

ENEE645

(4)

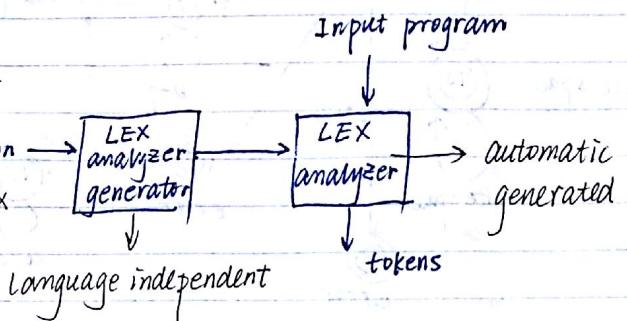
E.g. "if8" matches both identifier as well as keyword "if" followed by number "8".

According to the priority, it will be matched as keyword "if" followed by number "8".

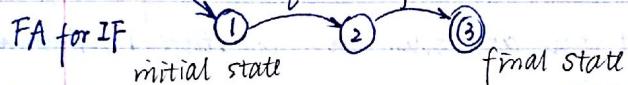
but solution: longest match rule

language dependent  
What we want: set of regular expression  
the one-to-one (token type) &  
(regular expression) table

E.g. "lex" in Unix



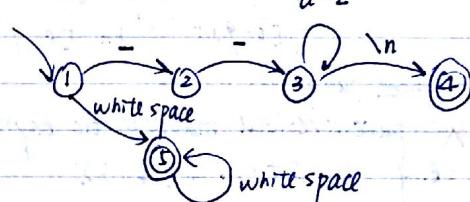
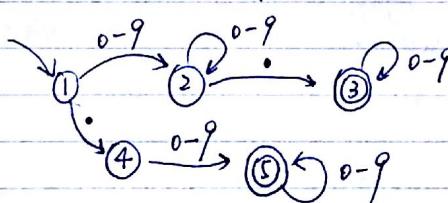
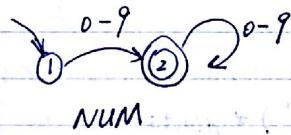
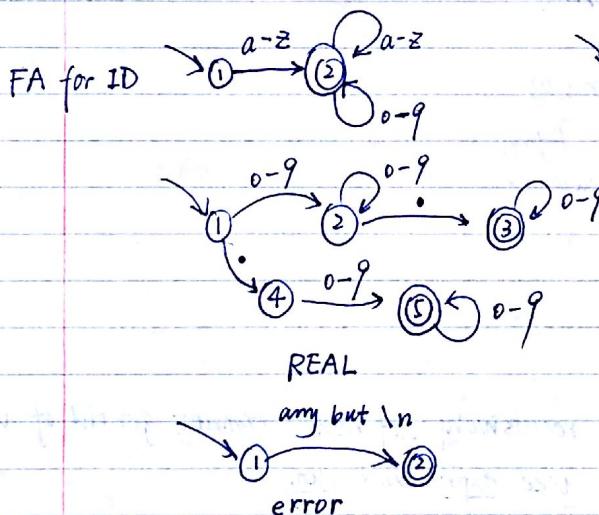
Finite Automata (FA) (finite state machine)



Regular expressions essentially can be transformed into FA

① start state

② final state

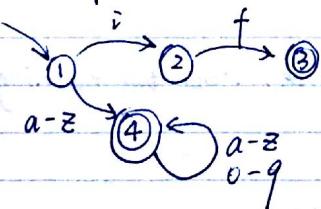


means space, new line, tab

/\* do nothing \*/

Each of the finite state machines can be implemented by switch in C.

Combine IF & ID FAs

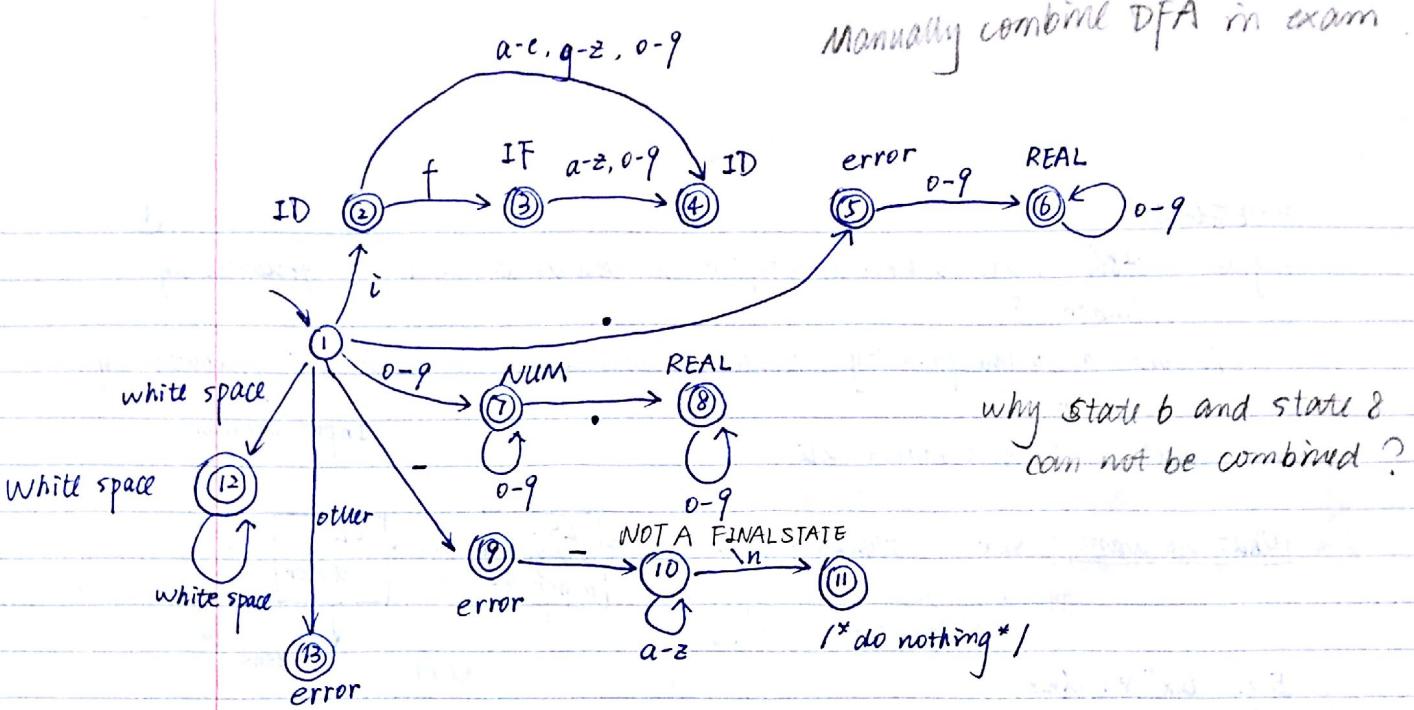


This cannot work, since it cannot determine which way to go when see a "i" in state 1.

Non-deterministic FA (NFA)  
→ Unimplementable

What we need is DFA. (Deterministic FA)

It is possible to convert NFA to DFA, (but will not be covered in this course)



In this combined FA, what happens when there is no transition for the character we see next?

→ Take any action in the current state and then move to state 1, meanwhile then, look at the same character again.

$$2.12 \quad \text{digits} = [0-9]^+ \quad \text{How to recognize sum}$$

$$\begin{aligned} \text{sum} &= (\text{digits } "+")^* \text{ digits} \quad \text{not regular expression} \\ &= ([0-9]^+ "+")^* [0-9]^+ \quad \text{regular expression} \end{aligned}$$

recognize fully

How to → parenthesized mathematic expression

$$\text{E.g. } ((15+17)+(14+3))$$

$$\text{digits} = [0-9]^+$$

$$\text{sum} = \text{expr } "+" \text{ expr}$$

$$\text{expr} = "(" \text{ sum } ")" \mid \text{digits}$$

> recursively defined, cannot get rid of variables like expr and sum

It is impossible to recognize fully parenthesized mathematic expression in regular expression. Because you have to track the number of open parenthesis, but which can be infinite.

Context-free grammars (CFGs)

the mean natural languages are not context free, since the meaning of one sentence may depend on the previous sentences.

CFGs → Terminals (usually tokens in language)

→ Non-terminals (variable) e.g. expr, sum

→ Production rules which define non-terminals in terms of terminals as well as non-terminals.

ENEE 645

(5)

e.g. 1. sum = expr "+" expr  
expr = "(" "sum ")" | digits → {  
2. expr = "(" ("sum") ")"  
3. expr = digits } eliminate alternation  $\{ \mid \}$

1, 2, 3 are three production rules.

1. sum → expr plus expr  
Non-terminals      terminals
2. expr → LPAREN sum RPAREN  
terminals      non-terminals
3. expr → digits

| X  
CFGs can be recognized by stack machine,  
which is infinite.

### Derivation

Example of a complete CFG (Page 41)

S → S ; S      S stands for statement or statements

S → id := E      assignment in the language

S → print(L)

E → id

E → num

E → E + E

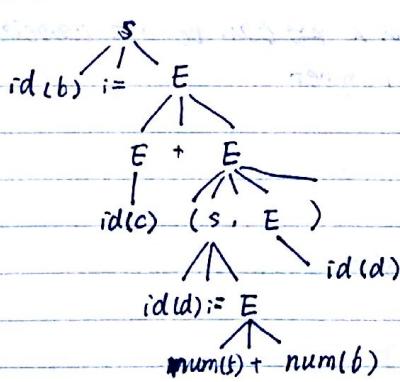
E → (S, E)

comma expression : b = c + (d = 5 + 6, d);       $\Rightarrow b = c + 11$   
statement      expression

L → E

L → L, E

parse tree



derivation

S

id := E

id := E + E

id := id + E

id := id + (S, E)

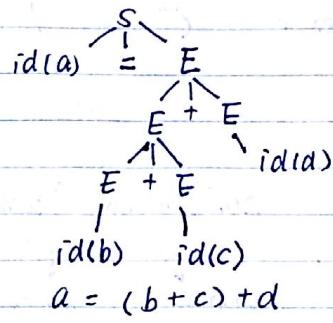
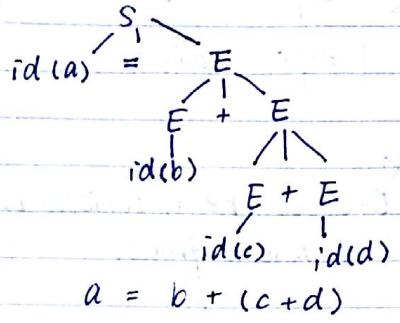
derivation } left most  
right most  
neither

there exists

Don't want ambiguous grammar where  $\exists$   
more than one parse tree for a context.

That grammar is ambiguous

E.g.  $a = b + c + d$



You may end up with different results when using floating point.  
In C, associate by left is always assumed.

2.19 A grammar is ambiguous if there exists a string, which has more than one parse trees.

grammar on P44 :

Expand on expressions (E)

$E \rightarrow id$

$E \rightarrow num$

$E \rightarrow E \times E$

$E \rightarrow E / E$

$E \rightarrow E + E$

$E \rightarrow E - E$

$E \rightarrow (E)^P$

which means the first expression will be correct

ambiguous example:  $a = 1 - 2 - 3$

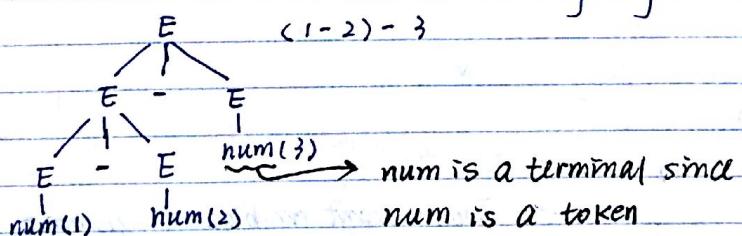
can be:  $a = (1 - 2) - 3 = -4$

or:  $a = 1 - (2 - 3) = 2$

We say an operator  $\oplus$  is associative if:  $(a \oplus b) \oplus c = a \oplus (b \oplus c)$

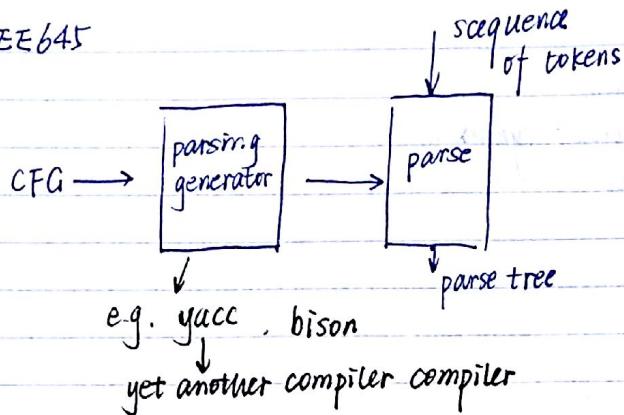
+ and  $\times$  are associative, - and / are not associative

However, when operating float, + and  $\times$  ~~can be~~ can be not associative.  
extremely tiny number



ENEE645

(6)

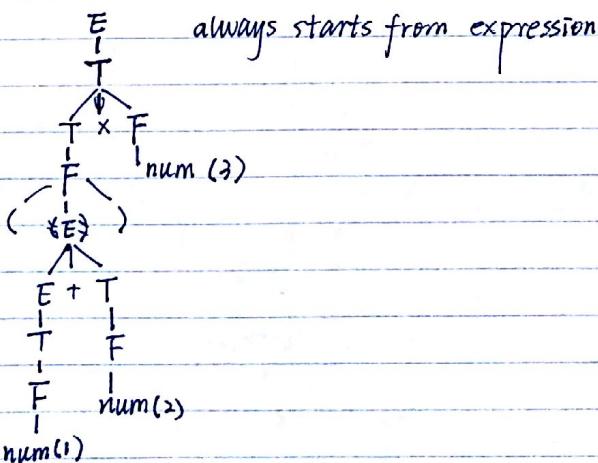


Replace  $E \rightarrow E - E$  by  $E \rightarrow E - T$  encapsulate associativity from left

Unambiguous grammar | CFG : Context-free Grammar

$E \rightarrow E + T$	$T \rightarrow T \times F$	$F \rightarrow id$
$E \rightarrow E - T$	$T \rightarrow T / F$	$F \rightarrow num$
$E \rightarrow T$	$T \rightarrow F$	$F \rightarrow (E)$ parenthesized expression
give you associativity from left		
lowest priority	next lower priority	highest priority
the three layers give the order of express operation		

$(1+2) \times 3$



Abstract syntax | This is what

what follows parsing is semantic actions.

These are fragments of C code that are executed for each parsing rule.

In a compiler, these include:

1. generating an abstract syntax tree
2. define symbol tables
3. do type checking

\*

the compiler's job is to convert source code into object code.

source code -> intermediate representation

intermediate representation -> object code

intermediate representation -> assembly language

intermediate representation -> machine language

selection of intermediate representation depends on target machine

ENEE 645

(7)

```
int *p, a;  
p = &a;  
= a;  
*p =
```

(compiler register, instead of hardware register)

↑  
must denote a from register to memory

pointers and arrays are the only two things that are considered as memory in compiler.

Assembly language: human readable one-to-one mapping machine code.

e.g. MIPS ISA

instructions are 32 bit

[opcode | op1 | op2 | op3] → machine code

ADD R3, R2, R6; → assembly code (nice for human, but horrible for computer)  
→ (if translate in ASCII, will be more than 104 bits)

LEX operation (The input of LEX is a sequence of ASCII characters)

Token type: P17

Token type

Example

(the name) ← ID (identifier)

foo, n14, a, last (if not a keyword)

NUM

73, 0, 00, 082 (=82)

REAL

66.1, 0.5, .5, 10.0, 10., 1e67, 5.5e-10

IF

if

FOR

for

COMMA

,

NOTEQ

!=

LPAREN

(

RPAREN

)

Language for rules of tokens is "regular expressions"

preprocessor will remove

an ID that has  
value "match0"

2.5 Example on P17

comments before LEX

Token:

ID(s) RPARE

```
float match0(char* s) /* find a zero */
```

FLOAT ID(match0) LPAREN CHAR STAR

{

LBRACE

```
    if (!strcmp(s, "0.0", 3))  
        return 0.;
```

IF LPAREN BANG ID(strcmp) ...

won't be recognized as ID(match) NUM because  
of the longest matching rule (token with longest match wins)

## Regular expressions

Language for regular expressions has:

- ① Symbols (e.g.  $a$  is symbol for itself)
- ② Alternation: ~~M/N~~  $M|N$  (either rule  $M$  or rule  $N$ )
- ③ Concatenation:  $M.N$  or  $MN$  (rule for  $M$  followed by rule  $N$ )
- ④ Epsilon ( $\epsilon$ ) matches empty string "" → the repetition of the rule
- ⑤ Repetition ( $M^*$ ) Concatenation of zero or more instances of  $M$
- ⑥  $M^+$  Concatenation of one or more instances of  $M$

$$M^+ \equiv MM^*$$

→ taking the rule " $0|1$ " and re concatenation of zero or more times.

E.g.  $(0|1)^*0$   $0100 \vee$  a member of this language

0 ✓

$(0|1)^*0$   $0101 \times$  not a member of this language

is called regular

expression

E.g.  $b^*(abb^*)^*(a|\epsilon)$  doesn't have to end with a

strings of  $a$ 's and  $b$ 's with no concatenation  $a$ 's

The difference between  $(a|\epsilon)$  and  $(a|b)$  is that the latter one cannot have empty string.

set of regular  
expression ↗

$(a|b)^*aa(a|b)^*$ : strings of  $a$ 's and  $b$ 's that contain  $aa$  somewhere

highest priority ↗

Token type

regular expression ↗

IF

if

ID

$[a-zA-Z][a-zA-Z0-9]^*$  (cannot start with a number, cannot have capital letters and underscore)

NUM

$[0-9]^+$  (zero and zero onwards)

REAL

$([0-9]^+ \cdot [0-9]^*)^*$  (not a concatenation, avoid . case)

lowest priority ↗

/\*do nothing\*/

$(\text{--} [a-zA-Z]^* \text{\\n}) \mid (\text{\\d} \mid \text{\\s} \mid \text{\\t})^*$  (nothing before and after the dot)

return no token

• → means any single character except new line

error()

## Ambiguity

E.g. ① "if" matches both keyword & identifier

solution: rule priority

ENEE 645

(3)

If you get tons of files (Separate compilation) compiling source files separately  
gcc -c file1.c -o file1.o  
...  
gcc file1.o file2.o ... -o file.exe

(Linking takes much less time than compiling.)

A makefile can detect the files that have been changed and just re-compile them (not all others) and then do the linking.

Reasons to learn about compilers:

- proliferation of non-standard architectures
- parallelism
  - { multi-core parallelism
  - hyper threading (two threads inside a core)
  - multi-media extension (MMX → → AVX)
  - gaming console : IBM cell
  - GPU
- embedded system (a simple computer system in a device whose primary purpose is not computing)
- exclusive explicit memory hierarchy
- cybersecurity
  - (the tool to detect cybersecurity S/W vulnerability use the same theory as compiler)
- ASIC (Application Specific Integrated Circuit)
- power reduction (make computers/phones consume less energy)

### Parallelism

→ Detect dependent and independent instructions

$$a = b + 1$$

...

$$c = a + 2$$

$$A[i] = \dots$$

...

$$\dots = A[j]$$

dependent instructions, cannot run in parallel  
(data dependency)

may or may not be dependent, since i might equal to j

$$*P = \dots$$

$$*P = *q$$

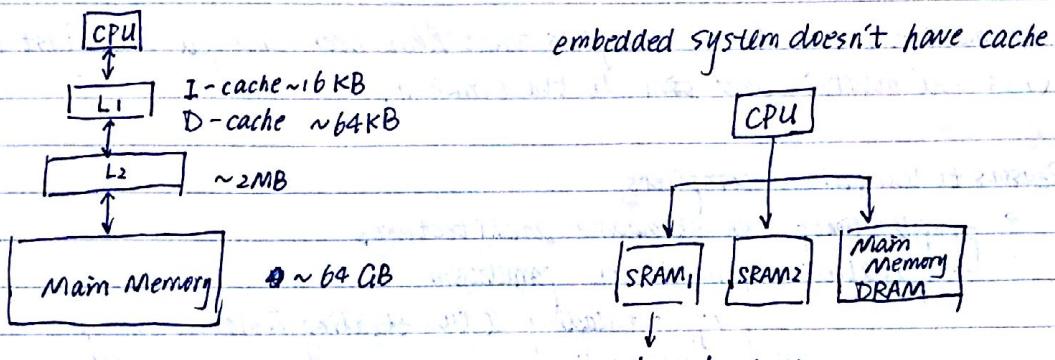
if P and q point to the same thing, there might be dependence

```

for i = 1 to 100 {
    A[i] = A[i] + 3;      This is parallel
}

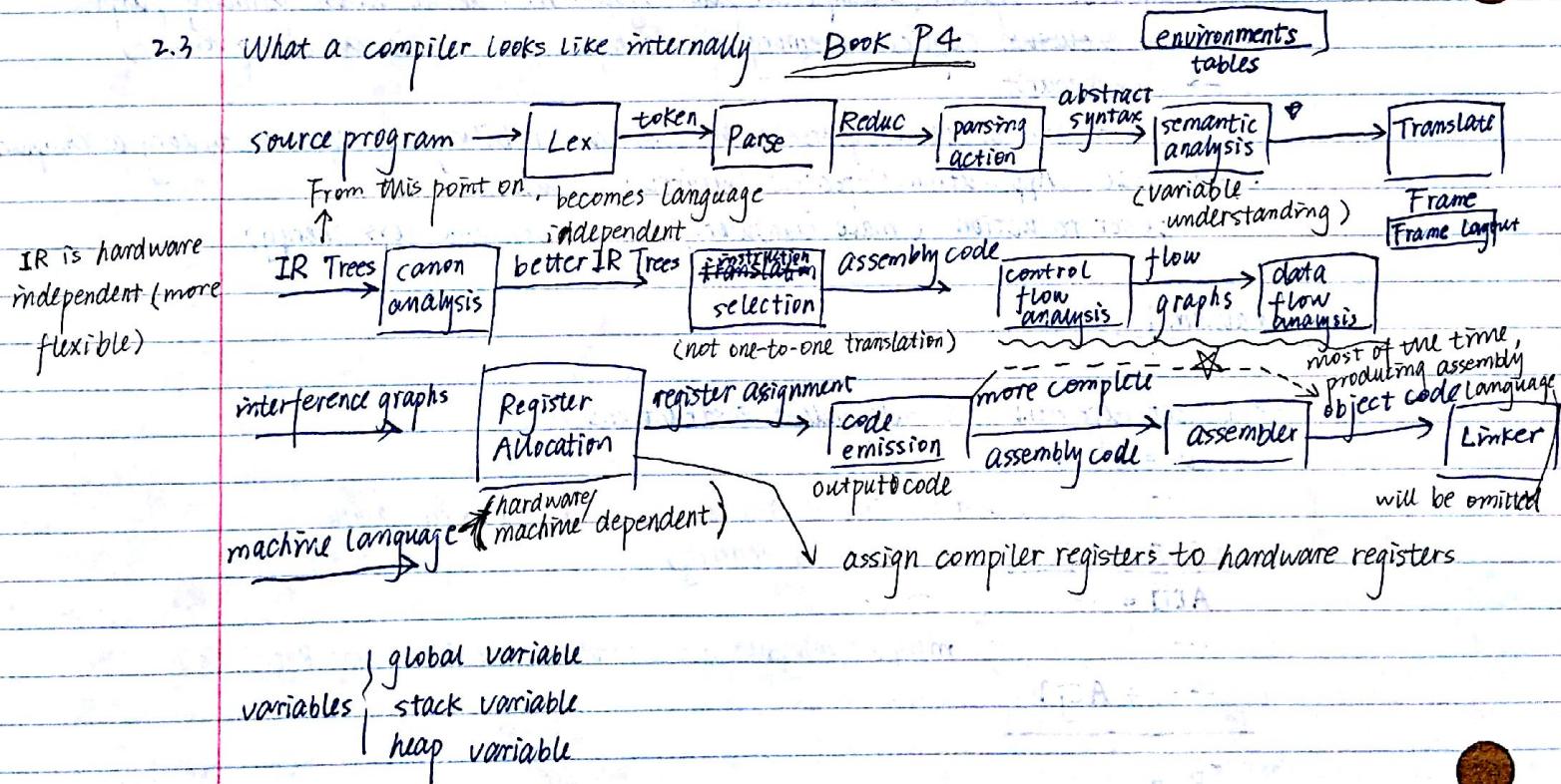
```

Affine theory  
↳ linear functions.



scratch pad memory  
Compiler will determine the addresses of variables  
(put the most frequently used to SRAM 1)

### 2.3 What a compiler looks like internally Book P4



if you name a location directly, that's called register, index accessing called memory

pointer can only point to memory location, but not register

ENEE 645

①

1.27 Prof. Rajeev B A.V.W. 1431  
piazza.com (for discussion)

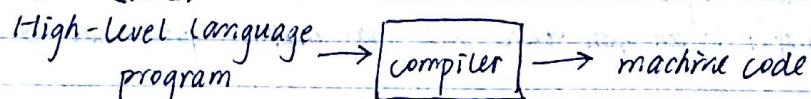
4-5 homeworks 20%

1 midterm exam 30%

1 software project 15% (use LLVM) open-source compiler

1 final exam 25%

Compiler is a piece of software that translates high-level language into machine code  
(HLL)



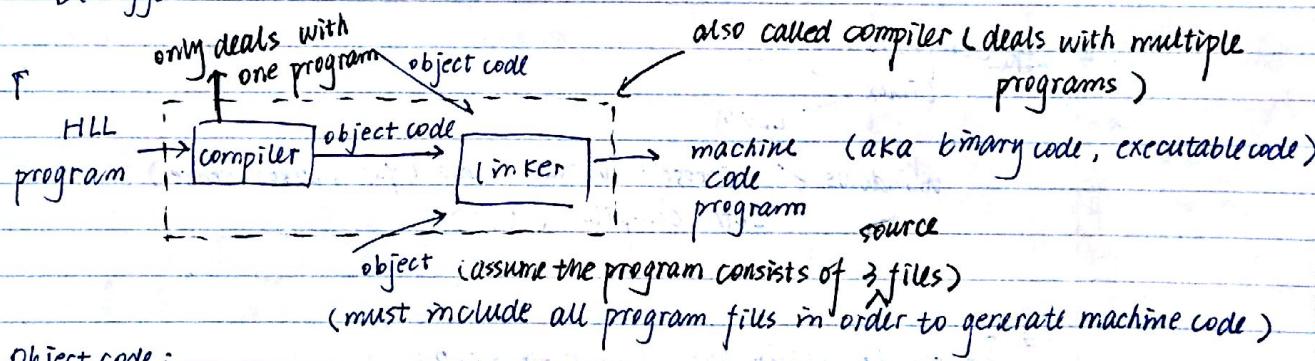
software development process:

Compiler (assembler may be part of compiler)

Linker

Loader

Debugger



Object code:

Like machine code, but

→ corresponds to only one source file

→ addresses of external symbols are missing

Linker: (Linker is running statically because it's run during compile time)  
→ Joins multiple object files together  
→ Resolves external addresses (Link different object files)

Loader:

→ Not a part of the compiler

→ Part of the OS.

cannot use standard libraries, as a result it  
statically linked, will cause redundancy in file system and main memory

- Copies executable from file system to main memory & transfers control to it.
- Links dynamically linked libraries (DLLs)  
= run time statically = compile time

Library: (are not necessarily part of compilers)

- a set of software routines that are often reuse in programs.

Reasons why written:

- convenience & reuse
- often call low-level resources
  - system calls (trap)
  - Assembly code

e.g. printf() is an open source code that makes system call →

Locations where linker filling addresses are specified in the relocation table.

## 1.29 Debugger

- Part of SDK (Software Development Kit)
- SDK = Compiler + Linker + Debugger + GUI

gcc Compilers:

Linux ← gcc  
LLVM

Windows ← Microsoft Visual Studio (95% market share)  
Intel compiler (ICC)

gcc

source to object file ( $xx.c \rightarrow xx.o$ )

command:  $gcc -c file.c -o file.o$   
stop at object file

interpretive language: Java .net php has intellectual property issue

source to executable file

command:  $gcc file.c -o file.exe$

if you have more than one file

$gcc file1.c file2.c ... -o file.exe$