

Programming Massively Parallel Processors: A Hands-on Approach, Second Edition

大规模并行处理器 编程实战 (第2版)

[美] David B. Kirk
Wen-mei W. Hwu 著
赵开勇 汪朝辉 程亦超 译

清华大学出版社



• 014006937

TP311.11

47-2

大规模并行处理器

编程实战

(第 2 版)

[美] David B. Kirk
Wen-mei W. Hwu 著
赵开勇 汪朝辉 程亦超 译

清华大学出版社



北航

C1693932

TP311.11
47-2

TEC300810

Programming Massively Parallel Processors: A Hands-on Approach, Second Edition

David B. Kirk, Wen-mei W. Hwu

EISBN: 978-0-12-415992-1

Copyright © 2013 by Elsevier Inc. All rights reserved.

Authorized Simplified Chinese translation edition published by Elsevier (Singapore) Pte Ltd Press and Tsinghua University.

Copyright © 2013 by Elsevier (Singapore) Pte Ltd and Tsinghua University Press. All rights reserved.

Published in China by Tsinghua University Press under special arrangement with Elsevier (Singapore) Pte Ltd.. This edition is authorized for sale in China only, excluding Hong Kong, Macau SAR and Taiwan. Unauthorized export of this edition is a violation of the Copyright Act. Violation of this Law is subject to Civil and Criminal Penalties.

本书简体中文版由 Elsevier (Singapore) Pte Ltd. 授予清华大学出版社在中国大陆地区(不包括香港、澳门特别行政区以及台湾地区)出版与发行。未经许可之出口，视为违反著作权法，将受法律之制裁。

北京市版权局著作权合同登记号 图字: 01-2013-0278

本书封面贴有 Elsevier 防伪标签，无标签者不得销售。

版权所有，侵权必究。侵权举报电话：010-62782989 13701121933

图书在版编目(CIP)数据

大规模并行处理器编程实战(第2版)/(美)柯克(Kirk, D.B.), (美)胡(Hwu, W.W.)著; 赵开勇, 汪朝辉, 程亦超 译. —北京: 清华大学出版社, 2013

书名原文: Programming Massively Parallel Processors: A Hands-on Approach, Second Edition

ISBN 978-7-302-34272-4

I. ①大… II. ①柯… ②胡… ③赵… ④汪… ⑤程… III. ①并行程序—程序设计 IV. ①TP311.11

中国版本图书馆 CIP 数据核字(2013)第 249417 号

责任编辑: 王军 刘伟琴

装帧设计: 牛静敏

责任校对: 曹阳

责任印制: 宋林

出版发行: 清华大学出版社

网 址: <http://www.tup.com.cn>, <http://www.wqbook.com>

地 址: 北京清华大学学研大厦 A 座 邮 编: 100084

社 总 机: 010-62770175 邮 购: 010-62786544

投稿与读者服务: 010-62776969, c-service@tup.tsinghua.edu.cn

质 量 反 馈: 010-62772015, zhiliang@tup.tsinghua.edu.cn

印 刷 者: 清华大学印刷厂

装 订 者: 三河市新茂装订有限公司

经 销: 全国新华书店

开 本: 185mm×230mm 印 张: 27 字 数: 554 千字

版 次: 2013 年 11 月第 1 版 印 次: 2013 年 11 月第 1 次印刷

印 数: 1~4000

定 价: 59.80 元

产品编号: 049320-01

译者序

本书是一本由浅入深地介绍异构高性能编程的优秀教材。

2007 年, NVIDIA 公司推出了 CUDA(Compute Unified Device Architecture, 统一计算设备架构)这一编程模型的 1.0 版本。作为译者, 我有幸于 2007 年底就开始接触 CUDA。这几年间, 异构高性能计算不断发展, 支持异构计算的架构不断更新和涌现, 越来越多的人开始关注异构高性能计算。在大数据时代, 异构高性能计算是处理大数据的一把利器。传统并行计算算法的研究在 20 世纪 60 年代和 70 年代非常活跃, 现在的很多并行算法实际上在那个时代都已出现。技术的不断发展, 芯片的高度集成, 使得以前需要用大规模计算集群完成的工作, 现在可在一个小芯片上完成, 从而使高性能并行计算得以普及。高性能计算小芯片(包括多核心 CPU、GPU、众核核心协处理器等)使万亿次计算走向了单个节点。除了高性能计算集群外, 个人高性能计算机和移动平台都会用到高性能编程技术。高性能编程技术在个人计算机和移动平台将成为未来的一个发展方向。针对移动平台的高性能计算, 将来也会是一个重要领域。且计算和传输本就是一体的, 为在降低传输带宽的同时达到高质量的传输, 越来越多的新协议需要大量的计算。单节点的高性能计算能力, 可以很好地解决计算部分的问题。

硬件的变革要求我们紧跟国际步伐, 在这几年的 TOP 500 排名中, 中国多次排名前列。但是, 对于高性能计算而言, 光有硬件架构往往是不够的。近几年硬件的变化基本是 18 个月左右革新一次。如果没有优秀软件的支持, 再好的硬件也很快就会过时, 投入的大量人力物力就会在几年之内迅速贬值。对于新的硬件架构必须开发新的软件, 全世界基本是在同一条起跑线上。我们希望更多的开发者加入到异构高性能开发中来。高性能异构计算也是解决大数据的一个强有力手段。

本书首先讲述硬件的体系结构, 以 GPU 为例讲解异构计算芯片的架构, 然后结合实际开发示例, 讲解软件系统结构, 以及异构高性能计算软件开发的思想和实践。异构高性能计算芯片的发展近几年非常活跃, 作为译者, 我们希望, 通过翻译本书, 能架起知识传播的桥梁, 将国外的最新技术传播给国内的读者。

感谢清华大学出版社对我们的信任, 把这本重要的书交给我们翻译。感谢参与校对和审稿的同仁们: 武汉某研究所的周学谦, 香港中文大学的王若薇, 武汉大学的方留杨, 香港科技大学的李佳俐, 北京师范大学的赵国兴, 中国科技大学近代力学系的冉伟, 香港浸会大学异构计算实验室的梅辛欣, 中科院超算中心的黄潘育, 华大基因的张帆。没有你们

II 大规模并行处理器编程实战(第2版)

的支持与帮助，本书将无法顺利付梓。在此，还要由衷地感谢我的导师褚晓文老师，感谢您对我的支持与鼓励。

异构高性能计算不断发展，翻译本书旨在为传播知识做一点绵薄之力。书中所涉及的领域比较广泛，翻译中如有不当之处，敬请广大读者批评指正。

再次向参与、支持和关切本书出版的同仁表示诚挚的谢意！

赵开勇

2013年7月4日

于香港浸会大学计算机系异构计算实验室

致 谢

有很多人为我们出版第 2 版提供了帮助。首先需要感谢本书新章节的作者。Yuan Lin 和 Vinod Grover 撰写了 OpenACC 这一章的初稿。Nathan Bell 和 Jared Hoberock 在 Chris Rodrigues 的基础上撰写了 Thrust 这章的初稿。Greg Ruetsch 和 Massimiliano Fatica 撰写了 CUDA FORTRAN 这章的初稿。David Callahan 撰写了 C++ AMP 这一章。Isaac Gelado 撰写了 MPI-CUDA 这章的初稿。Brent Oster 为 Kepler 章节提供了原始资料和代码实例。如果没有这些人的帮助，我们就不能深入地将这些内容提供给读者。

非常感谢 Izzat El Hajj，他不厌其烦地帮我们验证代码示例并改善插图和练习题的质量。

也要特别感谢 CUDA 之父 Ian Buck 和 Tesla GPU 计算架构的首席架构师 John Nickolls。他们的团队为这个课程提供了优秀的基础设施。在我们撰写第 2 版本的时候，John 与世长辞。我们非常怀念他。Nadeem Mohammad 组织 NVIDIA 修订和提供了附录 B。Bill Bean、Simon Green、Mark Harris、Nadeem Mohammad、Brent Oster、Peter Shirley、Eric Young 和 Cyril Zeller 对手稿进行了修订和校对。Calisa Cole 帮助设计了封面。Nadeem 的无私分享为这本书的完成起到了关键性的作用。

尤其要感谢 Jensen Huang 提供大量人力物力开设这门课程，为本书的出版奠定了基础。Tony Tamasi 的团队为本书提供了大量的评论和修订建议。Jensen 还花费了大量时间阅读早期的草稿，为我们提供了很有价值的反馈。David Luebke 为这个课程提供了 GPU 计算资源。Jonah Alben 提供了颇有价值的见解。Michael Shebanow 和 Michael Garland 给予了友情讲座和教学材料。

伊利诺伊的 John Stone 和 Sam Stone 为案例研究和 OpenCL 章节提供了大量的基础材料。John Stratton 和 Chris Rodrigues 为编程思维这一章提供了素材。I-Jui “Ray” Sung、John Stratton、Xiao-Long Wu 和 Nady Obeid 在他们的研究之外自愿担当助教，提供了实验素材并帮助改进课程材料。Jeremy Enos 孜孜不倦地工作，使学生有一个稳定的、用户友好的 GPU 计算集群用于实验和项目研究。

感谢 Dick Blahut 让我们在伊利诺伊大学有挑战性地开展了这个课程。他不断提醒我们需要写这一本书来帮助我们不断前行。Beth Katsinas 安排了一次与 Dick Blahut 和 NVIDIA 的副总 Dan Vivoli 的会议。通过这次会议，向 David 介绍了 Blahut，也促使 David 和 Wen-mei 在伊利诺伊大学开设这门课程。

IV 大规模并行处理器编程实战(第2版)

也要感谢伊利诺伊大学的 Thom Dunning 和密歇根大学的 Sharon Glotzer，作为计算机科学和工程虚拟大学的副执行长，他们盛情举办了暑期课程。Trish Barker、Scott Lathrop、Umesh Thakkar、Tom Scavo、Andrew Schuh 和 Beth McKown 都参与组织了暑期课程。Robert Brunner、Klaus Schulten、Pratap Vanka、Brad Sutton、John Stone、Keith Thulborn、Michael Garland、Vlad Kindratenko、Naga Govindaraju、Yan Xu、Arron Shinn 和 Justin Haldar 为暑期课程提供了教材和课程计划。

Nicolas Pinto 在 MIT 的课程中测试了早期版本的前几章，然后提供了非常有价值的反馈和修改建议。Steve Lumetta 和 Sanjay Patel 讲授了不同版本的课程，为我们提供了有价值的建议。John Owens 非常友善地为我们提供了他的幻灯片。Tor Aamodt、Dan Connors、Tom Conte、Michael Giles、Nacho Navarro 和全世界许许多多的教员以及他们的学生为我们提供了非常有价值的反馈信息。

尤其要感谢我们的同事 Kurt Akeley、Al Aho、Arvind、Dick Blahut、Randy Bryant、Bob Colwell、Ed Davidson、Mike Flynn、John Hennessy、Pat Hanrahan、Nick Holonyak、Dick Karp、Kurt Keutzer、Dave Liu、Dave Kuck、Yale Patt、David Patterson、Bob Rao、Burton Smith、Jim Smith 和 Mateo Valero。他们抽出时间来与我们一起分享积累多年的知识。

感谢所有为本书和相关课程做出贡献的人们，他们的慷慨和热情令我们感动。

David B. Kirk 和 Wen-mei W.Hwu

前　　言

我们很荣幸向大家介绍这本《大规模并行处理器编程实战(第2版)》。大众市场的计算系统，结合多核计算机处理单元(CPU)和多线程GPU计算，带来了万亿次计算能力的笔记本电脑和千万亿次的计算机集群。如此强大的计算能力，为我们在科研、工程、卫生和商业等领域使用计算机带来新的曙光。许多领域通过计算机系统将实验的规模、准确性、可控性和客观性等方面提升到前所未有的水平，将会在各自的领域实现突破。本书为实现这个愿景提供了一个关键因素：教授数以百万计的研究生和本科生学习并行编程，使他们拥有并行计算思维和并行编程技巧，从而应对无处不在的计算需求。

自从2010年第1版出版以来，我们收到了读者和老师的很多建议。许多人提出他们重视的一些现有功能。另一些人给我们提出了如何丰富内容的建议，从而使本书更有价值。此外，异构并行计算的硬件和软件已经越来越先进。在硬件领域，自从此书第1版出版以后，至少两款通用图像显示卡(GPU)计算架构——Fermi和Kepler已发布。在软件领域，CUDA 4.0和CUDA 5.0也可以让程序员使用Fermi和Kepler的新特性。因此，我们增加了8章新的内容并几乎重写了5章已有的内容。

总的来说，我们在保留第1版很多有用内容的基础上，增加了3部分主要的内容。第1个增加的部分是更系统地讲解了并行编程。具体内容包括：(1)新增了第8、9和10章来介绍常用的基本并行算法模式；(2)在第3、4、5和6章中增加了更多的背景材料；(3)在第7章中增加了数值稳定性的处理。这些增加的内容旨在消除“学生已经熟悉基本并行编程概念”的假设。我们还通过更多的例子来帮助读者理解。

第2个增加的部分是在异构并行计算系统中采用MPI-CUDA模型的并行编程实践。这一部分内容是大量读者希望增加的部分。鉴于GPU的成本效益明显以及单位瓦特的高吞吐量，现在许多高性能计算系统在每个节点中都使用了GPU。新增加的第19章在概念上解释了这些系统编程接口背后的框架。

第3个增加的部分是介绍了新的并行编程接口和工具，这些都可以用来加快数据并行编程。新增加的第15、16、17和18章介绍了OpenACC、Thrust、CUDA FORTRAN和C++AMP。与介绍这些工具的用户使用手册不同，我们更注重于从概念上理解使用上面这些工具解决的编程问题。

虽然我们做了这些改进，但也保留了第1版的特性，这似乎有助于它的受欢迎程度。

首先，我们尽可能保证本书简洁。虽然我们非常想增加更多有用的材料，但还是希望尽量减少页数。第二，我们尽量直观地解释问题。虽然一些概念是非常形式化的，特别是介绍基本并行算法时，但我们尽量使讲解直观而实用。

读者对象

本书的读者对象是研究生和本科生，他们来自各种科学学科和工程学科，他们的学科中需要用到计算思维能力和并行编程技巧，通过用已经普及的万亿次运算硬件能使他们在自己的领域中有所突破。我们假设读者至少已经具备了基本的 C 语言编程经验，因此不管是在计算机科学领域内还是领域外，他们都是比较高级的程序员。我们专门面向计算领域内的科学家，如机械工程、土木工程、电气工程、生物工程、物理、化学、天文学和地理学，这些科学家将通过计算在他们各自的领域中做进一步的研究。正因为这样，这些科学家既是他们领域内的专家，也是高级程序员。本书采用基本的 C 语言编程技巧，教大家用 C 语言进行并行编程。我们使用 CUDA C 并行编程环境(NVIDIA 公司开发的 GPU 上支持它，并可以在 CPU 上进行模拟)。大约有 3.75 亿左右的处理器被用户和专业人士在使用，其中超过 120 000 的程序员在积极地使用 CUDA。读者作为学习经验的一部分而开发的应用程序将会被一个庞大的用户社区使用。

如何使用本书

我们愿意提供在本书教学过程中的一些经验。从 2006 年开始，我们讲了多个类型的课程：采用为期一学期的模式和为期一周的密集模式。原来的 ECE498AL 课程在伊利诺伊大学香槟分校已成为一个永久的课程，称为 ECE408 或 CS483。当我们开设 ECE498AL 这门课的时候，开始了一些早期章节的撰写。前面 4 章在 2009 年的春天由 Nicolas Pinto 在 MIT 课程中进行了测试。从那时起，我们已经在众多的学校(包括 VSCSE 和 PUMPS 暑期学校)使用这本书。

分 3 个阶段的学习方法

为了能同时兼顾 ECE498AL 的课程和编程作业，我们把学习过程分成 3 个阶段：

第 1 个阶段——基于第 3 章的一次课致力于教会学生 CUDA 基本的存储器和线程模型、CUDA 对 C 语言的扩展和基本的编程与调试工具。在这次课后，学生要能在几个小时

内编出简单矢量相加的代码。随后的 4 个课时，让学生从概念上理解 CUDA 的内存模型、CUDA 线程执行模型、GPU 硬件特性和现代计算机体系架构。这些课时基于第 4、5 和 6 章的内容。

第 2 个阶段——一系列的课时涵盖并行计算中的浮点数问题和在高性能并行应用中需要的数据并行编程模式。这些课时基于第 7~第 10 章。经过这个阶段的学习，矩阵乘法的代码性能提高了近 10 倍。在这一阶段，学生还应该完成关于卷积、向量归约和前缀扫描的学习。

第 3 个阶段——一旦学生掌握了 CUDA 基本的编程技巧，剩下的课程内容涵盖计算思想、各种并行执行模型和并行编程规则。这些课的内容对应于第 11~第 20 章(这些课程的录音和视频可以从网上下载(<http://courses.engr.illinois.edu/ece408/>)。

把所有内容连贯起来：最终项目

当读者通过这门课、实验和本书的章节掌握了基础知识后，最终项目把所学到的经验联系起来。最终项目对这门课十分重要，它在这门课中的地位很显眼，大概需要花费两个月的时间。它采用了 5 个创新点：指导、专题讨论、检查、最终的报告和专题讨论会(在 ECE408 的网站上可以下载最终项目的更多信息，我们想提供对这些设计因素的深入思考)。

我们鼓励学生选取目前研究领域中一些具有挑战性的问题作为他们的最终项目。在这个过程中，指导老师会成立计算科学研究小组来提出问题，这些研究小组担当导师的角色。我们要求导师提供一两页的项目规格说明表，简单地描述应用程序的重要性，导师想要学生组在应用程序中完成的工作，了解与开发应用程序所需要具有的技术和技能(特定类型的数学、物理或化学课程)，以及一个网站和传统资源列表。学生可以利用这些表中对应的内容了解技术背景、总体信息构建模块，以及特定的实现方式和代码示例对应的具体的 URL 或 FTP 路径。这些项目说明表也应该为学生提供学习经验，在以后的工作当中可以定义他们自己的研究项目(在 ECE408 课程的网站上可以下载几个这种示例)。

我们也鼓励读者在项目的选择过程中与导师不断交流。只要学生和导师都同意了某个项目，他们就会密切地联系、频繁地切磋并不断地报告项目的进展情况。我们试图为学生和导师之间的协作关系提供便利，不仅是为学生也是为导师增长宝贵的经验。

项目专题讨论

整个班级对彼此的最终项目的意见主要通过专题讨论进行交流。我们通常会安排 6 个报告厅来进行专题讨论。专题讨论是为学生设计的。例如，如果一个学生已经确定了项目，

就可以把专题讨论作为一个提出初步想法、获取反馈意见和招募队友的场所。如果一个学生还没有确定自己的项目，他可以仅参加演讲和讨论并加入其中的一个项目组。在项目专题讨论期间，学生没有等级之分，这么做是为了营造一种和谐融洽的气氛，所以学生可以集中精力与导师、教学助理以及班上其他同学进行有意义的谈话。

有了专题讨论日程表后，导师和教学助理就可以花些时间给项目组提供反馈意见，以便学生可以提出问题。演讲的时间限制在10分钟内，因此在上课期间还有反馈和提问的时间。假设每次课90分钟，班级的大小最好限制在36人左右。为了严格控制时间并使反馈时间最大化，所有演示文稿预先加载到PC中。由于在专题讨论中并不是每个学生都会出席，因此每个班大概能容纳50个学生，可以根据需要提供额外的专题讨论时间。

导师和教学助理必须参加所有演讲，并给出有用的反馈意见。学生经常在下列问题上寻求帮助。第一，对于给定的时间，项目的规模是太大还是太小？第二，在这个领域中是否存在有利于这个项目的现有工作？第三，为并行执行而设立的目标计算是否适合CUDA编程模型？

设计文档

一旦学生决定选择某个项目并形成了团队后，我们就要求学生提交一个关于项目的设计文档。这有助于他们在正式进入项目前，思考整个项目的步骤。制定这种计划的能力对于他们将来事业的成功很重要。设计文档应该讨论项目的背景与目的，应用程序级目标和可能带来的影响，应用程序的主要特点，设计的概况，实现计划，它们的性能目标，经过论证的计划和可接受的测试，以及项目进度表。

在班级专题讨论会的前一周，教学助理对最终项目组进行一次项目检查。这个检查过程有助于确保学生没有偏离轨道，并且他们已经发现了设计过程早期存在的障碍。要求学生组带上他们最初的草案来检查，草案包括其应用程序的如下2个版本：(1)性能最优的CPU串行代码，采用SSE2和其他优化方案构建一系列健壮的基本代码，便于比较它们的速度；(2)性能最优的CUDA并行代码——这个版本是项目的主要成果。学生可以用这个版本认识到并行算法涉及额外的计算量时带来的系统开销。

学生组要准备好讨论每个版本代码的主要思想，任何浮点精度问题，任何针对应用程序之前结果的比较，以及显著地提高运行速度会对这个领域产生的影响。依我们的经验来看，在班级专题研讨会前一周进行检查，这是最佳的时间安排。如果时间太早，可能会把时间浪费在对不太成熟的项目和没有太大意义的会议上。如果时间太晚，学生就没有足够的时间根据反馈意见修订他们自己的项目。

项目报告

学生应该提交一份项目报告，这个报告讨论的是他们所在团队的主要发现。把 6 次课的时间放在一天，举行班级专题讨论会。在专题讨论会期间，学生所用的演讲时间正比于团队的规模。在演讲期间，学生要强调项目报告中做得最好的部分，这样也有益于整个班集体。演讲过程在学生的成绩中占的比例很高。每个学生必须单独回答一个直接向他提出的问题，因此同一个团队中每个学生的成绩也不相同。我们记录了之前的一些参与学习的情况，在 ECE408 网站上可以查看。专题研讨会为学生提供了制作简洁的演示文稿的大好机会，简洁的演示文稿能激发其他人阅读整篇文档。演讲结束后，学生要提交一份最终项目的完整报告。

在线补充

如果教师用本书作为课程的教材，就可以使用我们的实验安排、最终项目指南和示例项目说明。虽然本书提供了这些课程的知识性内容，但要想实现全部教学目标还应该查阅其他资料。我们也想邀请大家充分利用与本书相关的在线资料，这些资料可以在 ELSEVIER 出版社的网站 www.elsevierdirect.com 上下载。

最后，我们鼓励读者提供反馈意见。如果读者有关于改进本书的任何想法或者对在线资料的补充，希望能反馈给我们。当然，我们也想知道读者喜欢本书的哪一方面。

目 录

第 1 章	引言	1
1.1	异构并行计算	2
1.2	现代 GPU 的体系结构	6
1.3	为什么需要更高的速度和 并行化	8
1.4	应用程序的加速	9
1.5	并行编程语言和模型	11
1.6	本书的总体目标	12
1.7	本书的组织结构	13
	参考文献	16
第 2 章	GPU 计算的发展历程	19
2.1	图形流水线的发展	19
2.1.1	固定功能的图形流水线 时代	20
2.1.2	可编程实时图形流水线的 发展	23
2.1.3	图形与计算结合的处理器	25
2.2	GPGPU：一个中间步骤	27
2.3	GPU 计算	28
2.3.1	可扩展的 GPU	29
2.3.2	发展近况	29
2.3.3	未来发展趋势	30
	参考文献与课外阅读	30
第 3 章	CUDA 简介	35
3.1	数据并行性	36
3.2	CUDA 的程序结构	37
3.3	向量加法 kernel 函数	39
3.4	设备全局存储器与数据传输	41
3.5	kernel 函数与线程	46
3.6	小结	50
3.6.1	函数声明	50
3.6.2	启动 kernel 函数	50
3.6.3	预定义变量	51
3.6.4	运行时 API	51
3.7	习题	51
	参考文献	53
第 4 章	数据并行执行模型	55
4.1	CUDA 的线程组织	56
4.2	线程与多维数据的映射	59
4.3	矩阵乘法——一个更加复杂 的 kernel 函数	65
4.4	线程同步和透明的可扩展性	70
4.5	线程块的资源分配	73
4.6	查询设备属性	74
4.7	线程调度和容许时延	75
4.8	小结	78
4.9	习题	79
第 5 章	CUDA 存储器	81
5.1	存储器访问效率的重要性	82
5.2	CUDA 设备存储器的类型	83
5.3	减少全局存储器流量的 一种策略	89
5.4	分块矩阵乘法的 kernel 函数	93
5.5	存储器——限制并行性的 一个因素	98

5.6 小结	100
5.7 习题	101
第6章 性能优化	103
6.1 WARP 和线程执行	104
6.2 全局存储器的带宽	111
6.3 执行资源的动态划分	118
6.4 指令混合和线程粒度	120
6.5 小结	121
6.6 习题	121
参考文献	124
第7章 浮点运算	127
7.1 浮点格式	128
7.1.1 M 的规范化表示	128
7.1.2 E 的余码表示	129
7.2 能表示的数	130
7.3 特殊的位模式与 IEEE 格式 中的精度	134
7.4 算术运算的准确度和舍入	135
7.5 算法的优化	136
7.6 数值稳定性	137
7.7 小结	141
7.8 习题	141
参考文献	142
第8章 并行模式：卷积	143
8.1 背景	144
8.2 一个基本算法：一维 并行卷积	148
8.3 常数存储器和高速缓存	149
8.4 使用光环元素的分块 一维卷积	153
8.5 一个更简单的分块一维 卷积——通用高速缓存	158
8.6 小结	160
8.7 习题	161
第9章 并行模式：前缀和	163
9.1 背景	164
9.2 简单并行扫描	165
9.3 考虑工作效率	169
9.4 工作高效的并行扫描	170
9.5 任意输入长度的并行扫描	174
9.6 小结	177
9.7 习题	177
参考文献	178
第10章 并行模式：稀疏矩阵- 向量乘法	179
10.1 背景	180
10.2 使用 CSR 格式的并行 SpMV	183
10.3 填充与转置	184
10.4 用混合方法来控制填充	186
10.5 通过排序和划分来规则化	189
10.6 小结	191
10.7 习题	191
参考文献	192
第11章 应用案例研究： 高级 MRI 重构	193
11.1 应用背景	194
11.2 迭代重构	197
11.3 计算 $F^H D$	198
11.4 最终评估	214
11.5 习题	217
参考文献	218

第 12 章 应用案例研究：分子可视化和分析	219	15.4 基本的 OpenACC 程序	266
12.1 应用背景	220	15.4.1 并行构造	266
12.2 kernel 函数简单的实现方案	221	15.4.2 循环构造	267
12.3 线程粒度调节	225	15.4.3 kernels 构造	272
12.4 存储器合并	227	15.4.4 数据管理	275
12.5 小结	230	15.4.5 数据构造	276
12.6 习题	231	15.4.6 异步计算和数据传输	278
参考文献	232	15.5 OpenACC 的发展方向	279
第 13 章 并行编程和计算思想	233	15.6 习题	280
13.1 并行计算的目标	234		
13.2 问题分解	235		
13.3 算法选择	238		
13.4 计算思想	243		
13.5 小结	244		
13.6 习题	244		
参考文献	244		
第 14 章 OpenCL 简介	245		
14.1 背景	246		
14.2 数据并行性模型	247		
14.3 设备的体系结构	249		
14.4 kernel 函数	250		
14.5 设备管理和启动 kernel	251		
14.6 OpenCL 中的静电势图谱	254		
14.7 小结	258		
14.8 习题	258		
参考文献	259		
第 15 章 OpenACC 并行编程	261		
15.1 OpenACC 与 CUDA C 的比较	261		
15.2 执行模型	263		
15.3 存储器模型	265		
第 16 章 Thrust：一个面向效率的 CUDA 编程库	281		
16.1 背景简介	282		
16.2 动机	284		
16.3 Thrust 的基本特性	284		
16.3.1 迭代器和内存空间	286		
16.3.2 互操作性	286		
16.4 泛型编程	288		
16.5 抽象的益处	290		
16.5.1 编程效率	290		
16.5.2 鲁棒性	291		
16.5.3 真实性能	291		
16.6 最佳范例	293		
16.6.1 融合	293		
16.6.2 数组结构体	294		
16.6.3 隐式范围	296		
16.7 习题	297		
参考文献	298		
第 17 章 CUDA FORTRAN	299		
17.1 CUDA FORTRAN 和 CUDA C 的区别	300		
17.2 第一个 CUDA FORTRAN 程序	301		

17.3	CUDA FORTRAN 中的 多维数组	303	19.8	习题	363
17.4	用通用接口重载主机/ 设备端例程	304		参考文献	364
17.5	通过 iso_c_binding 调用 CUDA C	307	第 20 章	CUDA 动态并行	365
17.6	kernel 循环指令和归约操作 ..	309	20.1	背景	366
17.7	动态共享存储器	310	20.2	动态并行简介	367
17.8	异步数据传输	311	20.3	重要细节	368
17.9	编译和性能剖析	316	20.3.1	启动环境变量设置	369
17.10	在 CUDA FORTRAN 中 调用 Thrust	317	20.3.2	API 错误和启动失败	369
17.11	习题	321	20.3.3	事件	369
第 18 章	C++ AMP 简介	323	20.3.4	流	369
18.1	C++ AMP 核心特性	324	20.3.5	同步范围	370
18.2	C++ AMP 执行模式详解	329	20.4	内存可见性	371
	18.2.1 显式和隐式的数据复制 ..	330	20.4.1	全局内存	371
	18.2.2 异步操作	331	20.4.2	零拷贝内存	371
	18.2.3 本节小结	333	20.4.3	常量内存	371
18.3	加速器管理	333	20.4.4	局部内存	371
18.4	分块执行	335	20.4.5	共享内存	372
18.5	C++ AMP 图形特性	338	20.4.6	纹理内存	372
18.6	小结	340	20.5	一个简单示例	373
18.7	习题	341	20.6	运行时限制	376
第 19 章	异构集群编程	343	20.6.1	内存占用	376
19.1	背景简介	344	20.6.2	嵌套深度	376
19.2	运行示例	344	20.6.3	内存分配和生存周期	376
19.3	MPI 基础	346	20.6.4	ECC 错误	377
19.4	MPI 点对点通信模型	348	20.6.5	流	377
19.5	重叠计算和通信	355	20.6.6	事件	377
19.6	MPI 集合通信模型	362	20.6.7	启动池	377
19.7	小结	363	20.7	一个更复杂的示例	378
			20.7.1	线性贝塞尔曲线	378
			20.7.2	二次贝塞尔曲线	378
			20.7.3	贝塞尔曲线计算 (非动态并行版本)	378

20.7.4 贝塞尔曲线计算 (使用动态并行)	381
20.8 小结	384
参考文献	384
第 21 章 结论与展望	385
21.1 重点回顾	385
21.2 存储器模型的演变	386
21.2.1 大型虚拟和 物理地址空间	386
21.2.2 统一的设备存储空间	388
21.2.3 可配置的缓存和 暂时存储器	388
21.2.4 提高原子操作的速度	389
21.2.5 提高全局内存的 访问速度	389
21.3 kernel 函数执行控制过程 的演变	389
21.3.1 kernel 函数内部的 函数调用	389
21.3.2 kernel 函数中的 异常处理	390
21.3.3 多个 kernel 函数的 同步执行	390
21.3.4 可中断的 kernel 函数	391
21.4 内核的性能	391
21.4.1 双精度的速度	391
21.4.2 更好的控制流效率	392
21.5 编程环境	392
21.6 美好前景	392
参考文献	393
附录 A 矩阵乘法主机版的源代码	395
附录 B GPU 的计算能力	407