

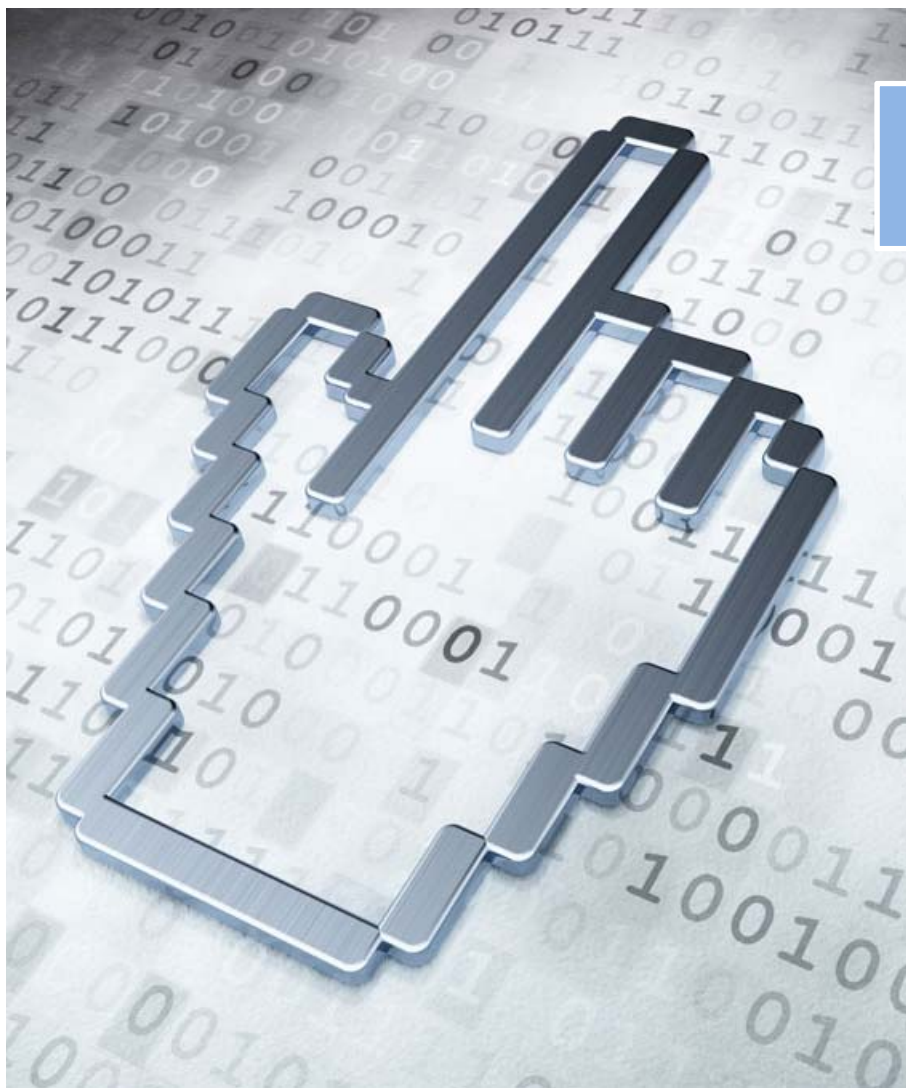


编译原理

第一章 绪论

哈尔滨工业大学 陈鄞





本章内容

- 1.1 什么是编译
- 1.2 编译系统的结构
- 1.3 为什么要学习编译原理
- 1.4 编译技术的应用

1.1 什么是编译？

类似于数学定义或自然语言的简洁形式

- 接近人类表达习惯
- 不依赖于特定机器
- 编写效率高

高级语言
(High Level Language)

引入助记符

- 依赖于特定机器，非计算机专业人员使用受限制

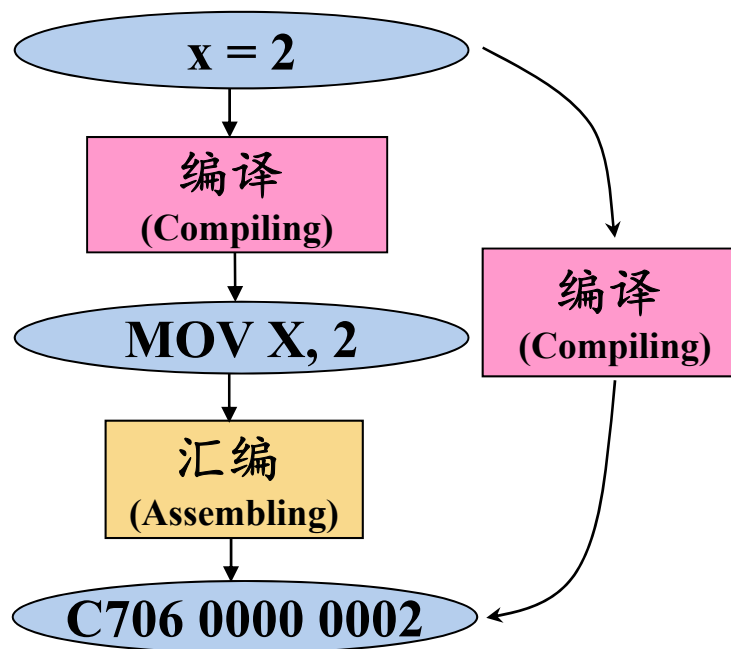
汇编语言
(Assembly Language)

- 编写效率依然很低

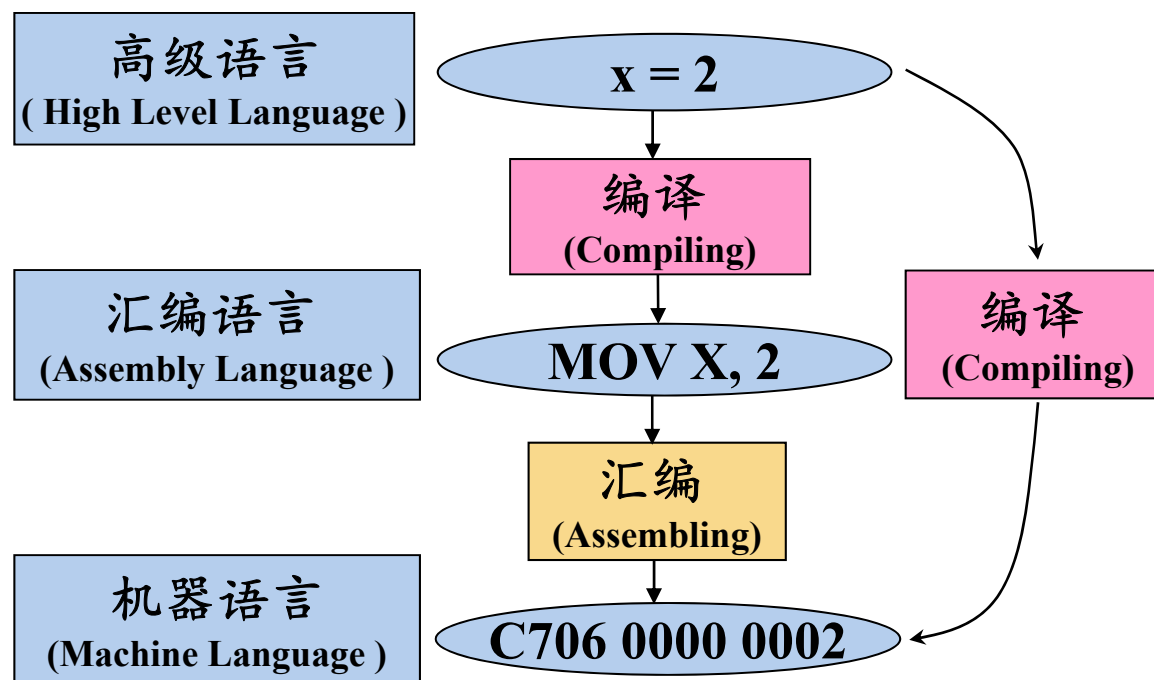
可以被计算机直接理解

- 与人类表达习惯相去甚远
- 难记忆
- 难编写、难阅读
- 易写错

机器语言
(Machine Language)

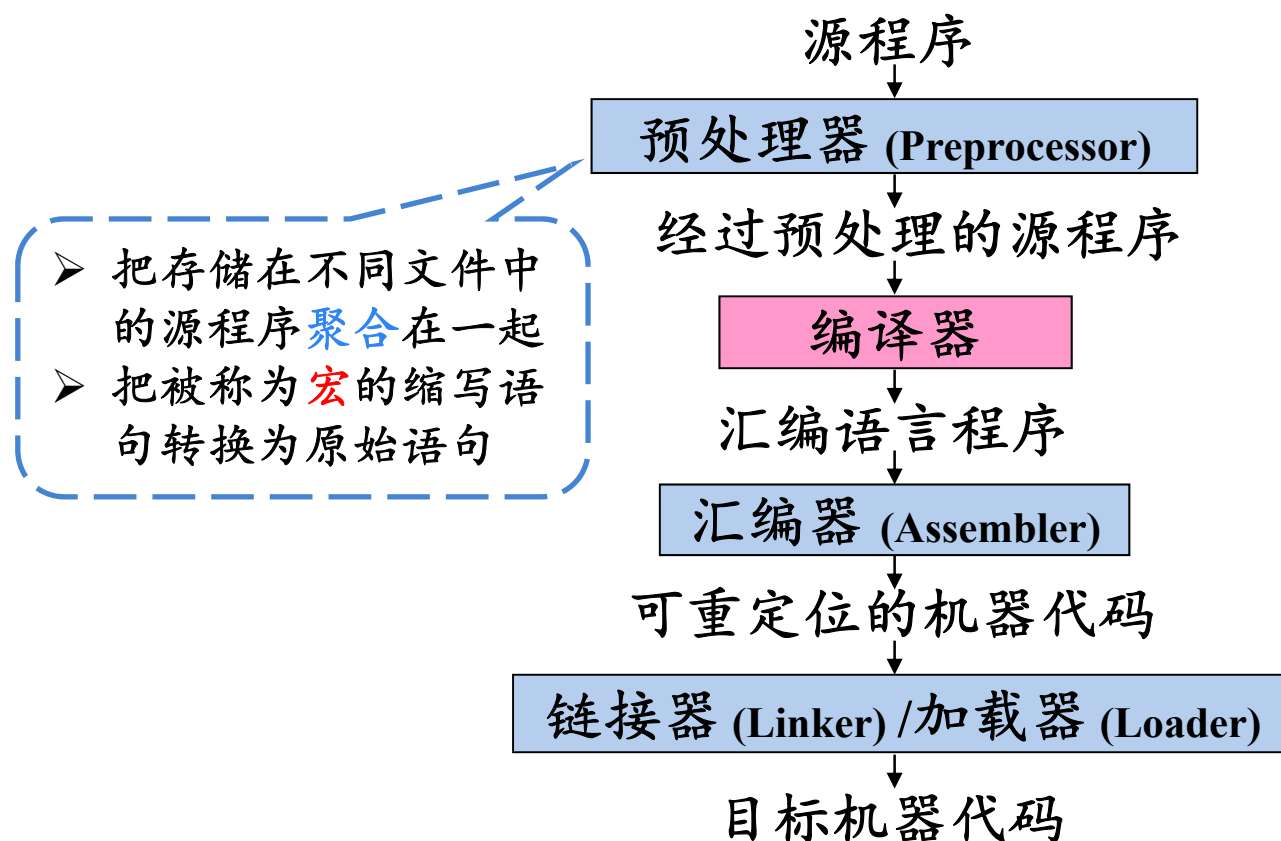


1.1 什么是编译？

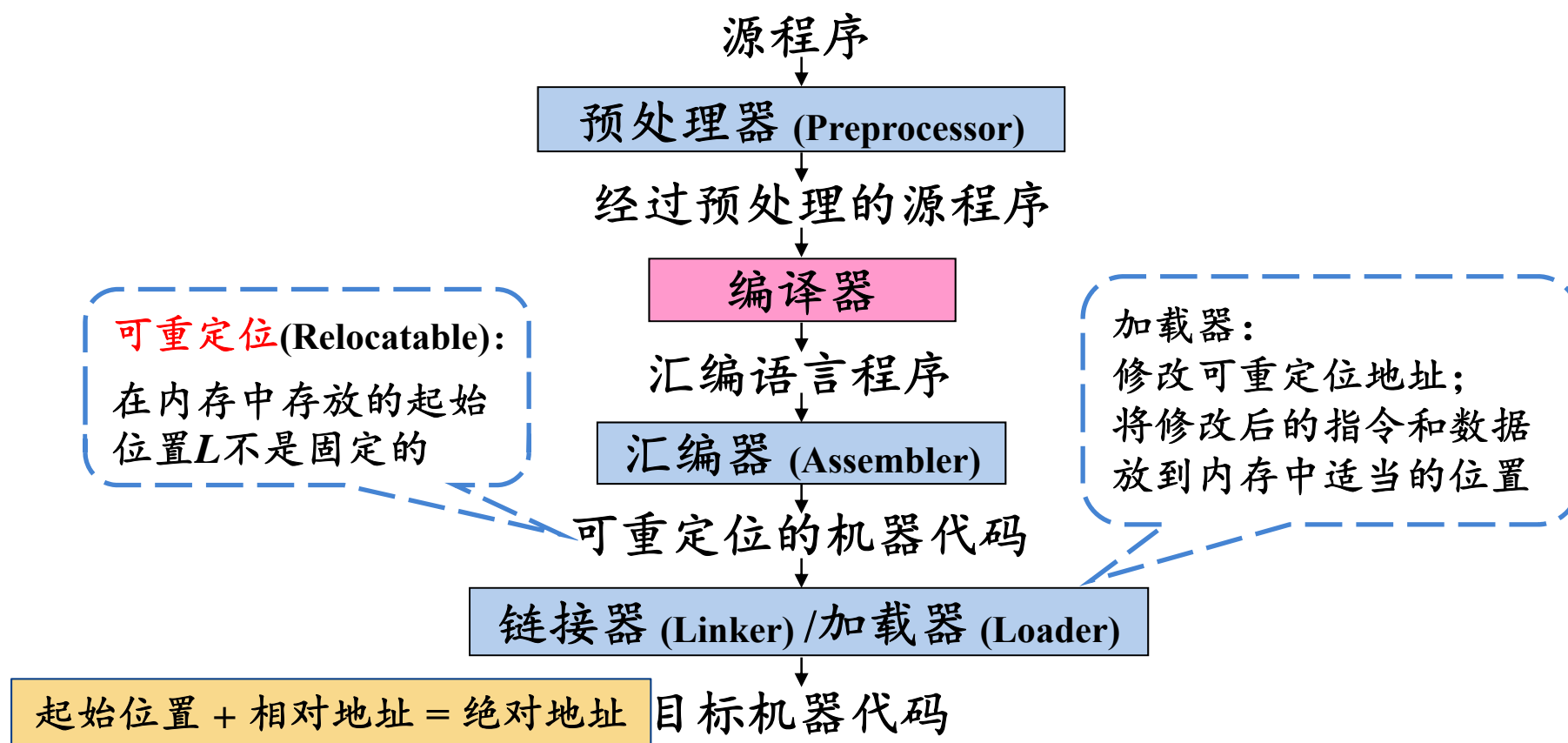


编译：将高级语言翻译成汇编语言或机器语言的过程
源语言 目标语言

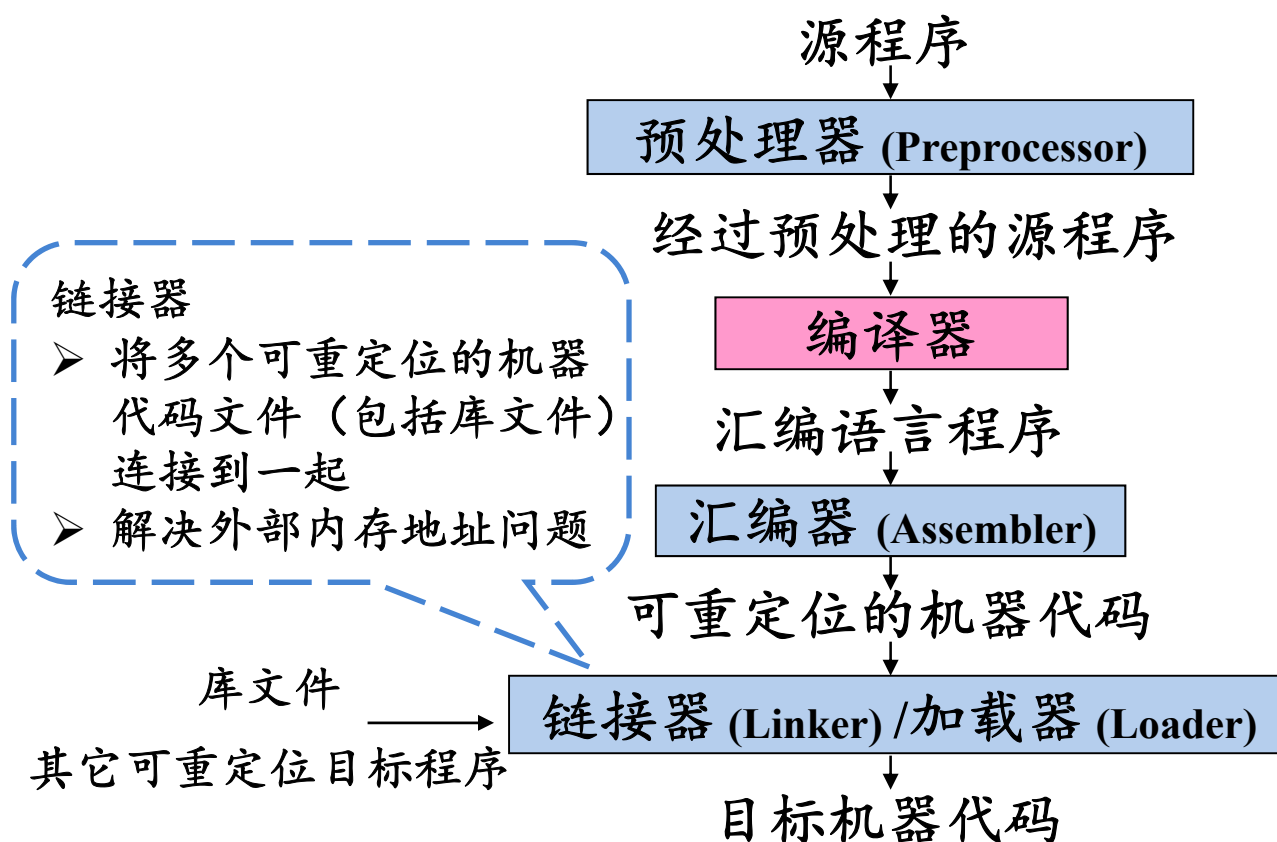
编译器在语言处理系统中的位置



编译器在语言处理系统中的位置



编译器在语言处理系统中的位置





提纲

1.1 什么是编译

1.2 编译系统的结构

1.3 为什么要学习编译原理

1.4 编译技术的应用

1.2 编译系统的结构

```
#include<stdio.h>
int main()
{
    int a,b,ge,shi,bai,m,n,i,number;
    printf("请输入(输入完两个数按一次回车键): \n");
    scanf("%d %d",&a,&b);
    while(a!=0,b!=0)
        scanf("%d %d",&a,&b);
    if(a>=100,b<=999)
        if(a>b)
            m=b,n=a;
        else
            m=a,n=b;
    for(i=m;i<=n;i++)
    {
        bai=i/100;
        shi=(i%100)/10;
        ge=i%10;
        if(i==bai*bai*bai+shi*shi*shi+ge*ge*ge)
            printf("%5d",i);
        number++;
        printf("\n");
    }
    if(number==0)
        printf("no\n");
}
```

高级语言程序

编译器

汇编语言程序/机器语言程序

机器是如何自动翻译的？

```
59     call sc
60     mov al, [si - 1]
61     call sc
62     mov al, [si]
63     call sc
64     mov al, ' '
65     call sc
66     ret
67 storechr endp
68 ;清屏, 无入口参数
69 clearcrt proc near
70     mov ax, 0600h
71     mov bh, 07h
72     mov cx, 00h
73     mov dx, 184fh
74     int 10h
```

```
84     mov ah, 02h
85     mov dl, cr
86     int 21h
87     mov ah, 02h
88     mov dl, lf
89     int 21h
90     .untilcxz
91     ret
92 nextlien endp
93     end
94
```

人工英汉翻译的例子

在房间里，他用锤子砸了一扇窗户。

[In the room], he broke a window <with a hammer>

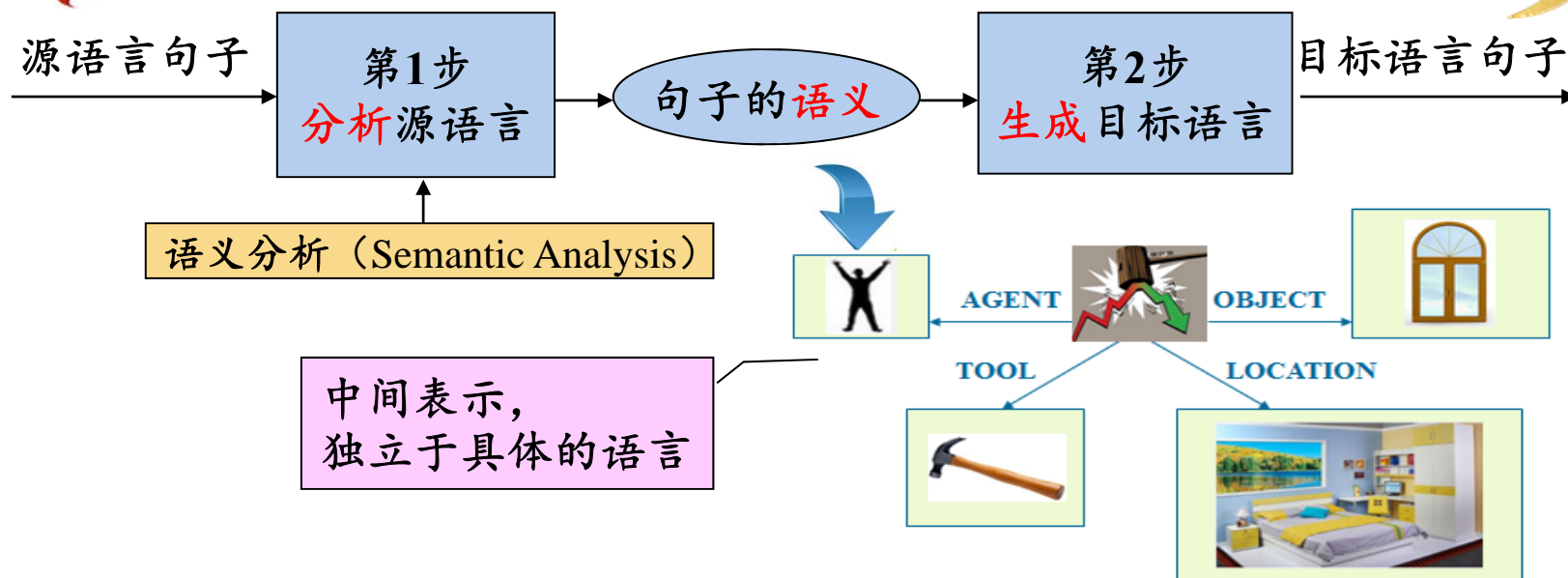
状语

主语

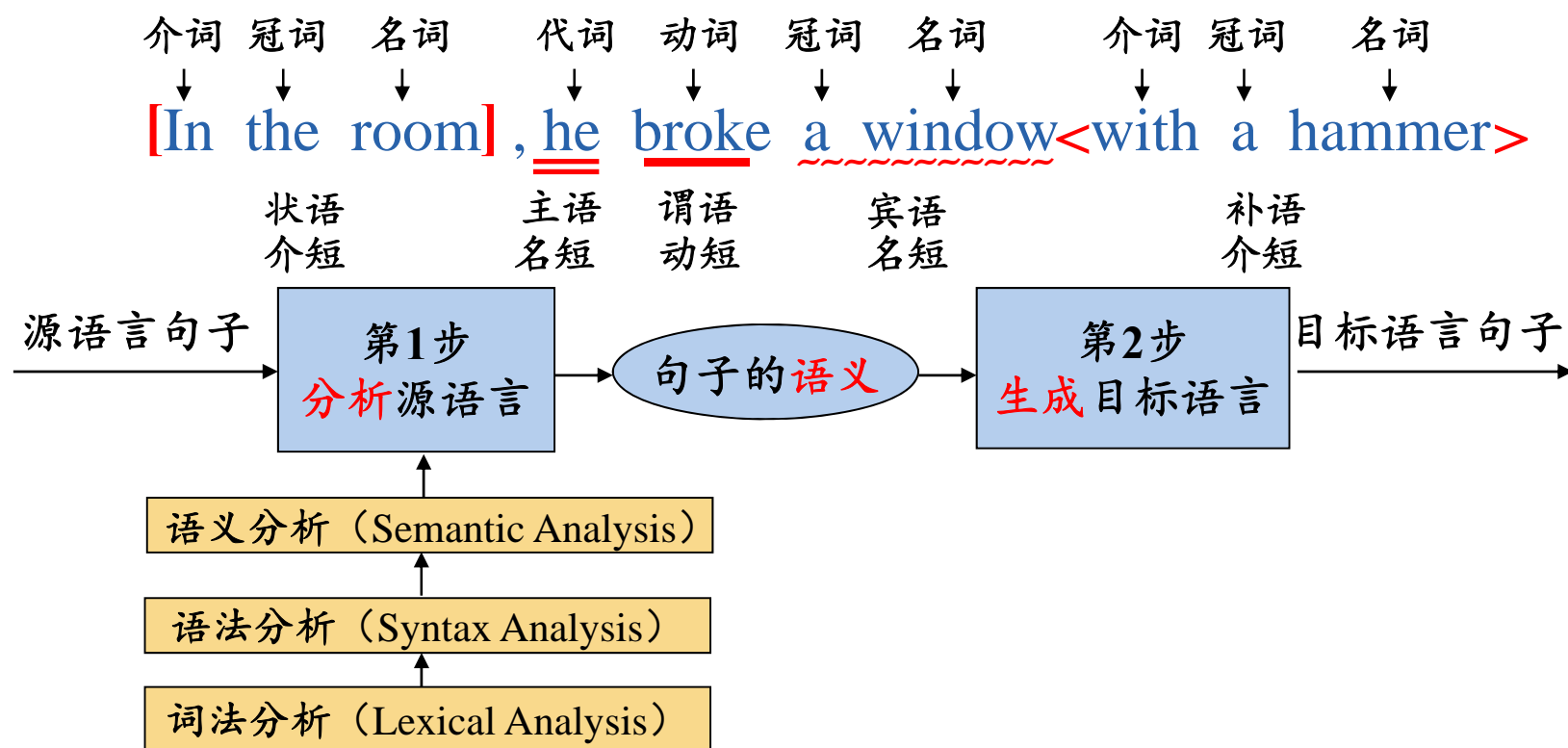
谓语

宾语

补语



人工英汉翻译的例子



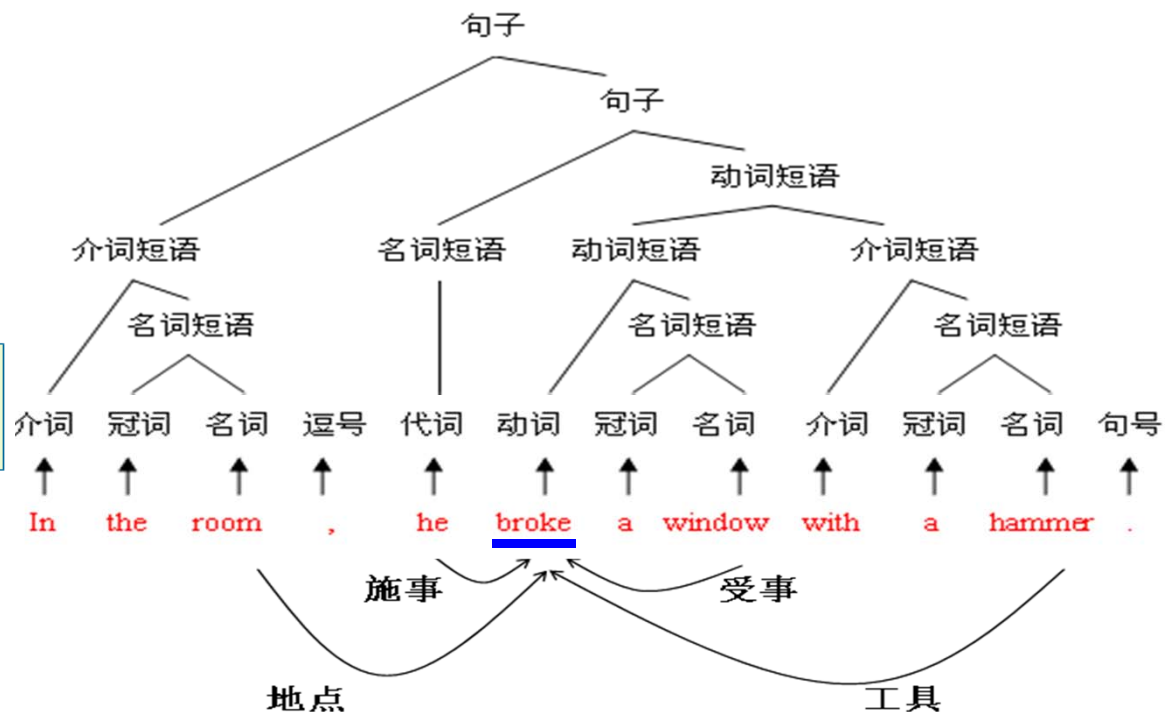
人工英汉翻译的例子

In the room , he broke a window with a hammer

➤ 词法分析

➤ 语法分析

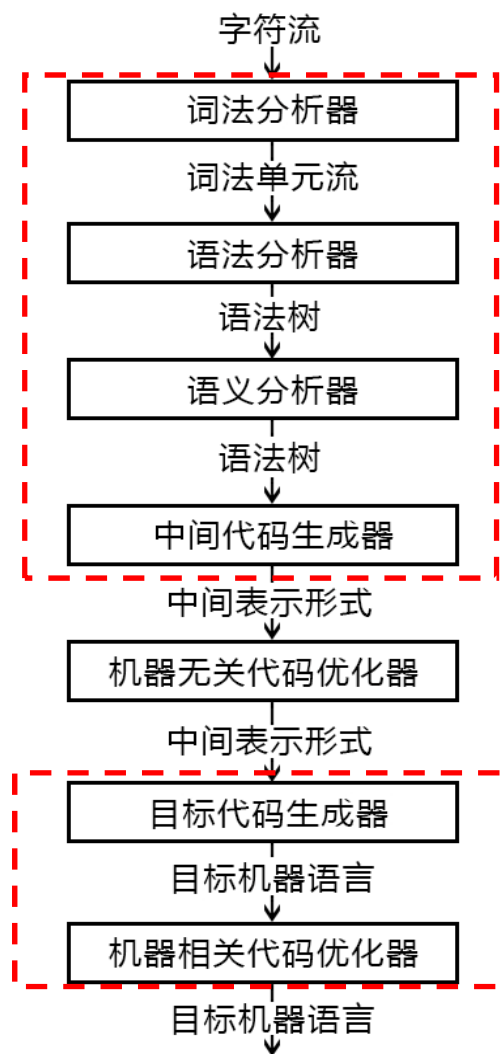
➤ 语义分析



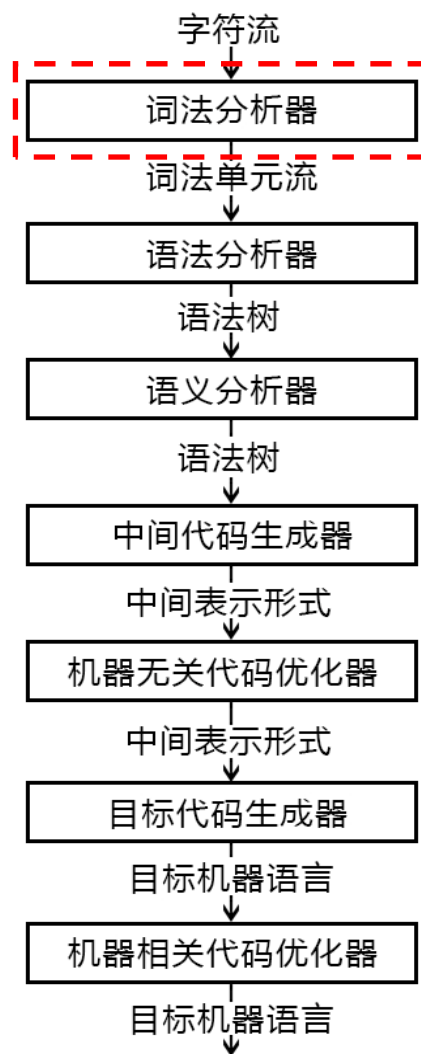
编译器的结构

分析部分/
前端(front end):
与源语言相关

综合部分/
后端(back end):
与目标语言相关



编译器的结构



词法分析/扫描(Scanning)

➤ 词法分析的主要任务

从左向右逐行扫描源程序的字符，识别出各个单词，确定单词的类型。

将识别出的单词转换成统一的机内表示——词法单元(token)形式

token: <种别码, 属性值>

	单词类型	种别	种别码
1	关键字	program、if、else、then、...	一词一码
2	标识符	变量名、数组名、记录名、过程名、...	多词一码
3	常量	整型、浮点型、字符型、布尔型、...	一型一码
4	运算符	算术 (+ - * / ++ --) 关系 (> < == != >= <=) 逻辑 (& ~)	一词一码 或 一型一码
5	界限符	; () = { } ...	一词一码

例：词法分析后得到的token序列

➤ 输入 while(value!=100){num++;}

➤ 输出

1	while	< WHILE ,	-	>
2	(< SLP ,	-	>
3	value	< IDN ,	value	>
4	!=	< NE ,	-	>
5	100	< CONST ,	100	>
6)	< SRP ,	-	>
7	{	< LP ,	-	>
8	num	< IDN ,	num	>
9	++	< INC ,	-	>
10	;	< SEMI ,	-	>
11	}	< RP ,	-	>

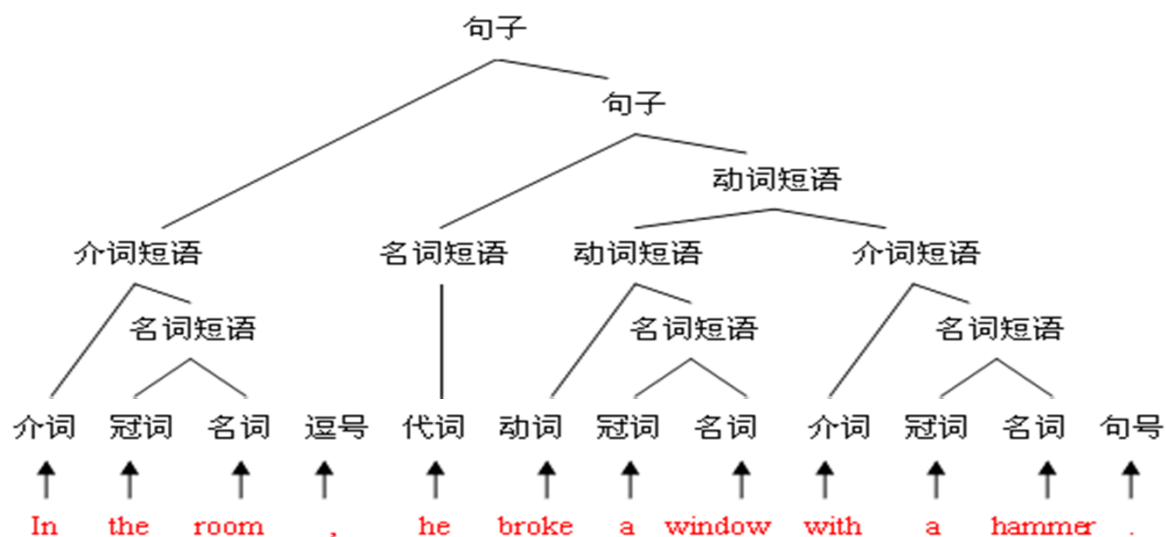
如何实现
词法分析器？

编译器的结构



语法分析 (*parsing*)

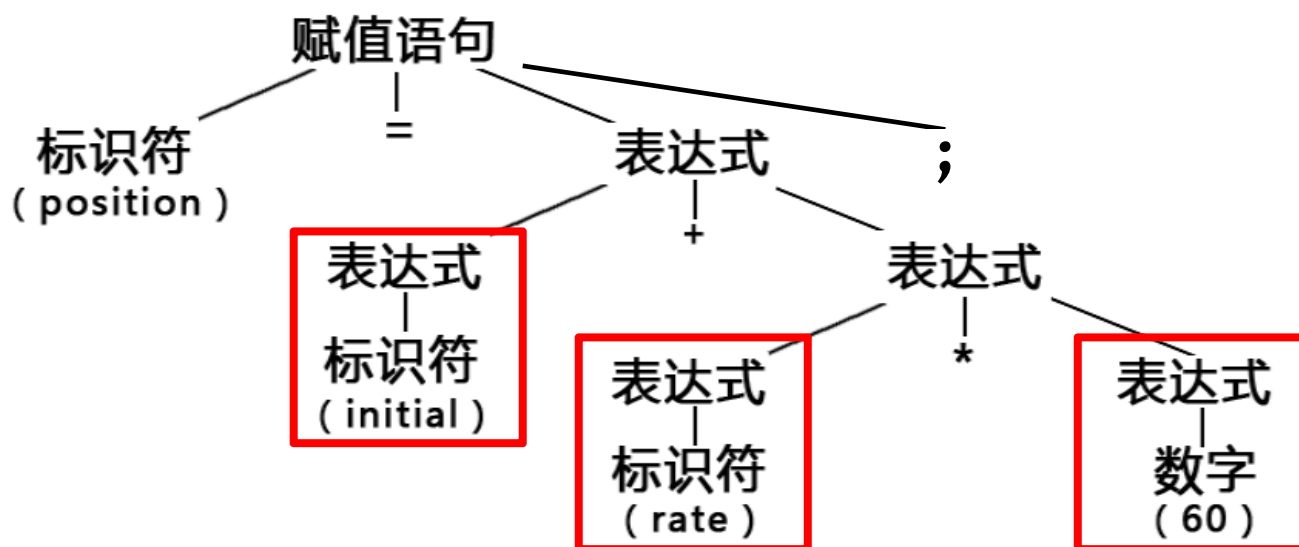
- 语法分析器(parser)从词法分析器输出的token序列中识别出各类短语, 并构造语法分析树(parse tree)
- 语法分析树描述了句子的语法结构



例1：赋值语句的分析树

position = initial + rate * 60 ;

$\langle \text{id, position} \rangle \Rightarrow \langle \text{id, initial} \rangle \langle + \rangle \langle \text{id, rate} \rangle \langle * \rangle \langle \text{num, 60} \rangle \langle ; \rangle$



例2：变量声明语句的分析树

➤ 文法：

$\langle D \rangle \rightarrow \langle T \rangle \langle IDS \rangle ;$

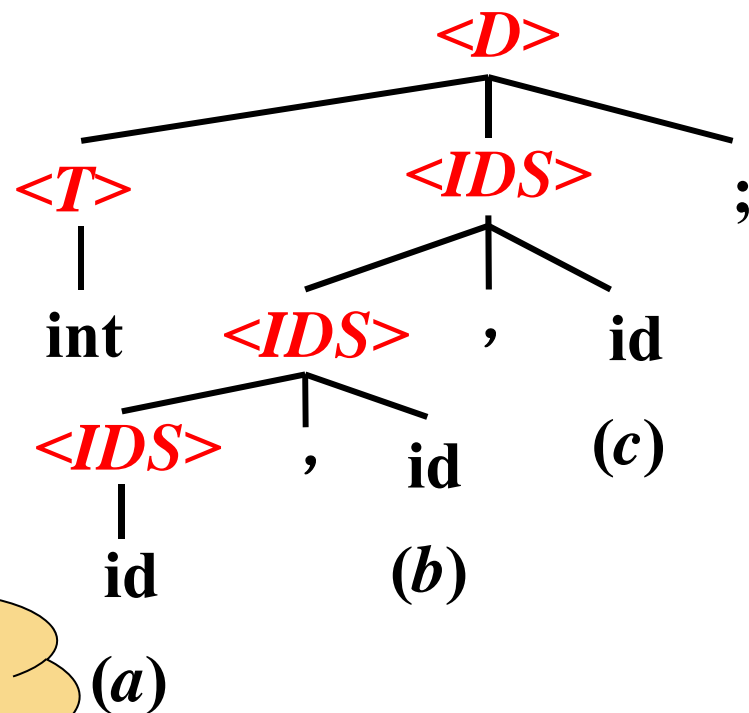
$\langle T \rangle \rightarrow \text{int} \mid \text{real} \mid \text{char} \mid \text{bool}$

$\langle IDS \rangle \rightarrow \text{id} \mid \langle IDS \rangle, \text{id}$

➤ 输入：

$\text{int } a, b, c ;$

如何根据语法规则为
输入句子构造分析树？



编译器的结构



语义分析的主要任务

- 收集标识符的属性信息

- 种属 (Kind)

- 简单变量、复合变量（数组、记录、...）、过程、...

语义分析的主要任务

➤ 收集标识符的属性信息

➤ 种属 (Kind)

➤ 类型 (Type)

➤ 整型、实型、字符型、布尔型、指针型、...

语义分析的主要任务

➤ 收集标识符的属性信息

➤ 种属 (Kind)

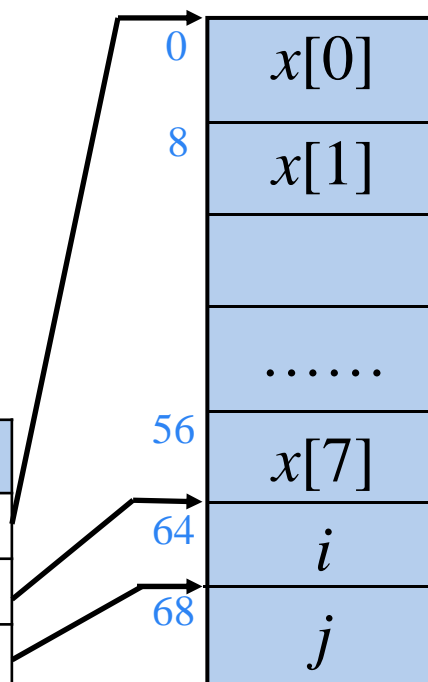
➤ 类型 (Type)

➤ 存储位置、长度

例:

```
begin
    real x[8];
    integer i,j;
    .....
end
```

名字	相对地址
<i>x</i>	0
<i>i</i>	64
<i>j</i>	68



语义分析的主要任务

- 收集标识符的属性信息
 - 种属 (Kind)
 - 类型 (Type)
 - 存储位置、长度
 - 值
 - 作用域
 - 参数和返回值信息
 - 参数个数、参数类型、参数传递方式、返回值类型、...

语义分析的主要任务

➤ 收集标识符的属性信息

➤ 种属 (Kind)

➤ 类型 (Type)

➤ 存储位置、长度

➤ 值

➤ 作用域

➤ 参数和返回值信息

符号表 (Symbol Table)

NAME	TYPE	KIND	VAL	ADDR
SIMPLE	整	简变		
SYMBLE	实	数组		
TABLE	字符	常数		
⋮	⋮	⋮	⋮	⋮

符号表是用于存放标识符的属性信息的数据结构

语义分析的主要任务

符号表中为什么要设计字符串表这样一种数据结构？

➤ 收集标识符的属性信息

➤ 种属 (Kind)

➤ 类型 (Type)

➤ 存储位置、长度

➤ 值

➤ 作用域

➤ 参数和返回值信息

符号表 (Symbol Table)

NAME	TYPE	KIND	VAL	ADDR
6	整	简变		
6	实	数组		
5	字符	常数		
⋮	⋮	⋮	⋮	⋮

字符串表

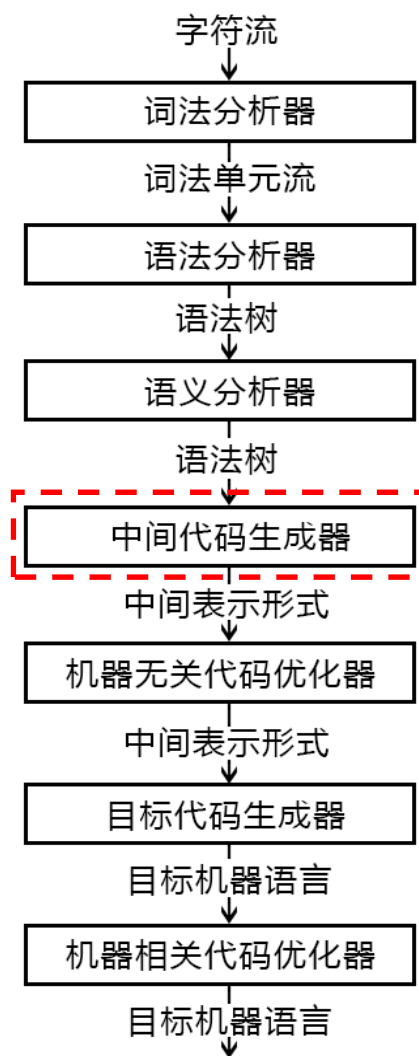


符号表是用于存放标识符的属性信息的数据结构

语义分析的主要任务

- 收集标识符的属性信息
- 语义检查
 - 变量或过程未经声明就使用
 - 变量或过程名重复声明
 - 运算分量类型不匹配
 - 操作符与操作数之间的类型不匹配
 - 数组下标不是整数
 - 对非数组变量使用数组访问操作符
 - 对非过程名使用过程调用操作符
 - 过程调用的参数类型或数目不匹配
 - 函数返回类型有误

编译器的结构



常用的中间表示形式

➤ 三地址码 (Three-address Code)

- 三地址码由 类似于汇编语言 的指令序列组成，
每个指令 最多有三个操作数(operand)

➤ 语法结构树/语法树 (Syntax Trees)

常用的三地址指令

序号	指令类型	指令形式
1	赋值指令	$x = y \text{ op } z$ $x = \text{op } y$
2	复制指令	$x = y$
3	条件跳转	if $x \text{ relop } y$ goto n
4	非条件跳转	goto n
5	参数传递	param x
6	过程调用 函数调用	call p, n $y = \text{call } p, n$
7	过程返回	return x
8	数组引用	$x = y[i]$
9	数组赋值	$x[i] = y$
10	地址及 指针操作	$x = \& y$ $x = * y$ $*x = y$

地址可以具有如下形式之一

- 源程序中的名字 (name)
- 常量 (constant)
- 编译器生成的临时变量 (temporary)

三地址指令的表示

➤ 四元式 (Quadruples)

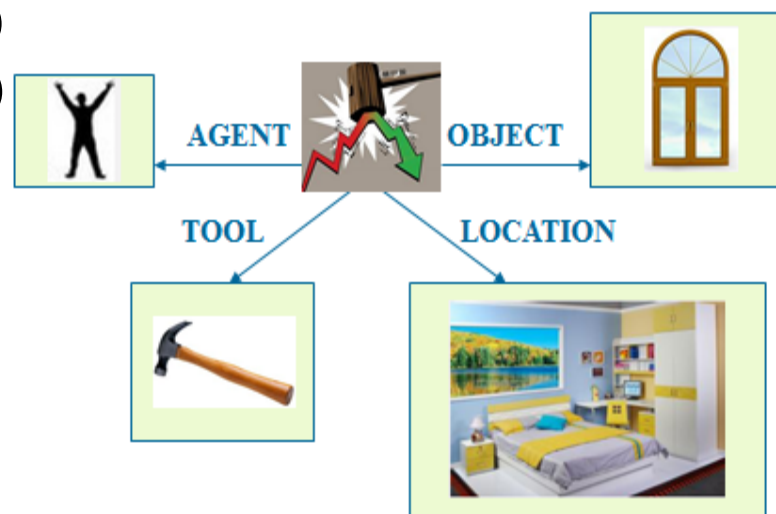
➤ (op, y, z, x)

➤ 三元式 (Triples)

➤ 间接三元式 (Indirect triples)

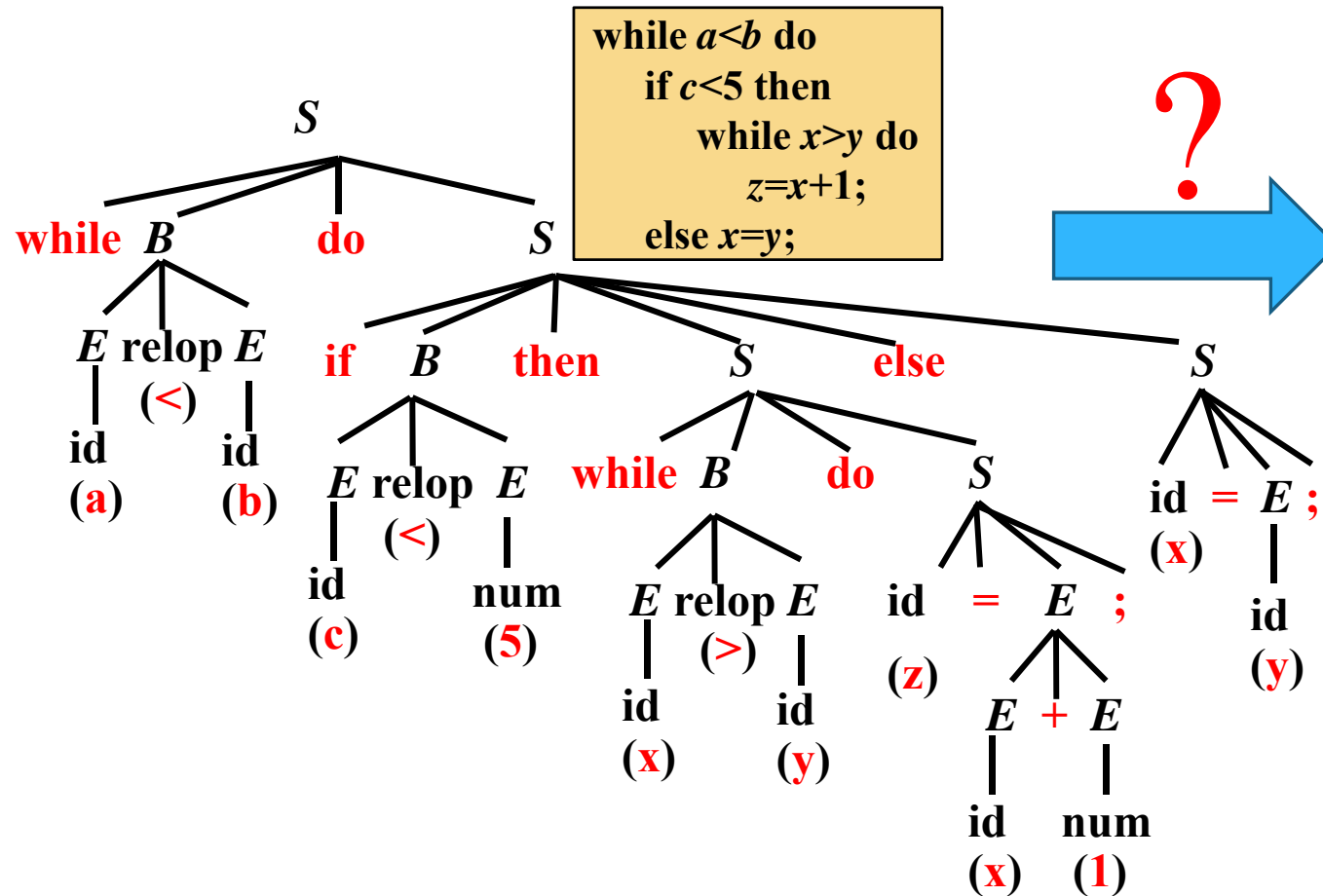
三地址指令的四元式表示

- $x = y \text{ op } z$ (**op** , y , z , x)
- $x = \text{op } y$ (**op** , y , $_$, x)
- $x = y$ (**=** , y , $_$, x)
- if $x \text{ relop } y$ goto n (**relop** , x , y , n)
- goto n (**goto** , $_$, $_$, n)
- param x (**param** , $_$, $_$, x)
- $y = \text{call } p, n$ (**call** , p , n , y)
- return x (**return** , $_$, $_$, x)
- $x = y[i]$ (**=[]** , y , i , x)
- $x[i] = y$ (**[]=** , y , x , i)
- $x = \&y$ (**&** , y , $_$, x)
- $x = *y$ (**=*** , y , $_$, x)
- $*x = y$ (***=** , y , $_$, x)



三地址指令序列唯一确定了
运算完成的顺序

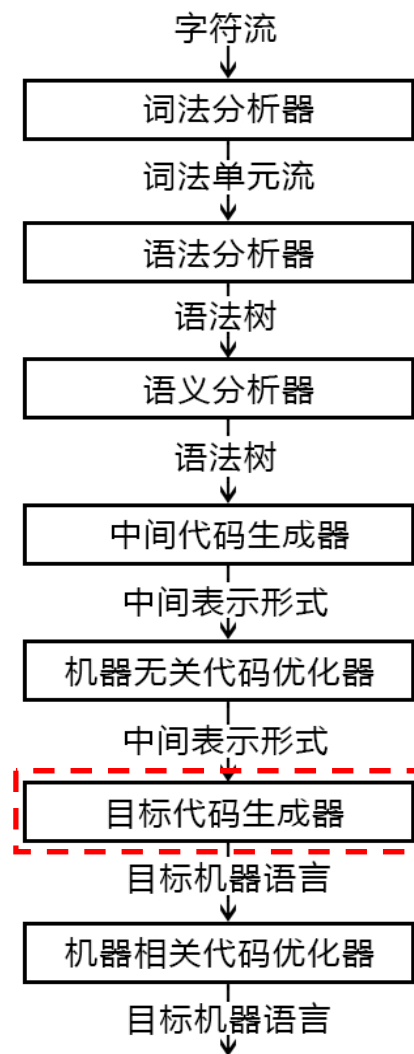
中间代码生成的例子



100: ($j <$, a , b , 102)
 101: (j , - , - , 112)
 102: ($j <$, c , 5 , 104)
 103: (j , - , - , 110)
 104: ($j >$, x , y , 106)
 105: (j , - , - , 100)
 106: ($+$, x , 1 , t_1)
 107: ($=$, t_1 , - , z)
 108: (j , - , - , 104)
 109: (j , - , - , 100)
 110: ($=$, y , - , x)
 111: (j , - , - , 100)
 112:

编译器的结构

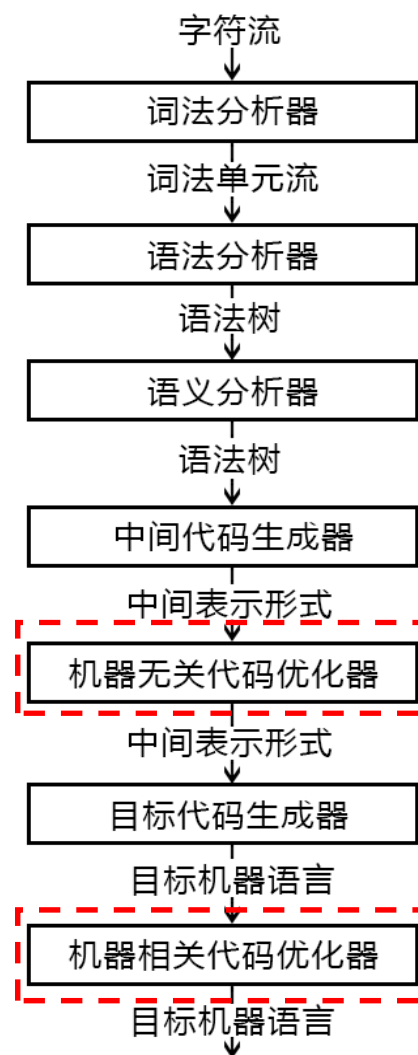
- 目标代码生成以源程序的**中间表示形式**作为输入，并把它映射到**目标语言**
- 目标代码生成的一个重要任务是为程序中使用的变量**合理分配寄存器**



编译器的结构

➤ 代码优化

- 为改进代码所进行的等价程序变换，使其运行得更快一些、占用空间更少一些，或者二者兼顾





提纲

1.1 什么是编译

1.2 编译系统的结构

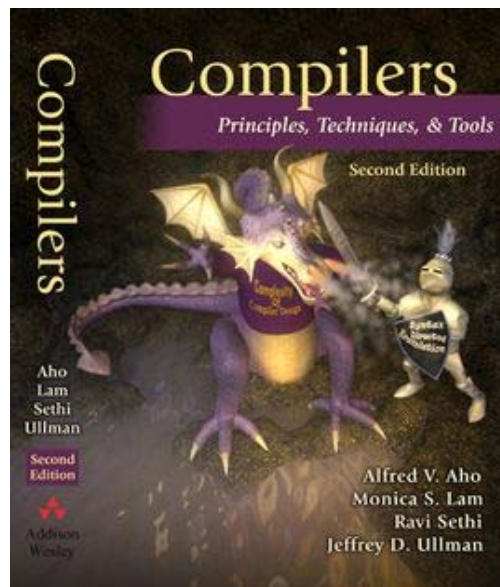
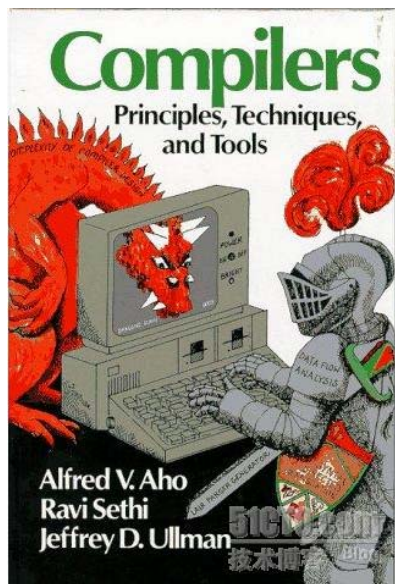
1.3 为什么要学习编译原理

1.4 编译技术的应用

1.3 为什么要学习编译原理

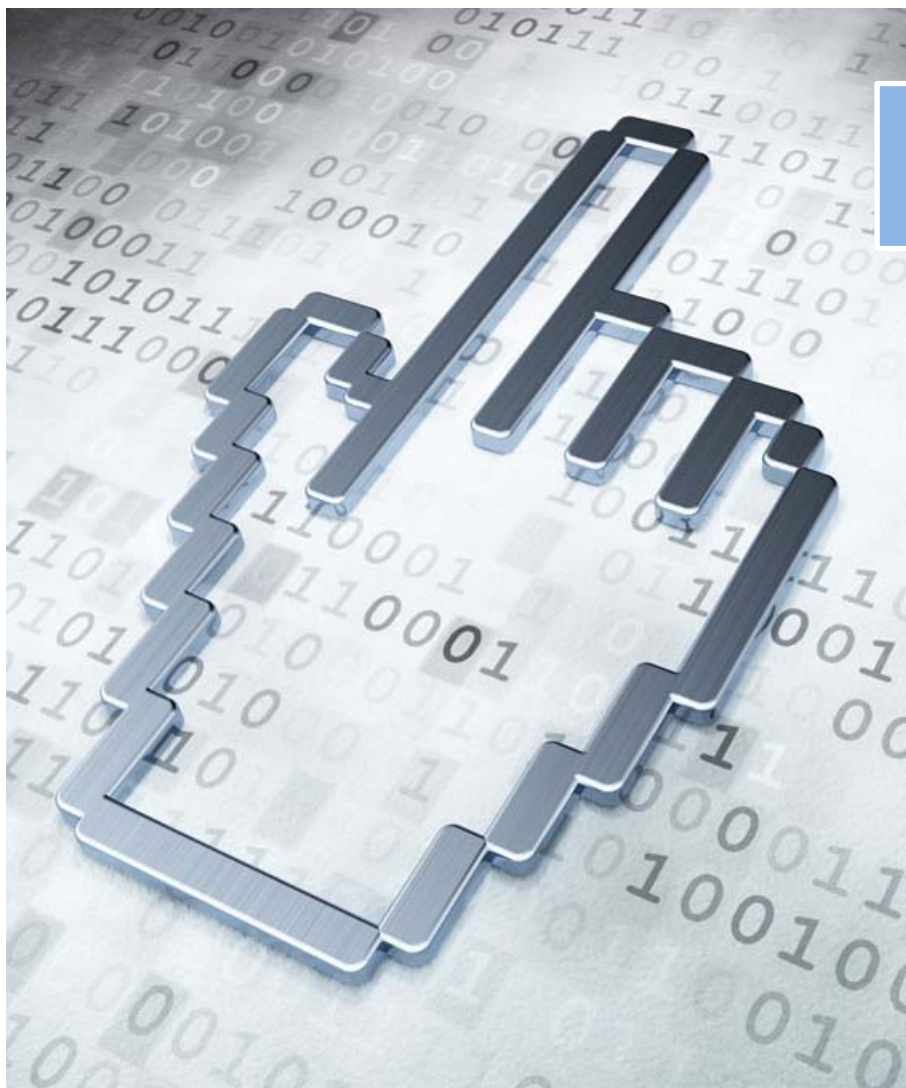
编写编译器的原理和技术具有十分普遍的意义，以至于在每个计算机科学家的研究生涯中，本课程中的原理和技术都会反复用到。

——Alfred V.Aho



通过本课程的学习

- 更深刻地理解高级语言程序的内部运行机制
- 教给我们如何严谨地去思考、编写程序
- 编译原理蕴含了计算机科学解决问题的思路和方法，即“形式化→自动化”
- 所涉及的理论和方法在很多领域都会被用到
 - 自然语言处理、模式识别、人工智能、.....
- 很多应用软件都会用到编译技术



提纲

1.1 什么是编译

1.2 编译系统的结构

1.3 为什么要学习编译原理

1.4 编译技术的应用

1.4 编译技术的应用

- 结构化编辑器 (Structure editors)
 - 引导用户在语言的语法约束下编制程序
 - 能自动地提供关键字和与其匹配的关键字

1.4 编译技术的应用

- 结构化编辑器 (Structure editors)
- 智能打印机 (Pretty printers)
 - 对程序进行分析，打印出结构清晰的程序
 - 注释以一种特殊的字体打印
 - 根据各个语句在程序的层次结构中的嵌套深度进行缩进

1.4 编译技术的应用

- 结构化编辑器 (Structure editors)
- 智能打印机 (Pretty printers)
- 静态检测器 (Static checkers)
 - 静态定位程序中的错误
 - 释放空指针或已释放过的指针
 - 检测出程序中永远不能被执行的语句

1.4 编译技术的应用

- 结构化编辑器 (Structure editors)
- 智能打印机 (Pretty printers)
- 静态检测器 (Static checkers)
- 文本格式器 (Text formatters)
 - 文本格式器处理的字符流中除了需要排版输出的字符以外，还包含一些用来说明字符流中的段落、图表或者上标和下标等数学结构的命令

1.4 编译技术的应用

- 结构化编辑器 (Structure editors)
- 智能打印机 (Pretty printers)
- 静态检测器 (Static checkers)
- 文本格式器 (Text formatters)
- 数据库查询解释器 (Database Query Interpreters)
 - 数据库查询语句由包含了关系和布尔运算的谓词组成。查询解释器把这些谓词翻译成数据库命令，在数据库中查询满足条件的记录。

1.4 编译技术的应用

- 结构化编辑器 (Structure editors)
- 智能打印机 (Pretty printers)
- 静态检测器 (Static checkers)
- 文本格式器 (Text formatters)
- 数据库查询解释器 (Database Query Interpreters)
- 高级语言的翻译工具



本章小结

什么是编译

编译系统的结构

为什么要学习编译原理

编译技术的应用

课程主要内容

- 1. 绪论 (2学时)
- 2. 语言及其文法 (2学时)
- 3. 词法分析 (4学时)
- 4. 语法分析 (9学时)
- 5. 语法制导翻译 (6学时)
- 6. 中间代码生成 (7学时)
- 7. 运行时的存贮组织 (3学时)
- 8. 代码优化 (5学时)
- 9. 代码生成 (2学时)



结束

