > connection speed

↳ primary limiting factor in PC speed

Interpreter

Older

used in Scripting Languages

Better portability

skips the compilation part, goes straight to execution

uses more space

runs programs as is (little to no preprocessing)

Compiled code

ready to Run

often faster

not cross platform

inflexible

Source code is Private

Compilers

newer

Better execution time

uses less space

compiles the code before execution

does extensive preprocessing

Interpreted Code

Cross-platform

interpreter

simpler to test

often slower

easier to debug

Source code is public

Compiled

C
C++

Objective-C

Interpreted

PHP

JavaScript

Hybrids

Python

Java

C#

VB.NET

Java

C++

a h Oh
Open it up

Readability

more readable has a simple design and easy to read

bad in readability since it has pointers which are hard to read

Writeability

although Java is great in writeability we gotta give it to C++ since it uses pointers

C++ is "Him!!" in terms of writeability because it uses Pointers

Reliability

Java is better 1. uses type checking 2. exception handling and indexing 3. check the OOB

bad because as mistakes can happen when we use arithmetic Pointers

Portability

Java is better you can compile it to an intermediate language (Byte Code)

has to recompile it

Cost of execution

higher cost of execution since its a strong checking language ex. OOB while C++ doesn't check

C++ translate directly into machine code, while Java → intermediate language → machine code ↑

NFA

Can have multiple transition for the same input

Σ -moves

Slower to execute (more paths to consider)

Simpler and smaller

One transition per input per state

No Σ -moves

Faster to execute (no choices to consider)

more complex and can be exponentially larger

$2^n - 1$ states

States : S

Start state: $s \in S$

Final states: F subset of S

The transition function:

$a(x) = \{y \mid x \in X \wedge x \xrightarrow{a} y\}$

States : subset of S

Start state: ϵ -closure(s)

Final state: $\{X \mid X \cap F \neq \emptyset\}$

The transition function:

$x \xrightarrow{a} y \text{ if } y = \epsilon\text{-closure}(a(X))$