

北京大学计算机学院
2024年秋季学期 《编译原理》



《编译原理》课程简介

Course Introduction



基本信息

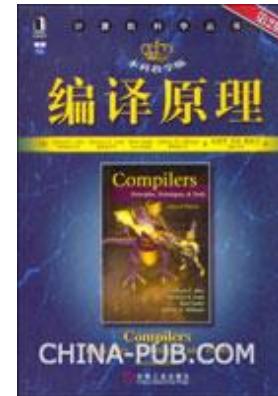
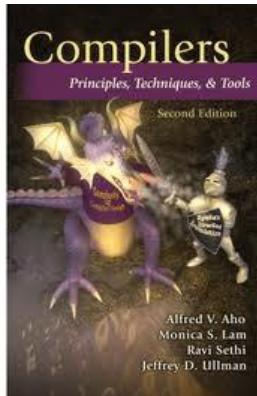
- 课程网页
 - 北京大学教学网 (<http://course.pku.edu.cn>)
- 主讲教师：张路、刘先华、王迪
 - Office: 理科1号楼1542
 - Tel: 13366302060
 - Email: zhanglucs@pku.edu.cn
- 助教
 - 刘晨、肖越易
- 交流：微信群



教材及参考书

教材：龙书第2版

- Alfred V. Aho, Monica S. Lam, Ravi Sethi, Jeffrey D. Ullman, **Compilers: Principles, Techniques, & Tools** (2nd Edition), Addison-Wesley, 2007.
- 中文翻译版或本科教学版：《编译原理(第2版)》 赵建华等译, 机械工业出版社, 2009.



□ 参考书

- 《编译原理》孙家騤 编著, 北京大学出版社, 2008年。
- 【虎书】Andrew Appel, **Modern Compiler Implementation in Java** (2nd Edition), Cambridge University Press, 2002
- 陈火旺等, 程序设计语言编译原理. 国防工业出版社, 1984。
- 张素琴, 吕映芝等, 编译原理 (第2版). 清华大学出版社, 2005。



课程目标

- 掌握编译的基本理论
- 掌握常用的编译技术
 - 词法分析、正则表达式、有限自动机
 - 上下文无关文法
 - 自底向上和自顶向下的语法分析
 - 基本的编译优化技术
- 掌握编译器的构造技术
 - 编译器的手工构造
 - 编译器生成器的实现



课程定位

- 计算科学技术专业和软件工程专业核心课程
- 强调**形式化描述技术**和**自动生成技术**
 - 兼顾语言的描述方法、设计、应用
 - 能形式化就能自动化
- 强调对编译原理和技术的**宏观理解**
 - 不把注意力过于分散到具体的算法细节
 - 不过于偏向于任何源语言或目标机器
- 强调实际构造编译器的能力
- 进一步培养“**计算机思维能力**”



学习本课程的意义

- Alfred V. Aho [龙书作者]: “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，本书中的原理和技术都会反复用到。”
- 编译技术涉及的是一个比较适当的抽象层面上的数据变换（既抽象，又实际）。
 - 一些具体的表示和变换算法。
- 一个相当规模的系统的设计（含总体结构）。
- “自顶向下的方法”和“自底向上的方法”系统设计方法（思想、方法、实现全方位讨论）



评分标准

- 期末笔试 (30%)
 - 范围：全部内容
- 期中笔试 (20%)
 - 范围：第1-4章
- 书面作业 (10%)
- 实习 (40%)

- 出勤情况
 - 可能不定期点名（缺席一次扣一分）
- 注：成绩比例可能按需调整



课程要求

- **课堂：**按时上课、认真听讲，积极参与
 - 遇到问题欢迎随时打断老师
- **作业：**一定要认真独立完成
 - 不明白的地方要追根究底、做到举一反三
- **考试：**鼓励知识的灵活运用，不提倡死记硬背

- 讲课内容比较多，进展较快，理论性比较强
 - 平时如果不复习并加深理解，后面将听不懂
 - 然后就会导致心理上的畏惧感
 - 进而影响实习部分



其他

- 重修、旁听、留学生以及其他有特殊困难的同学
 - 课间或答疑时间找我交流一下情况

- 有特殊情况不能按时上课的同学请提前请假
 - 事后请假一律无效

北京大学计算机学院
2024年秋季学期 《编译原理》



第1章 编译器概述

Introduction to Compilers



主要内容

□ 编译的起源

- 程序设计语言的发展
- 编译的基本概念
- 编译器vs解释器

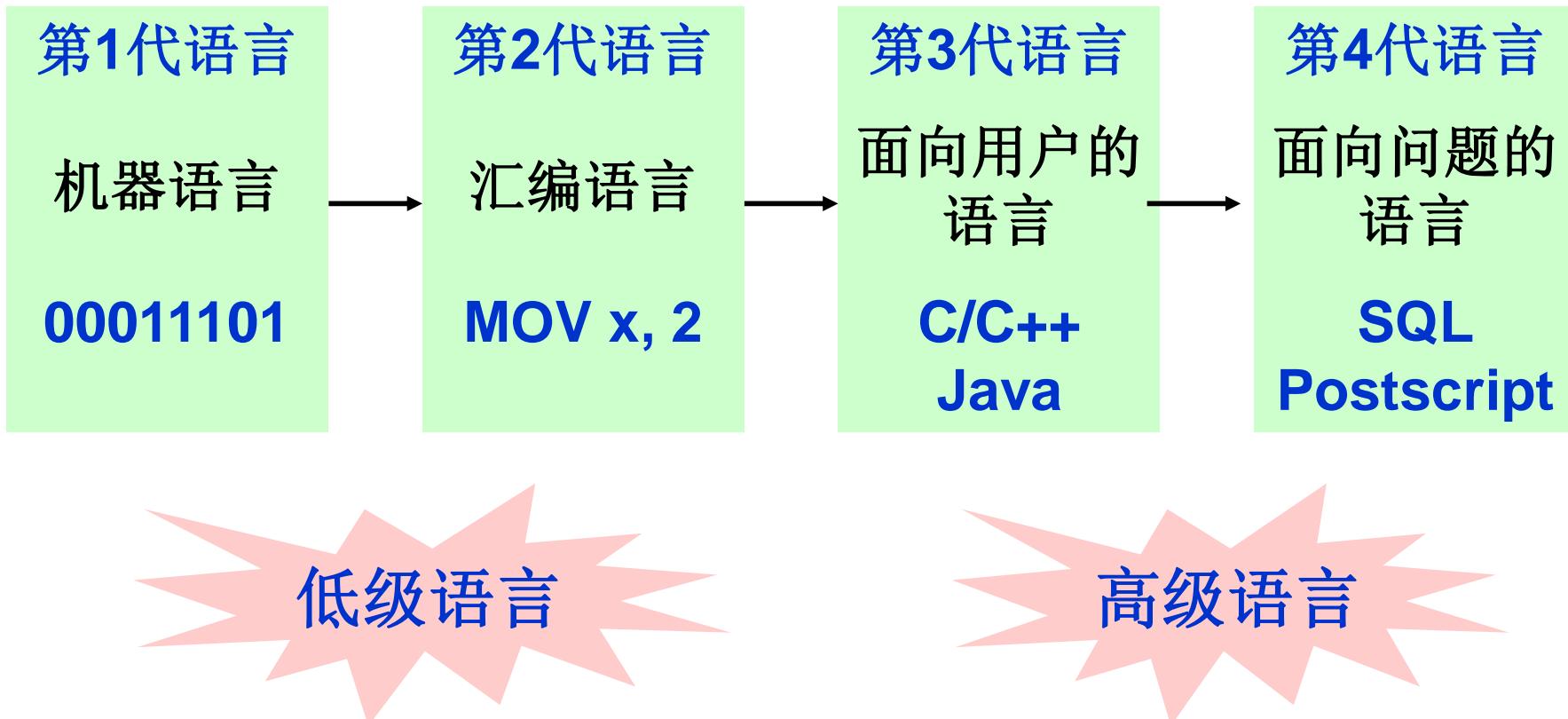
□ 编译器的结构

□ 编译技术的应用

□ 程序设计语言基础



程序设计语言的发展





低级语言 vs. 高级语言

□ 低级语言 (Low level Language)

- 字位码、机器语言、汇编语言
- 特点：与特定的机器有关，功效高，但使用复杂、繁琐、费时、易出错

□ 高级语言 (High level language)

- Fortran、Pascal、C/C++、Java 语言等
- 特点：不依赖具体机器，移植性好、对用户要求低、易使用、易维护等。

用高级语言编写的程序，计算机不能立即执行，必须通过一个“翻译程序”进行加工，转化为与其等价的机器语言程序，机器才能执行。

这种翻译程序，被称为“编译程序”。

最早的编译器



- 第一个编译器 (1952年)
 - A-0 System on UNIVAC I
 - 编写者: **Grace Murray Hopper**
 - 美国海军少将
 - COBOL发明人
 - 创造了术语“debug”
- 第一个完整的编译器 (1957年)
 - John Backus 领导的 **FORTRAN**编译器 (IBM)
- 第一个高级语言书写的编译器 (1962年)
 - **Lisp**编译器 (MIT, Tim Hart and Mike Levin)



如何构造第一个编译器？

- 在机器W上，如何为语言X构造一个用X编写的编译器？
- 假设机器W上已经存在某种语言Y的编译器
 - 首先用Y语言编写X语言的编译器，
 - 然后用该编译器编译用X语言编写X的编译器
- 假设机器W上不存在任何语言的编译器
 - **bootstrapping**



基本概念

□ 源程序 (Source Program)

- 用汇编语言或高级语言编写的程序称为源程序。

□ 目标程序 (Object Program)

- 用目标语言所表示的程序。
- 目标语言：可以是介于源语言和机器语言之间的“中间语言”，可以是某种机器的机器语言，也可以是某机器的汇编语言。

□ 翻译程序 (Translator)

- 将源程序转换为目标程序的程序称为翻译程序。
- 它是指各种语言的翻译器，包括汇编程序和编译程序，是汇编程序、编译程序以及各种变换程序的总称。



编译的基本功能

□ 编译（Compiling）就是把高级语言程序翻译成与其等价的目标语言程序的过程。

- 等价：针对所有相同输入会给出相同的输出
- 在应用程序与硬件体系结构之间的接口
- 应当对所生成的程序进行一定的优化





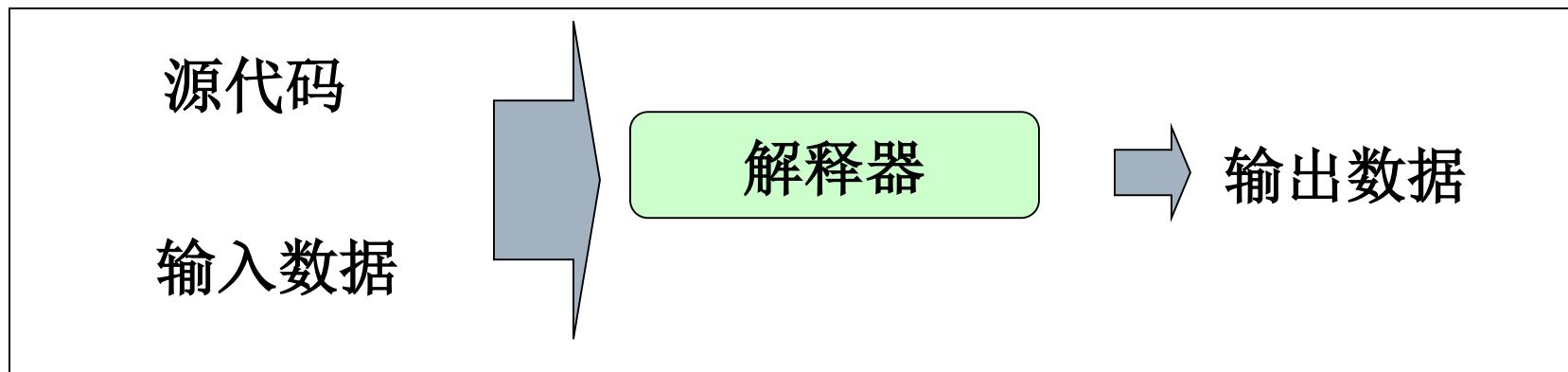
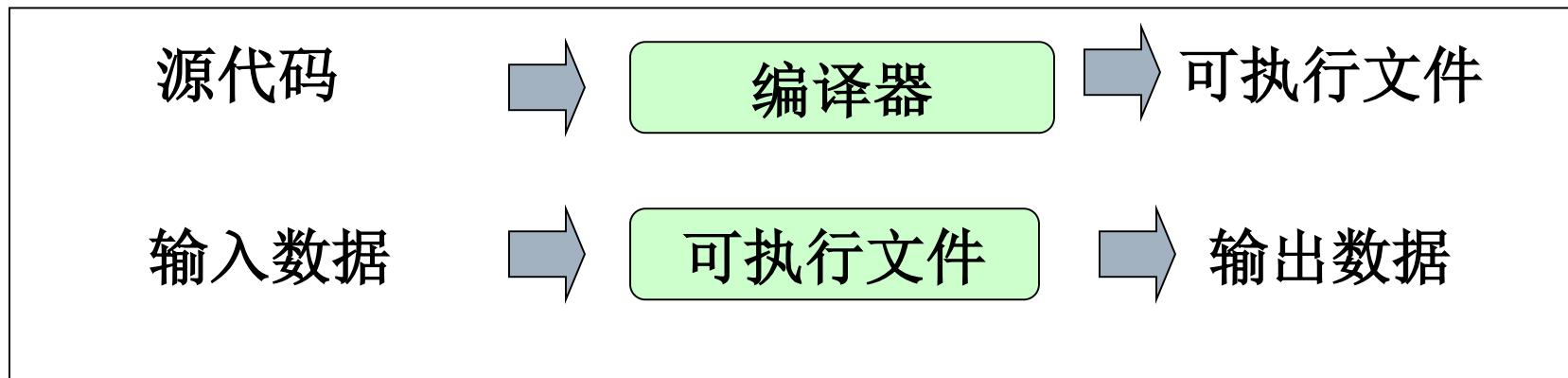
编译器 vs. 解释器 (1/5)

- 编译器 (Compiler) : 把 (人编写的) 源程序转化为 (机器可读的) 可执行程序
- 解释器 (Interpreter) : 在转换源程序的同时执行程序
 - 不生成目标代码，而是直接执行源程序所指定的运算
 - 解释器也需要对源程序进行词法，语法和语义分析，中间代码生成



编译器 vs. 解释器 (2/5)

□ 早期状态





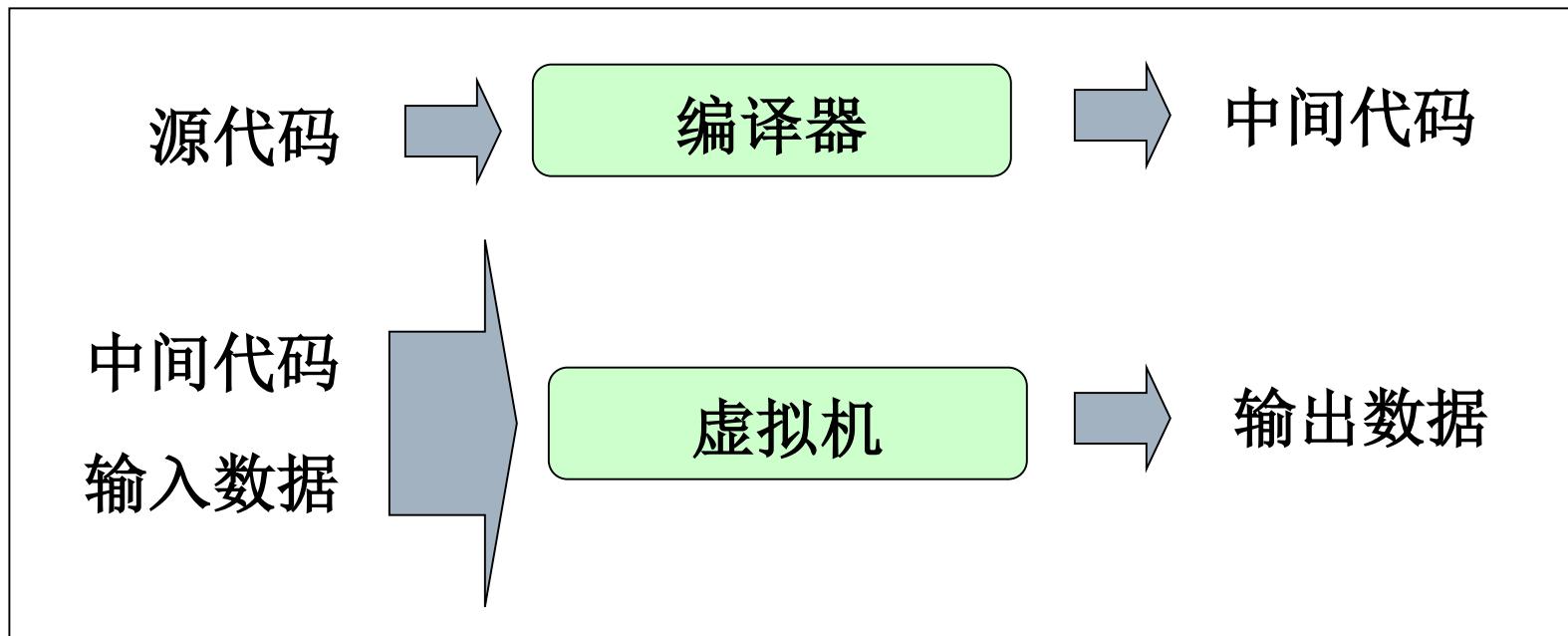
编译器 vs. 解释器 (3/5)

- 大多数语言通常属于二者之一：
 - 编译器：FORTRAN, COBOL, C, C++, Pascal, PL/1
 - 解释器：Lisp, BASIC, APL, Perl, Python, Smalltalk
- 但是：在现实中并不总是这样来区分的
 - 虚拟机 (e.g., Java)
 - 可执行代码的运行时链接
 - JIT (Just-in-time) 编译



编译器 vs. 解释器 (4/5)

- 实际使用中，在二者之间并没有明确的边界
 - 常见的情形是二者混合起来





编译器 vs. 解释器 (5/5)

编译器

□ 优点

- 占空间小
- 执行速度快

□ 缺点

- 处理速度慢
 - 已解决：分别编译
- 调试困难
 - 通过IDE的支持来改进

解释器

□ 优点

- 易于调试
- 开发快捷

□ 缺点

- 通常不适合大规模项目
- 占用较多空间
- 执行速度较慢
 - 解释器必须常驻内存



如何编译这个程序？

一个C语言源程序

```
void cal(float b,
         float c) {
    float a;

    a = b+c*60;
    printf("%f\n", a);
}
```

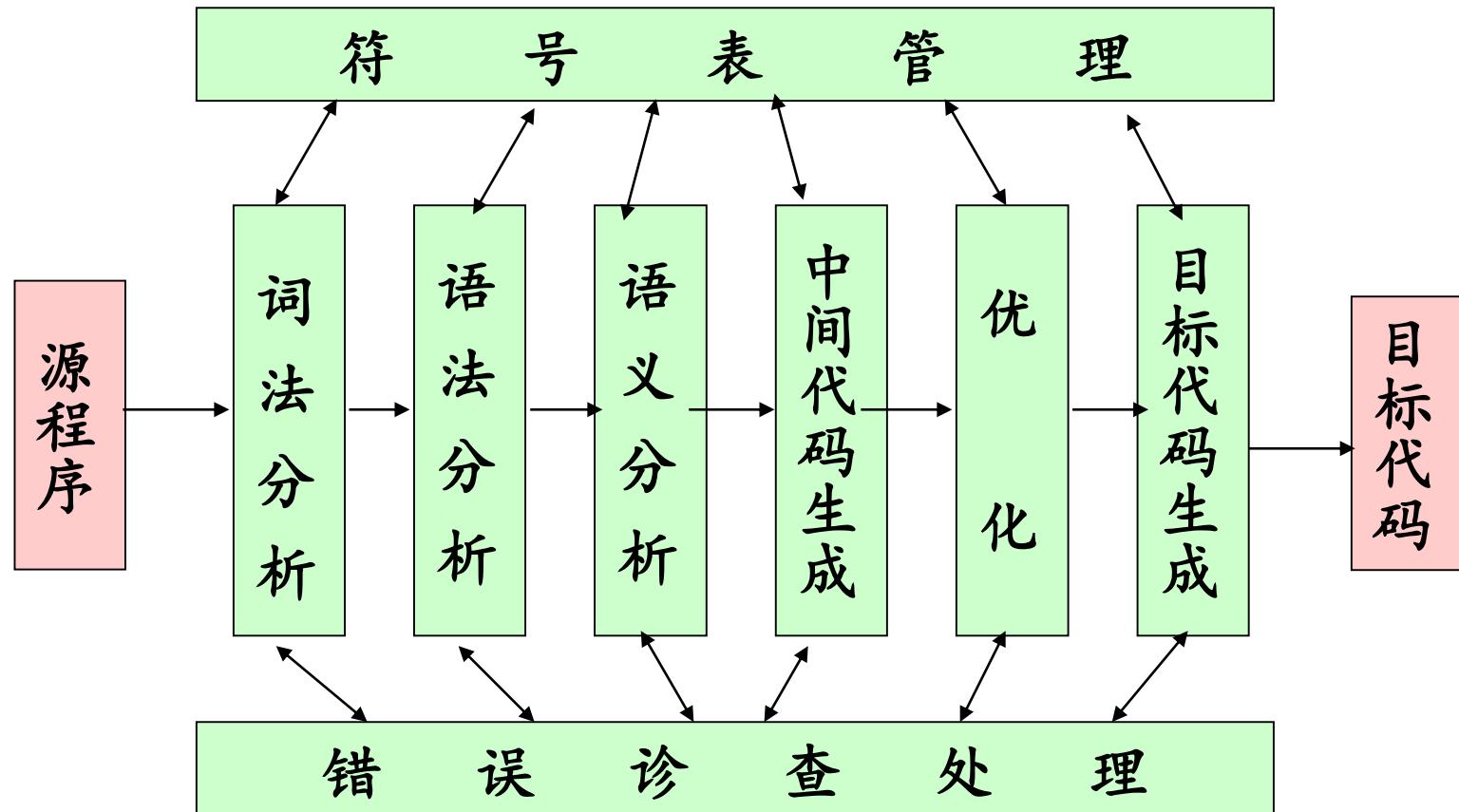


主要内容

- 编译的起源
 - 程序设计语言的发展
 - 编译的基本概念
 - 编译器vs解释器
- 编译器的结构
- 编译技术的应用
- 程序设计语言基础



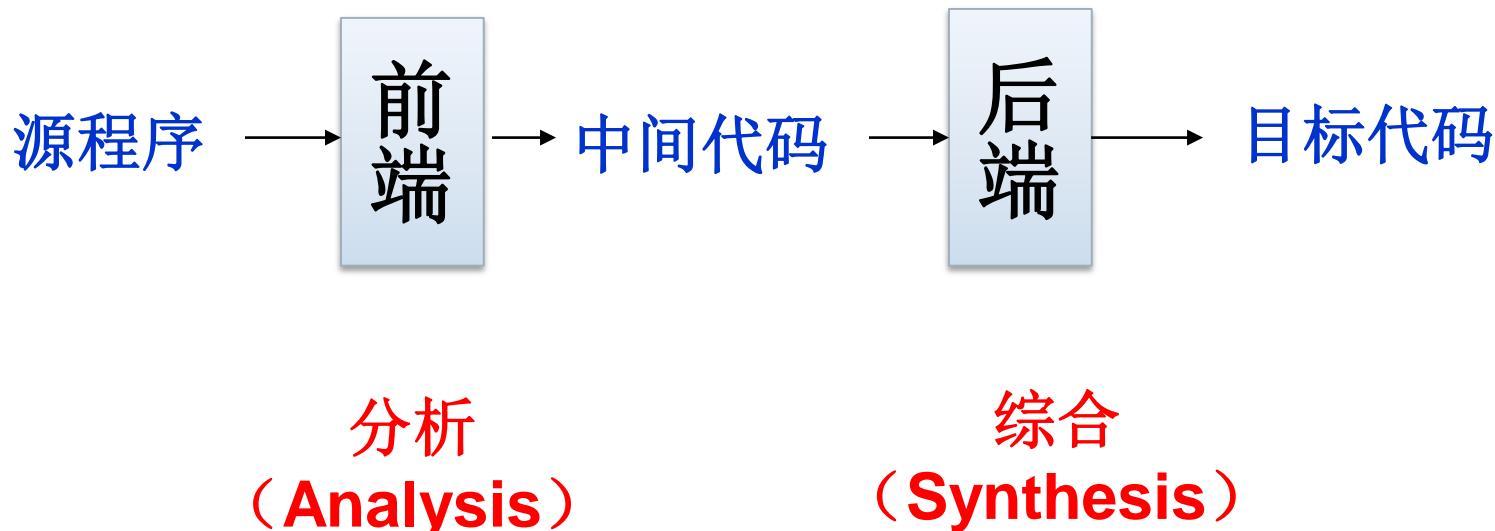
编译器的结构





编译程序（器）的组织

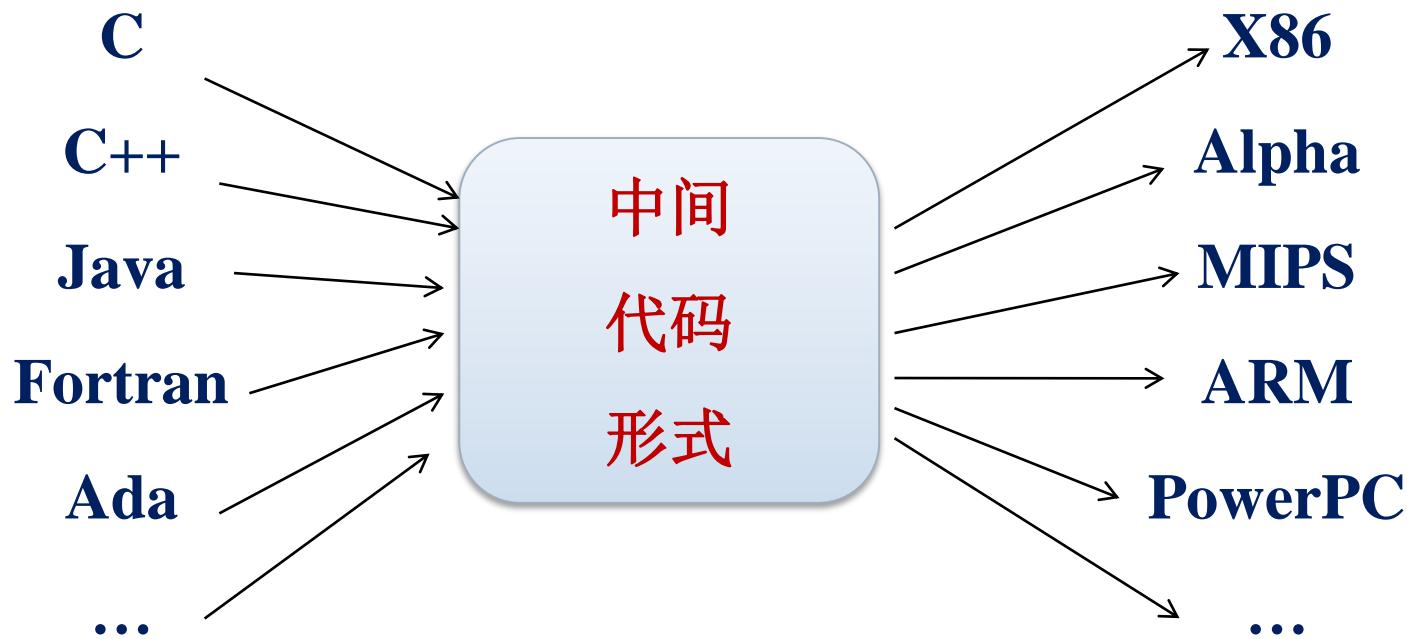
- 前端（frontend）：从源代码到中间代码的翻译，与具体机器无关
- 后端（backend）：把中间代码翻译为具体的机器指令（目标代码）





现代编译器（例如GCC）

- 多种语言（C, C++, Java, Fortran, Ada...）
- 多种目标机器（x86, Alpha, MIPS, ARM...）





编译器流程示例

回到这个C语言源程序

```
void cal(float b,
         float c) {
    float a;

    a = b+c*60;
    printf("%f\n", a);
}
```



编译器流程示例——词法分析

词法分析的结果是代表单词的token流

```
<保留字, void>
<标识符, cal>
<分隔符, (>
<保留字, float>
.....
<标识符, a>
<算符, =><标识符, b> <算符, +>
<标识符, c> <算符, *> <常数, 60>
.....
<字符串, "%f\n">
<分隔符, , >
<标识符, a>
.....
```

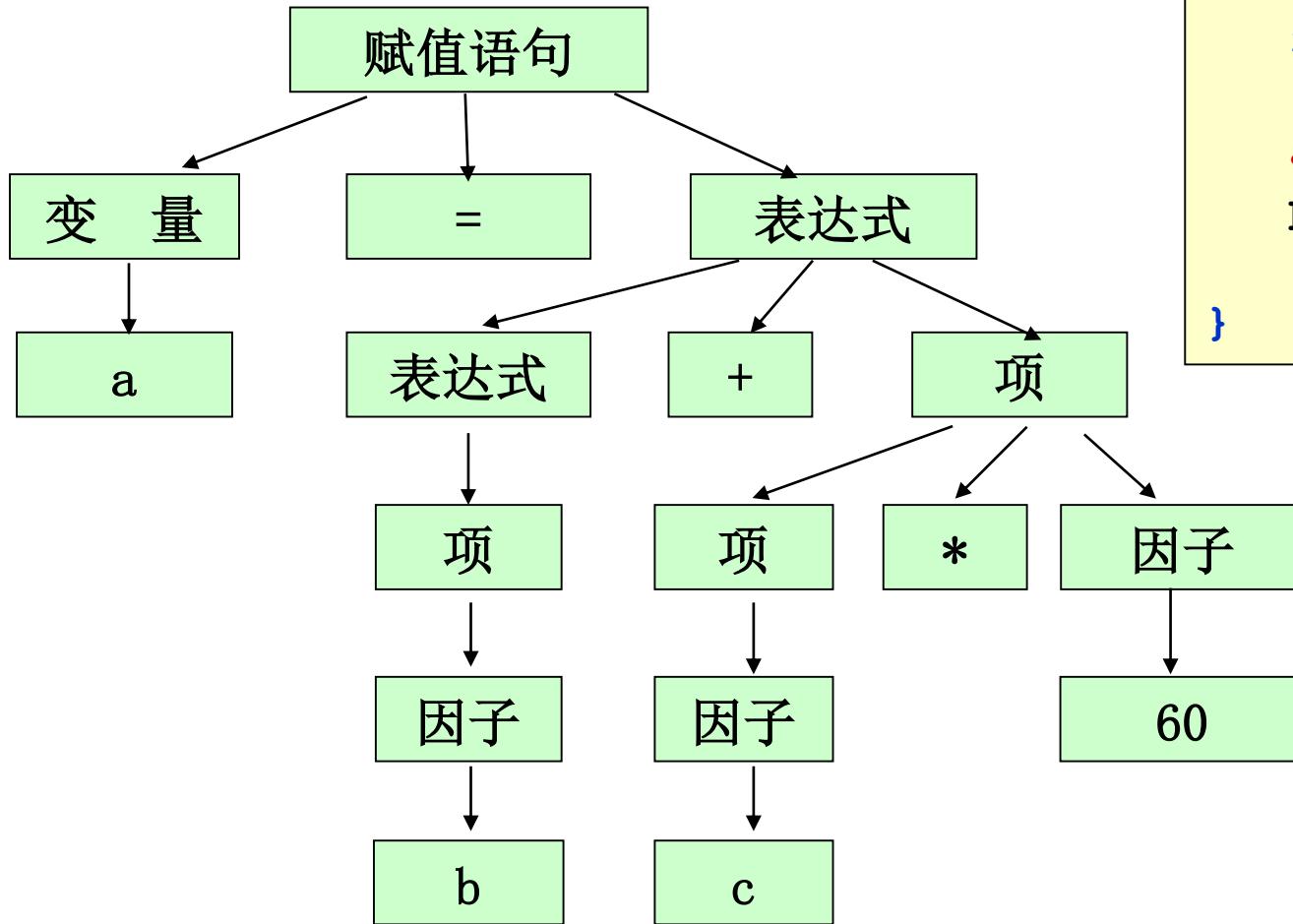
```
void cal(float b,
         float c) {
    float a;

    a = b+c*60;
    printf("%f\n",
           a);
}
```



编译器流程示例——语法分析

为程序中赋值语句生成的分析树

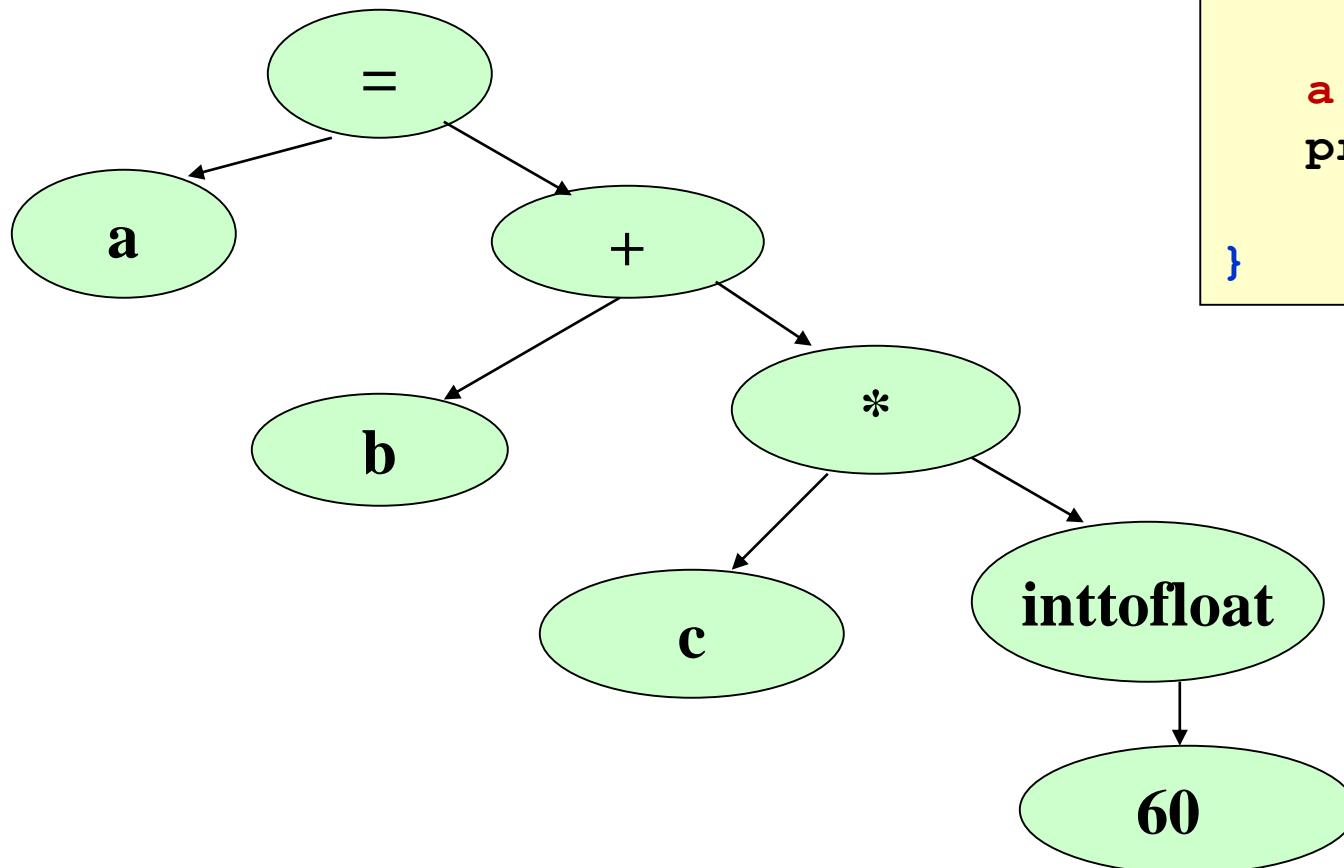


```
void cal(float b,  
        float c){  
    float a;  
  
    a = b+c*60;  
    printf("%f\n",  
           a);  
}
```



编译器流程示例——语义分析

为程序中赋值语句生成的抽象语法树

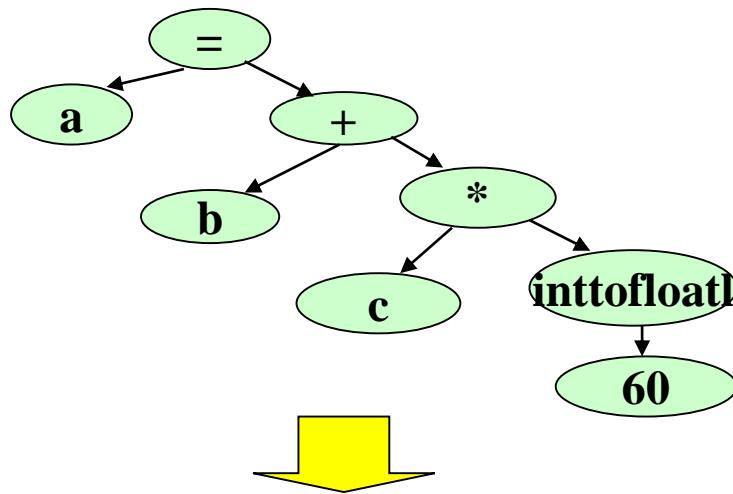


```
void cal(float b,  
        float c){  
    float a;  
  
    a = b+c*60;  
    printf("%f\n",  
          a);  
}
```



编译器流程示例——中间代码生成

为程序中赋值语句生成的中间代码表示



```
void cal(float b,  
        float c){  
    float a;  
  
    a = b+c*60;  
    printf("%f\n",  
          a);  
}
```

```
temp1 = inttofloat(60);  
temp2 = c * temp1;  
temp3 = b + temp2;  
a = temp3;
```



编译器流程示例——代码优化

中间代码

```
temp1 = inttofloat(60);  
temp2 = c * temp1;  
temp3 = b + temp2;  
a = temp3;
```

```
void cal(float b,  
         float c){  
    float a;  
  
    a = b+c*60;  
    printf("%f\n",  
           a);  
}
```



优 化



```
temp1 = c * 60.0  
a = b + temp1
```



编译器流程示例——目标代码生成

```
temp1 = c * 60.0  
a = b + temp1
```



目标代码生成



```
LDF    R2, c  
MULF   R2, R2, #60.0  
LDF    R1, b  
ADDF   R1, R1, R2  
STF    a, R1
```

```
void cal(float b,  
         float c){  
    float a;  
  
    a = b+c*60;  
    printf("%f\n",  
           a);  
}
```

符号表

名字	种类	类型	其他信息
cal	函数		
a	变量	float	
b	变量	float	
c	变量	float	
printf	函数		
...			

```
void cal(float b,  
        float c){  
    float a;  
  
    a = b+c*60;  
    printf("%f\n",  
           a);  
}
```



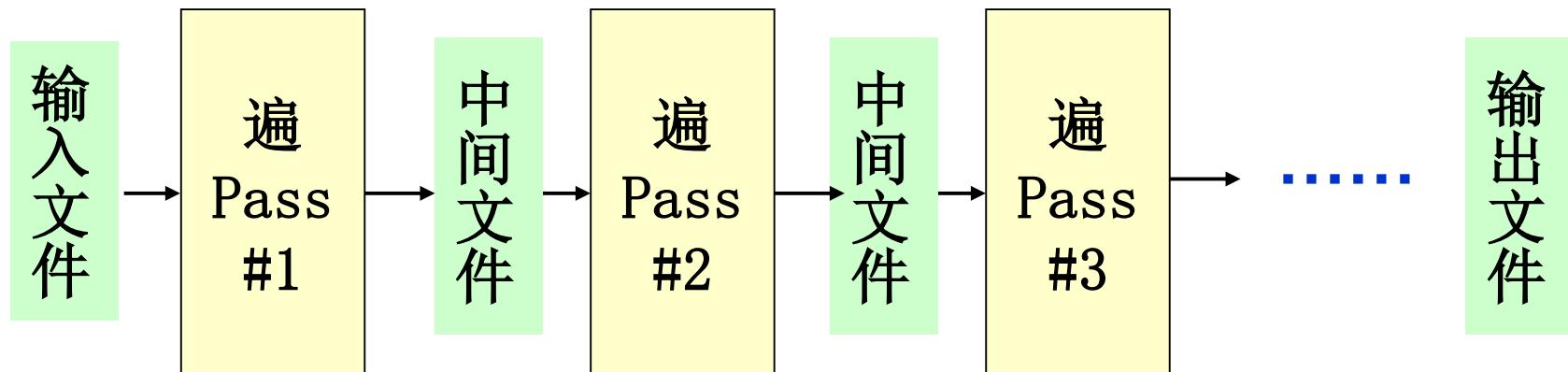
错误诊断与处理

- 编译程序在各个阶段应诊断和报告源程序中的错误
 - 词法错误
 - 语法错误
 - 语义错误
- 编译程序应报告出错位置（文件、行、列），并给出相应的纠错提示信息。
- 出错处理能力的大小是衡量编译程序质量好坏的一个重要指标。



遍 (PASS)

- 对源程序或其等价的中间形式（通常以文件形式存在）从头到尾扫视一次，完成预定的处理任务。





编译器构造工具

- 目标：辅助用户生成复杂的、容易出错的编译分析程序
- 词法分析工具
 - Lex, Flex
- 语法分析工具
 - YACC, Bison, ANTLR
- 研究/教学型编译器
 - SUIF, Soot, JikeRVM, LLVM,



主要内容

- 编译的起源
 - 程序设计语言的发展
 - 编译的基本概念
 - 编译器vs解释器
- 编译器的结构
- 编译技术的应用
- 程序设计语言基础



编译技术的应用 (1/3)

□ 高级语言的实现

- 流行编程语言的大多数演变都是朝着提高抽象级别的方向，每一轮编程语言新特征的出现都刺激编译器优化的新研究
 - 支持用户定义的**聚合数据类型**和高级控制流，如数组和记录、循环和过程调用：C、Fortran
 - **面向对象**的主要概念是数据抽象和性质继承，使得程序更加模块化并易于维护：Smalltalk、C++、C#、Java
 - **类型安全**的语言：Java没有指针，也不允许指针算术。它用无用单元收集机制来自动地释放那些不再使用的变量占据的内存
 - 面向大规模**并行编程**的语言与编译优化



编译技术的应用 (2/3)

- 针对计算机体系结构的优化
 - 计算机体体系结构的演化推动了新的编译器技术的研究
 - 并行化: 多核、众核、云计算
 - 内存分层 (memory hierarchy)
- 新型计算机体体系结构的设计
 - 计算机系统的性能不仅仅取决于它的原始速度, 还取决于编译器是否能生成充分利用其特征的代码
 - 在现代计算机体体系结构的研究中, 在处理器的设计阶段就开发编译器, 并将编译生成的代码在模拟器上运行, 以评价拟采用体系结构的特征



编译技术的应用 (3/3)

□ 提高软件开发效率的工具

源于代码优化技术的程序分析一直在改进软件开发效率

- 类型检查：

类型检查是一种捕捉程序中前后不一致的有效技术

- 边界检查

数据流分析技术可用来定位缓冲区溢出

- 内存管理

自动的内存管理删除内存泄漏等内存管理错误

□ 其他应用

- 软件测试工具，程序理解工具

- 程序变换：二进制翻译、硬件综合(hardware synthesis)

- 等等



主要内容

- 编译的起源
 - 程序设计语言的发展
 - 编译的基本概念
 - 编译器vs解释器
- 编译器的结构
- 编译技术的应用
- 程序设计语言基础



程序设计语言基础 (1/6)

□ 静态/动态 (static vs dynamic)

- **静态**: 语言策略支持编译器静态决定某个问题
- **动态**: 只允许在程序运行时刻作出决定
- 例: Java类声明中的**static**指明了变量的存放位置可静态确定。

□ 作用域 (scope)

- x的一个声明的**作用域**指程序中的一个区域，其中对x的使用都指向这个声明。
- **静态作用域**: 通过静态阅读程序即可决定作用域
- **动态作用域**



程序设计语言基础 (2/6)

□ 环境与状态

■ 环境

(environment):

是从名字到内存地址的映射

■ 状态 (state): 从内存地址到它们的值的映射

```
...
int i;                                /* 全局 i */
...
void f(...){ 
    int i;                            /* 局部 i */
    ...
    i = 3;                           /* 对局部 i 的使用 */
    ...
}
...
x = i + 1;                          /* 对全局 i 的使用 */
```

图 1-9 名字 i 的两个声明



程序设计语言基础 (3/6)

□ 静态作用域和块结构

- C语言(以及C++/Java/C#等)使用静态作用域。
 - C语言程序由顶层的变量、函数声明组成
 - 函数内部可以声明变量（局部变量/参数），这些声明的作用域在它出现的函数内
 - 一个顶层声明的作用域包括其后的所有程序。
- 作用域规则基于程序结构，声明的作用域由它在程序中的位置决定。
- 也可以通过**public**、**private**、**protected**等限定词进行明确控制



程序设计语言基础 (4/6)

```
main() {  
    int a = 1;  
    int b = 1;  
    {  
        int b = 2;  
        {  
            int a = 3;           B3  
            cout << a << b;  
        }  
        {  
            int b = 4;           B4  
            cout << a << b;  
        }  
        cout << a << b;  
    }  
    cout << a << b;  
}
```

B_1

B_2

B_3

B_4

□ 块作用域的例子

声 明	作用域
int a = 1;	$B_1 - B_3$
int b = 1;	$B_1 - B_2$
int b = 2;	$B_2 - B_4$
int a = 3;	B_3
int b = 4;	B_4



程序设计语言基础 (5/6)

□ 参数传递机制

- **传值调用 (call by value):** 对实在参数求值/拷贝，再存放到被调用过程的形参的内存位置上。
- **引用调用 (call by reference):** 实际传递的是实在参数的地址。
- **名调用 (call by name):** 早期使用，现在已经废弃。



程序设计语言基础 (6/6)

□ 别名(Alias)

- 两个指针指向同一个位置的情况
- 在函数调用时，看起来不同的形式参数实际上可能互为别名

□ 如何发现别名？

- Alias Analysis (或 Points-to Analysis)

□ 为什么区分别名很重要？



本章小结

- 编译是把一种语言编写的源程序翻译成为另外一种语言的目标程序的过程
 - 编译器与解释器的区别
- 编译器具有非常良好的系统架构，对于大规模程序开发来说具有很好的借鉴意义
 - 编译器的阶段划分
 - 前端、后端、遍
- 编译技术在计算机科学的各个领域得到了广泛的应用
- 了解程序设计语言的基础知识



Next Class

- 第3章 “词法分析”
- 自学：龙书第2章
 - “一个简单的语法制导翻译器”
 - 阅读本章内容，并对照附录代码学习