

自制 编程语言

[日] 前桥和弥 著

刘卓 徐谦 吴雅明 译



TURING

图灵程序设计丛书

自制 编程语言

[日] 前桥和弥 著

刘卓 徐谦 吴雅明 译



人民邮电出版社
北京

图书在版编目(CIP)数据

自制编程语言 / (日) 前桥和弥著; 刘卓, 徐谦,
吴雅明译. -- 北京: 人民邮电出版社, 2013.12

(图灵程序设计丛书)

ISBN 978-7-115-33320-9

I. ①自… II. ①前… ②刘… ③徐… ④吴… III.
①C语言—程序设计 IV. ①TP312

中国版本图书馆CIP数据核字(2013)第238549号

内 容 提 要

本书手把手地教读者用C语言制作两种编程语言: crowbar与Diksam。crowbar是运行分析树的无类型语言, Diksam是运行字节码的静态类型语言。这两种语言都具备四则运算、变量、条件分支、循环、函数定义、垃圾回收等功能, 最终版则可以支持面向对象、异常处理等高级机制。所有源代码都提供下载, 读者可以一边对照书中的说明一边调试源代码。这个过程对理解程序的运行机制十分有帮助。

本书适合有一定基础的程序员和编程语言爱好者阅读。

-
- ◆ 著 [日] 前桥和弥
译 刘卓 徐谦 吴雅明
责任编辑 乐馨
执行编辑 金松月
责任印制 焦志炜
- ◆ 人民邮电出版社出版发行 北京市丰台区成寿路11号
邮编 100164 电子邮件 315@ptpress.com.cn
网址 <http://www.ptpress.com.cn>
北京天宇星印刷厂印刷
- ◆ 开本: 800×1000 1/16
印张: 24.75
字数: 565千字 2013年12月第1版
印数: 1-4 000册 2013年12月北京第1次印刷
著作权合同登记号 图字: 01-2013-4381号
-

定价: 79.00元

读者服务热线: (010)51095186转604 印装质量热线: (010)81055316

反盗版热线: (010)81055315

广告经营许可证: 京崇工商广字第0021号

版权声明

PROGRAMMING GENGO WO TSUKURU by Kazuya Maebashi

Copyright © 2009 Kazuya Maebashi

All rights reserved.

Original Japanese edition published by Gijyutsu-Hyoron Co.,Ltd.,Tokyo

This Simplified Chinese language edition published by arrangement with
Gijyutsu-Hyoron Co.,Ltd.,Tokyo in care of Tuttle-Mori Agency, Inc.,Tokyo

本书中文简体字版由 Gijyutsu-Hyoron Co.,Ltd. 授权人民邮电出版社独家出版。未经出版者书面许可，不得以任何方式复制或抄袭本书内容。

版权所有，侵权必究。

（图灵公司感谢李典对本书的审读）

译者序

能翻开这本书的人，想必对编程都有着浓厚的兴趣。大部分编程爱好者都会利用业余时间写一些小程序、开源项目作为消遣，却很少有人会想要自己创造一门编程语言，这是为什么呢？

在翻译本书之前，如果别人问我要不要尝试自制编程语言，我一定会觉得他疯了。因为在潜意识里，我一直认为制作编程语言应该是 C 语言之父丹尼斯·里奇这样的业界大牛才能完成的浩大工程，作为一个普通程序员只要安于本分，用好已有的语言就已经足够了。

在翻译完本书后，我才发现自己真的是大错特错。原来创造一门编程语言，只需要一些 C 语言基础、一些正则表达式知识、加上不断思索的大脑就可以做到。如果你还觉得难以置信，那么就请看看在这本不算厚的书中，作者居然已经创造了两门编程语言，并且都具备高级编程语言的所有特性。

其实一开始的问题已经有了答案：很多看似难如登天的事情，一旦真的下决心去做，你会发现难度并没有想象中那么高，只是我们往往缺少一颗勇于挑战的心罢了。

本书记录了作者一步一步从零创造出编程语言的全过程，作者并不是什么行业精英，而是像你我一样的普通开发者。整本书中也没有用特别复杂的算法或酷炫的编程技巧，但是就凭借着一行行简单朴实的编程语句，作者最终完成了一个普通开发者看来几乎不可能完成的任务。阅读完本书后，除了自制编程语言的知识，我相信读者还能收获到一些更重要的东西。

本书原文讲到了日文编码的知识，为了更好的将内容精髓呈现给读者，我们大胆地将涉及日文编码的部分全部更改为中文编码的知识，译者刘卓还对此编写了很多原创的补充内容，力求能与原书保持同样的水平。如有错误或疏漏，还请读者随时指正。

读完全书后，你会对编程语言的原理和实现方式有一个全面深入的了解，比如你会明白为什么 Java 中 String 类型明明是对象类型却不能改变其内容，C 语言中为什么 `a++ + ++b` 这样看似合理的语句却会报错等。以前

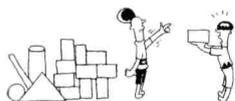
知其然而不知其所以然的问题都会得到答案，这对日后进行更高阶的开发有很大的帮助。

更重要的是，你可以获得自制编程语言的能力，从而可以去做很多以前敢想却没有能力做的事情，比如我现在就在构思能否创造一门以文言文和中国古代文化为基础的编程语言：易经八卦就是天然的二维矩阵，《九章算术》则有不少基础算法……相信读者还会有更加天才有趣的想法出现。如果能运用本书中的知识最终将其实现，那么这将对翻译工作最好的肯定。

最后，在这里代表其他二位译者一并感谢在翻译过程中给予我们帮助和支持的家人、同事，让这本书最终得以问世。

徐谦

2013年中秋



前言



这本书是为那些想独立制作一门编程语言的人而写的。

一听到这个话题，有的人会想：太疯狂了，制作编程语言肯定很有难度吧？有人会怀疑：制作编程语言能有什么用呢？其实这些都是误解。

制作编程语言在技术层面上其实并不难，只要掌握一些基础知识即可。而且，制作编程语言对于我们深入理解日常使用的 C、Java、JavaScript 等语言都有帮助。在一些应用程序的内置脚本语言中，我们也经常会因为种种限制从而萌生制作替代语言的想法。因此，自制编程语言并不是少数极客的个人癖好，它对大多数程序员都颇具实用价值。

日本关于制作编程语言的书籍已经很多了，其中一些还被选定为大学教科书。这些书中常出现有限状态机、NFA、LL(1)、LR(1)、SLA 等专业词汇，同时还大量使用 \cap 、 \in 等数学符号，对于不熟悉这部分理论知识的人（包括我自己在内）来说非常难以读懂。针对这种现状，本书会偏重实践，避免枯燥的理论。

本书将分别制作两种编程语言：crowbar 与 Diksam。crowbar 是运行分析树的无类型语言，Diksam 是运行字节码的静态类型语言。无论哪种语言，都具备四则运算、变量、条件分支、循环、函数定义、垃圾回收等功能，最终版则可以支持面向对象、异常处理等高级机制。总之，作为现代编程语言所必须具备的功能都基本覆盖了（唯一可能没实现的就是多线程了吧）。所有源代码都提供下载，读者可以一边对照书中的说明一边调试源代码，这样应该不难理解整个程序的运行机制。

当然，要一次实现如此多功能的编程语言，对于初学者而言可能有点吃力，因此本书会详细介绍 crowbar、Diksam 的制作步骤，请放心。

在制作编程语言的过程中，我体会到了一种无法用语言形容的快乐。其实无论在日本或其他地区，世界上还有很多人都在尝试自制编程语言，这正是编程语言不断增加的原因。如果以本书为契机，有朝一日你也向本已混乱的巴比伦之塔再添一门新语言的话，作为本书作者，这将是无上的光荣。

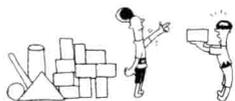
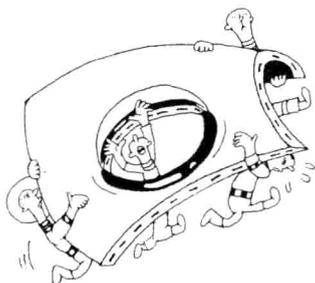


在本书的撰写过程中，得到了很多朋友的帮助与支持：

感谢百忙之中通读原稿并给出很多改进意见的吉田敦、间野健二、藤井壮一、山本将；感谢对本书原型，即网页版“自制编程语言”提出意见的朋友；感谢对博客连载“自制编程语言日记”提出意见的读者朋友，以及实际使用 crowbar 与 Diksam 并提出意见的朋友。最后还要感谢每次对我延迟交稿仍然充满耐心的技术评论社的熊谷裕美子编辑。多亏大家的鼎力支持，本书才终能完成，在此我表示深深的谢意。

2009 年 5 月 7 日 01:06 J.S.T

前桥和弥



目录CONTENTS

第1章 引子.....001

1.1 为什么要制作编程语言.....002

1.2 自制编程语言并不是很难.....003

1.3 本书的构成与面向读者.....004

1.4 用什么语言来制作.....006

1.5 要制作怎样的语言.....007

1.5.1 要设计怎样的语法.....007

1.5.2 要设计怎样的运行方式.....009

补充知识 “用户”指的是谁?.....012

补充知识 解释器并不会进行翻译.....012

1.6 环境搭建.....012

1.6.1 搭建开发环境.....012

补充知识 关于bison与flex的安装.....014

1.6.2 本书涉及的源代码以及编译器.....015



第2章 试做一个计算器.....017

2.1 yacc/lex是什么.....018

补充知识 词法分析与解析器是各自独立的.....019

2.2 试做一个计算器.....020

2.2.1 lex.....021

2.2.2 简单正则表达式讲座.....024

2.2.3 yacc.....026

2.2.4 生成执行文件.....033

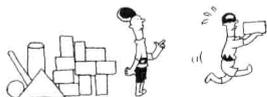
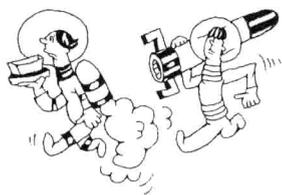
2.2.5 理解冲突所代表的含义.....034

2.2.6 错误处理.....040

2.3 不借助工具编写计算器	041
2.3.1 自制词法分析器	041
补充知识 保留字(关键字)	046
补充知识 避免重复包含	047
2.3.2 自制语法分析器	048
补充知识 预读记号的处理	053
2.4 少许理论知识——LL(1)与LALR(1)	054
补充知识 Pascal/C中的语法处理诀窍	056
2.5 习题：扩展计算器	056
2.5.1 让计算器支持括号	056
2.5.2 让计算器支持负数	058

第3章 制作无类型语言 crowbar

3.1 制作crowbar ver.0.1语言的基础部分	062
3.1.1 crowbar是什么	062
3.1.2 程序的结构	063
3.1.3 数据类型	064
3.1.4 变量	064
补充知识 初次赋值兼做变量声明的理由	066
补充说明 各种语言的全局变量处理	067
3.1.5 语句与结构控制	067
补充知识 elif、elsif、elseif的选择	068
3.1.6 语句与运算符	069
3.1.7 内置函数	069
3.1.8 让crowbar支持C语言调用	070
3.1.9 从crowbar中调用C语言(内置函数的编写)	071
3.2 预先准备	071





3.2.1	模块与命名规则	072
3.2.2	内存管理模块 MEM	073
	补充知识 valgrind	075
	补充知识 富翁式编程	075
	补充知识 符号表与扣留操作	076
3.2.3	调试模块 DBG	076
3.3	crowbar ver.0.1 的实现	077
3.3.1	crowbar 的解释器——CRB_Interpreter	077
	补充知识 不完全类型	080
3.3.2	词法分析——crowbar.l	081
	补充知识 静态变量的许可范围	084
3.3.3	分析树的构建——crowbar.y 与 create.c	085
3.3.4	常量折叠	089
3.3.5	错误信息	089
	补充知识 关于 crowbar 中使用的枚举型定义	091
3.3.6	运行——execute.c	092
3.3.7	表达式评估——eval.c	096
3.3.8	值——CRB_Value	104
3.3.9	原生指针型	105
3.3.10	变量	106
3.3.11	字符串与垃圾回收机制——string_pool.c	108
3.3.12	编译与运行	110

第 4 章 数组和 mark-sweep 垃圾回收器

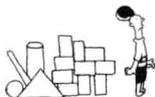
4.1	crowbar ver.0.2	114
4.1.1	crowbar 的数组	114
4.1.2	访问数组元素	115



4.1.3	数组是一种引用类型	116
	补充知识 “数组的数组”和 multidimensional arrays	116
4.1.4	为数组添加元素	118
4.1.5	增加 (模拟) 函数调用功能	118
4.1.6	其他细节	118
4.2	制作 mark-sweep GC	119
4.2.1	引用数据类型的结构	119
4.2.2	mark-sweep GC	121
	补充知识 引用和 immutable	123
4.2.3	crowbar 栈	124
4.2.4	其他根	127
4.2.5	原生函数的形式参数	128
4.3	实现 GC 本身	129
4.3.1	对象的管理方法	129
4.3.2	GC 何时启动	129
4.3.3	sweep 阶段	132
	补充知识 GC 现存的问题	133
	补充知识 Coping GC	134
4.4	其他修改	136
4.4.1	修改语法	136
4.4.2	函数的模拟	137
4.4.3	左值的处理	139
4.4.4	创建数组和原生函数的书写方法	142
4.4.5	原生指针类型的修改	144

第 5 章 中文支持和 Unicode

5.1	中文支持策略和基础知识	148
-----	-------------------	-----



5.1.1	现存问题.....	148
5.1.2	宽字符（双字节）串和多字节字符串.....	149
	补充知识 wchar_t 肯定能表示1个字符吗?	150
5.1.3	多字节字符/宽字符之间的转换函数群.....	150
5.2	Unicode.....	153
5.2.1	Unicode的历史.....	153
5.2.2	Unicode的编码方式.....	154
	补充知识 Unicode可以固定（字节）长度吗?	156
5.3	crowbar book_ver.0.3的实现.....	156
5.3.1	要实现到什么程度?	156
5.3.2	发起转换的时机.....	157
5.3.3	关于区域设置.....	158
5.3.4	解决0x5C问题.....	158
	补充知识 失败的#ifdef.....	160
5.3.5	应该是什么样子.....	160
	补充知识 还可以是别的样子——Code Set Independent.....	161



第6章

	制作静态类型的语言 Diksam.....	163
6.1	制作Diksam Ver 0.1语言的基本部分.....	164
6.1.1	Diksam的运行状态.....	164
6.1.2	什么是Diksam.....	165
6.1.3	程序结构.....	165
6.1.4	数据类型.....	166
6.1.5	变量.....	166
6.1.6	语句和流程控制.....	167
6.1.7	表达式.....	167
6.1.8	内建函数.....	168

6.1.9 其他	168
6.2 什么是静态的/执行字节码	169
6.2.1 静态类型的语言	169
6.2.2 什么是字节码	169
6.2.3 将表达式转换为字节码	170
6.2.4 将控制结构转换为字节码	173
6.2.5 函数的实现	173
6.3 Diksam ver.0.1的实现——编译篇	175
6.3.1 目录结构	175
6.3.2 编译的概要	176
6.3.3 构建分析树 (create.c)	176
6.3.4 修正分析树 (fix_tree.c)	179
6.3.5 Diksam的运行形式——DVM_Executable	185
6.3.6 常量池	186
补充知识 YARV的情况	187
6.3.7 全局变量	188
6.3.8 函数	189
6.3.9 顶层结构的字节码	189
6.3.10 行号对应表	190
6.3.11 栈的需要量	190
6.3.12 生成字节码 (generate.c)	191
6.3.13 生成实际的编码	193
6.4 Diksam虚拟机	197
6.4.1 加载/链接DVM_Executable到DVM	200
6.4.2 执行——巨大的switch case	202
6.4.3 函数调用	204



第7章 为 Diksam 引入数组.....207

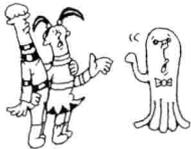
- 7.1 Diksam 中数组的设计208
 - 7.1.1 声明数组类型的变量208
 - 7.1.2 数组常量.....209
 - 补充知识** D 语言的数组210
- 7.2 修改编译器210
 - 7.2.1 数组的语法规则210
 - 7.2.2 TypeSpecifier 结构体.....212
- 7.3 修改 DVM.....213
 - 7.3.1 增加指令213
 - 补充知识** 创建 Java 的数组常量215
 - 补充知识** C 语言中数组的初始化217
 - 7.3.2 对象217
 - 补充知识** ArrayStoreException218
 - 7.3.3 增加 null219
 - 7.3.4 哎!还缺点什么?219



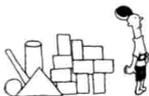
第8章 将类引入 Diksam.....221

- 8.1 分割源文件222
 - 8.1.1 包和分割源代码222
 - 补充知识** #include、文件名、行号.....225
 - 8.1.2 DVM_ExecutableList225
 - 8.1.3 ExecutableEntry226
 - 8.1.4 分开编译源代码227
 - 8.1.5 加载和再链接.....230
 - 补充知识** 动态加载时的编译器.....233
- 8.2 设计 Diksam 中的类.....233





8.2.1	超简单的面向对象入门.....	233
8.2.2	类的定义和实例创建.....	237
8.2.3	继承.....	239
8.2.4	关于接口.....	241
8.2.5	编译与接口.....	242
8.2.6	Diksam 怎么会设计成这样?.....	243
8.2.7	数组和字符串的方法.....	245
8.2.8	检查类的类型.....	246
8.2.9	向下转型.....	246
8.3	关于类的实现——继承和多态.....	247
8.3.1	字段的内存布局.....	247
8.3.2	多态——以单继承为前提.....	249
8.3.3	多继承——C++.....	250
8.3.4	Diksam 的多继承.....	252
	补充知识 无类型语言中的继承.....	254
8.3.5	重写的条件.....	254
8.4	关于类的实现.....	256
8.4.1	语法规则.....	256
8.4.2	编译时的数据结构.....	258
8.4.3	DVM_Executable 中的数据结构.....	260
8.4.4	与类有关的指令.....	262
	补充知识 方法调用、括号和方法指针.....	263
8.4.5	方法调用.....	264
8.4.6	super.....	266
8.4.7	类的链接.....	266
8.4.8	实现数组和字符串的方法.....	267
8.4.9	类型检查和向下转型.....	267
	补充知识 对象终结器 (finalizer) 和析构函数 (destructor).....	268



第9章	应用篇	271
9.1	为 crowbar 引入对象和闭包	272
9.1.1	crowbar 的对象	272
9.1.2	对象实现	273
9.1.3	闭包	274
9.1.4	方法	276
9.1.5	闭包的实现	278
9.1.6	试着跟踪程序实际执行时的轨迹	281
9.1.7	闭包的语法规则	284
9.1.8	普通函数	284
9.1.9	模拟方法 (修改版)	285
9.1.10	基于原型的面向对象	286
9.2	异常处理机制	286
9.2.1	为 crowbar 引入异常	286
9.2.2	setjmp()/longjmp()	289
	补充知识 Java 和 C# 异常处理的不同	293
9.2.3	为 Diksam 引入异常	295
	补充知识 catch 的编写方法	296
9.2.4	异常的数据结构	297
9.2.5	异常处理时生成的字节码	299
9.2.6	受查异常	301
	补充知识 受查异常的是与非	303
	补充知识 异常处理本身的是与非	304
9.3	构建脚本	305
9.3.1	基本思路	306
9.3.2	YY_INPUT	307
9.3.3	Diksam 的构建脚本	308
9.3.4	三次加载/链接	308

