



HZ BOOKS  
华章科技



ELSEVIER

CUDA Application Design and Development

# 高性能CUDA 应用设计与开发

方法与最佳实践

(美) Rob Farber 著  
于玉龙 唐堃 译 郭禾 王宇新 审校



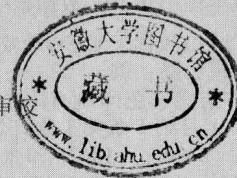
机械工业出版社  
China Machine Press

CUDA Application Design and Development

# 高性能CUDA 应用设计与开发

## 方法与最佳实践

于 2011 年 1 月第 1 版 2011 年 1 月第 1 次印刷 审核



机械工业出版社  
China Machine Press

本书是广受推崇的系统学习高性能 CUDA 应用开发与设计的经典著作，是美国国家安全实验室资深高性能编程专家多年工作经验结晶，橡树岭国家实验室资深专家鼎力推荐！本书不仅从硬件角度深入解读了 CUDA 的设计理念和 GPGPU 硬件的体系结构，而且从软件角度系统讲解了 CUDA 应用设计与开发的思想、方法、技巧、准则、注意事项和最佳实践。

第 1 章首先介绍了 CUDA 的核心概念和编程思想，以及构建与调试 CUDA 应用所需的工具和方法，然后讲解了有效提高程序性能的 CPU 编程准则；第 2 章讲解了 CUDA 在机器学习与优化中的核心概念与应用，并给出了完整的通用框架；第 3 章介绍了 CUDA 的性能分析工具套件以及性能分析的方法，同时讨论了 PCA 和 NLPCA 两种数据挖掘方法；第 4 章讲解了 CUDA 的执行模型，深刻揭示了 GPU 的工作方式和原理；第 5 章介绍了 CUDA 提供的多种 GPU 内存，以及各种内存的优缺点；第 6 章讲解了高效利用内存的技术；第 7 章介绍了 GPU 提供的多种并行方式及其应用；第 8 章首先讨论了多种 CUDA 后端设备，以及 CUDA 如何与 Python、Java、R 等高级语言交互；第 9 章讲解了 CUDA 与图形渲染混合编程；第 10 章讲解了在云计算和集群环境中使用 CUDA 的方法和技术细节；第 11 章介绍了 CUDA 在高维数据处理、力导向图、交互式工作流、量子化学等现实问题中的应用；第 12 章为学习 CUDA 设计了一个综合性的针对实时视频流的应用案例。

#### CUDA APPLICATION DESIGN AND DEVELOPMENT

Rob Farber

ISBN：978-0-12-388426-8

Copyright © 2011 by NVIDIA Corporation and Rob Farber. All rights reserved.

Authorized Simplified Chinese translation edition published by the Proprietor.

Copyright © 2013 by Elsevier (Singapore) Pte Ltd. All rights reserved.

Printed in China by China Machine Press under special arrangement with Elsevier (Singapore) Pte Ltd. This edition is authorized for sale in China only, excluding Hong Kong 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.

本书简体中文版由机械工业出版社与 Elsevier (Singapore) Pte Ltd. 在中国大陆境内合作出版。本版仅限在中国境内（不包括中国香港特别行政区及中国台湾地区）出版及标价销售。未经许可之出口，视为违反著作权法，将受法律之制裁。

封底无防伪标均为盗版

版权所有，侵权必究

本书法律顾问 北京市展达律师事务所

本书版权登记号：图字：01-2012-2655

#### 图书在版编目 (CIP) 数据

高性能 CUDA 应用设计与开发：方法与最佳实践 / (美) 罗布著；于玉龙，唐堃译. —北京：机械工业出版社，2013.1

(高性能计算系列丛书)

书名原文：CUDA Application Design and Development

ISBN 978-7-111-40446-0

I. 高… II. ①罗… ②于… ③唐… III. 计算机图形学 IV. TP391.41

中国版本图书馆 CIP 数据核字 (2012) 第 275905 号

机械工业出版社（北京市西城区百万庄大街 22 号 邮政编码 100037）

责任编辑：秦 健

北京市荣盛彩色印刷有限公司印刷

2013 年 1 月第 1 版第 1 次印刷

186mm × 240mm · 18.25 印张

标准书号：ISBN 978-7-111-40446-0

定价：59.00 元

凡购本书，如有缺页、倒页、脱页，由本社发行部调换

客服热线：(010) 88378991；88361066

购书热线：(010) 68326294；88379649；68995259

投稿热线：(010) 88379604

读者信箱：hzjsj@hzbook.com

# 译者序

目前，并行计算是提高算法性能的主要方式，无论是大规模集群系统还是个人电脑，甚至是智能手机，都将性能的提升寄托于并行计算。GPU 芯片曾主要用于计算机可视化和图形图像处理，其内部具有大量的并行计算单元，有着天然的并行优势。近年来，NVIDIA 和 AMD 等厂商相继推出以 GPU 为协处理器的通用并行计算解决方案，CUDA 是其中的典型代表。

CUDA 是 NVIDIA 推出的一种利用 GPU 进行通用并行计算的整套解决方案，包括硬件支持、程序语言扩展、编译器、调试器等整套开发工具链，它对矩阵数值计算、图像视频处理、机器学习、计算机可视化等问题都有良好的加速效果。CUDA 提供的 CUDA C 编程语言扩展，利用 Kernel 函数等机制，将大规模并行计算逻辑通过一种简单、有效且合理的方式进行表达，成为充分利用并行硬件实现程序性能加速的关键。同时，NVIDIA 还提供了充分的配套开发工具链，如 nvcc 编译器、cuda-gdb 调试器、gprof 和 computeprof 性能分析工具、Parallel Nsight 性能跟踪工具等。此外，CUDA 还包含了一些常用的支持库，如用于基础开发的 Driver、Runtime、Thrust API；用于线性代数数值计算的 CUBLAS；用于图像处理的 NPP 等。CUDA 甚至可以用于开发不使用 GPU 加速的应用程序。

本书是一本介绍 CUDA 程序设计的优秀图书。对于那些具有一定程序开发实践基础，并渴望学习并行计算、GPU 程序设计的读者来说，本书是快速入门 CUDA 并学会编写高效 GPU 并行计算程序的理想选择。本书注重理论与实践的结合，与其他介绍 CUDA 的图书相比，具有如下特点：

首先，对 CUDA 的介绍深入浅出、全面细致。对于基础知识的介绍可以使读者快速了解 CUDA，迅速上手编程；深入的分析又能够帮助读者了解不同的程序设计方式带来性能差异的原因，从而设计出更高效的程序，对于实践这点十分重要。

其次，在介绍 CUDA 的编程语言语法、各种资源以及工具的使用方法的基础上，立足于实践，选择了如机器学习、可视化系统等诸多有代表性的示例，将书中介绍的各种技术学以致用，不仅可以帮助读者更深入地理解 CUDA 编程思想，也有利于将 CUDA 技术更好地应用于实际工作中。

此外，本书内容广泛，涵盖了 CUDA 技术的各个方面，其中不仅介绍了基本的单

GPU 系统中 CUDA 程序的设计，还涉及了单机多 GPU 系统、多机 GPU 集群系统，甚至云计算环境中 CUDA 技术的应用等。本书还介绍了 CUDA 第三方支持产品和 CUDA 的衍生产品，例如可将 CUDA 程序编译为运行在常规 CPU 系统上的 PGI 等。

通过翻译本书，我们深深体会到作者对于 CUDA 和并行计算相关领域知识的深刻理解以及为编写本书所付出的努力，从中收获颇丰。我们相信读者也能够享受本书所带来的愉悦。本书的翻译力求在保证准确和接近原文风格的前提下，专业技术词汇采用国内专业通用术语，尽可能做到易于理解。

本书的第 1、3、5、6、9、10 章由于玉龙翻译，第 2、4、7、8、11、12 章由唐堃翻译，并由郭禾、王宇新和于玉龙共同完成全文终稿定稿。

感谢机械工业出版社的各位编辑在本书翻译过程中提供的帮助；感谢沈阳药科大学王薇同学在图书译稿的校对工作中提供的大量帮助。还要感谢我们的科研团队——大连理工大学软件学院高性能计算与图像处理实验室、大连理工大学计算机科学与技术学院系统结构实验室的各位同仁为本书翻译工作提供的各种支持。另外，还要特别感谢已经毕业的程童同学，为我们实验室开展 GPU 并行计算所做出的开创性工作，以及对本书翻译工作的支持。最后，感谢家人在图书翻译工作中给予我们的支持和鼓励。

由于时间仓促，翻译过程中的瑕疵和错误在所难免，欢迎各位读者批评指正。您的批评与建议会成为我们进步的基石。

# 序 言

作为一种革新的技术，GPU 近年来在科学计算领域备受关注。这项技术已有大量的科学应用，并在性能与能效方面取得了显著提升。这些应用大多由那些致力于使用 GPU 的先驱者开发和使用。最近，这项技术面临一个关键问题：它能否普遍适用于科学计算领域中各种各样的算法，以及能否为更广泛的人群而不仅是那些先驱者使用。软件开发是阻碍这项技术得以广泛应用的关键，其中包括大规模并行 CUDA 代码的编写与优化、新的性能与正确性分析工具的使用、CUDA 支持库的使用以及对 GPU 硬件架构的理解。

通过书籍、教程等方式，专家们可以将他们在此领域掌握的知识和方法与其他使用者分享，这在一定程度上解决了 CUDA 广泛应用的难题。本书就是这样的一本书。在这本书中，作者详尽地解释了重要算法的实现方法，如量子化学、机器学习以及计算机视觉等。本书不仅讲述了 GPU 编程的基本方法，还介绍了如何改写算法以便从 GPU 架构中最大化获益。此外，这本书提供了许多案例研究用以解释与充实 GPU 的重要概念，如 CUDA 线程、GPU 存储层次结构以及多 GPU 的扩展性（书中使用一个 MPI 示例例证了其扩展性可以近似线性地增至使用 500 个 GPU）。

最后，任何一种编程语言都不能独立存在。可以说，任何一个成功的编程语言，都伴随由强大的编译器、性能与正确性分析工具以及优化的支持库组成的整套系统环境。这些软件开发中的实用工具是快速开发应用程序的关键。本书的诸章节描述了如何使用 CUDA 编译器、调试器、性能分析工具、库以及与其他语言的互操作，本书在这些方面没有令人失望。

我享受从这本书中学到的一切，我确信您也会。

Jeffrey S. Vetter

美国橡树岭国家实验室杰出研究员

佐治亚理工学院教授

# 前　　言

在技术领域以及学术、职业生涯中，把握时机是极其重要的。我们这一代程序员非常幸运，能够率先利用这种经济实惠而且普及的大规模并行计算设备。GPGPU（通用图形处理单元）技术正在世界范围内影响着计算领域的各个方面——小到手机，大到超级计算机。它还在改变着商业应用、科学计算、云计算、计算机可视化系统、游戏和机器人等领域，甚至重新定义了我们所熟知的计算机编程方式。Teraflop（每秒万亿次浮点计算）如今能以一种经济可行的方式惠及全世界大多数人。无论是青少年、学生、家长、教师、专业人士，还是小型科研机构、大型公司都能负担得起 GPGPU 设备，又能获得免费的软件开发工具（SDK）。NVIDIA 公司估计已生产并销售了超过 3 亿个可编程 GPGPU 设备。

这 3 亿个用于 CUDA（Compute Unified Device Architecture，统一计算设备架构）开发的 GPGPU 设备，在商业应用领域中展现出极大的市场潜力，并且拓展了科学计算的可行范围。最重要的是，CUDA 以及大规模并行 GPGPU 设备正在改变我们对计算的认知。CUDA 程序不再是某个时刻只进行一个或少量的操作，而是同时进行数以万计的操作！

本书将介绍 CUDA 编程思想以及如何利用数以万计的线程使得应用（无论是商业还是学术、科研方面）的执行性能获得指数级的提升。进而，本书将介绍单一应用如何在单机或集群环境中使用单个或多个 GPGPU。此外，本书还将介绍如何利用 CUDA 来开发针对多核处理器平台的应用，使 CUDA 可以用于所有的应用开发场景，甚至是没有 GPU 的环境中！

本书不仅仅关注程序语法和 API 函数，还将涉及 CUDA 的设计理念，并在体系结构方面探讨 GPGPU 硬件带来惊人性能的原因。本书也将涉及 CUDA 开发中的准则和注意事项，以利于编写简洁、易读并具有可维护性的代码。本书主要以最新的 CUDA 4.x 为准。

使用与修改代码是学习过程中的必要环节，因此，本书提供了可供编译和修改的示例代码。这些示例展示了如何在 Fermi 架构的 GPGPU 设备（NVIDIA 20 系列）中编写具有高性能的代码，旨在让读者学会编写高效的代码而不是仅仅让代码能够运行。书中的示例可以在任何支持 CUDA 的 GPGPU 上编译运行，因此使用旧版 GPGPU 的读

者也可以从本书中获益。本书在适当的位置引用了我所撰写的全面 CUDA 教程之“Doctor Dobb’s Journal”（Dobb 博士的日志）专栏中的文字，来说明 CUDA 各版本之间的改进以及如何在各代 GPGPU 架构中获得高计算性能。

相关教学资料、附加示例以及读者评论在 <http://gpucomputing.net> 的 wiki 中可以找到，可以从以下任何一个网址中访问：

- 作者的名字：<http://gpucomputing.net/RobFarber>
- 本书名（连成一个单词）：<http://gpucomputing.net/CUDAapplicationdesignanddevelopment>
- 作者的专栏名：<http://gpucomputing.net/supercomputingforthemasses>

购买本书的读者还可以下载书中示例的源代码：<http://booksite.mkp.com/9780123884268>。

本书内容安排如下：

第 1 章介绍了 CUDA 的基本概念以及构建、调试 CUDA 应用所需要的工具。针对 Thrust C++ 和 C Runtime API，该章提供了一些简单的示例。此外，还介绍了三个有效提高程序性能的 GPU 编程准则。

第 2 章利用第 1 章所介绍的技术，提供了完整的机器学习以及相对于常规单核处理器加速了 341 倍的优化通用框架。同时，为了满足对领域知识和 GPU 程序开发能力有双重需求的读者，该章还介绍了机器学习与数值优化的核心概念。

第 3 章主要关注性能分析技术（profiling），它是高性能编程的一个重要技术。该章介绍了 CUDA 性能分析工具，并应用于第 2 章的案例中。而且将揭示一些 Thrust API 的惊人瓶颈。该章还讨论了数据挖掘技术以及主成分分析（PCA）和非线性主成分分析（NPCA）两种数据挖掘方法。无论是用户还是编程者，该章都是值得阅读的。

第 4 章主要介绍 CUDA 执行模型。该章的内容是任何希望获得 GPU 最高性能的人所必须了解的。为了帮助读者理解 GPU 的工作方式，该章还提供了示例以及性能分析的输出结果；这些资料有助于读者掌握如何使用现有工具来观察程序运行的内部细节。

第 5 章介绍并讨论了 CUDA 提供的多种 GPU 内存以及各种内存的优缺点。

第 6 章。鉴于最快与最慢的 GPU 内存之间有超过 3 个数量级的性能差距，因此合理利用 GPU 内存是获得高性能的唯一可行途径。该章将讨论使应用高效利用内存的技术，并通过观察性能分析输出结果理解这些技术。该章还给出了一个通用功能的示例，指导编写类似于关键 API 的基础通用函数。

第 7 章。GPU 提供了多种并行方式，包括多 GPU、异步核函数执行以及统一虚拟地址空间（Unified Virtual Address，UVA）。该章通过示例和对应的性能分析输出结果来介绍这些 GPU 并行方式的应用。

第 8 章。历经发展，CUDA 已适用于包括 GPU 和多核处理器在内的所有应用开发平台。该章将讨论多种 CUDA 后端设备，并给出可有效在异构多 GPU 环境中运行的示例以及响应的性能分析输出结果。该章还设计了 CUDA 支持库，以及如何在其他高级语言（如 Python、Java、R 和 FORTRAN 等）中应用 CUDA 与 GPU 计算接口。

第 9 章。在关注利用 CUDA 加速计算任务的同时，往往容易忽略 GPU 技术也是优秀的可视化平台。该章回归 GPU 的设计初衷，讨论 CUDA 如何显著地加速可视化与游戏应用。该章将给出一个完整的示例，使读者可以创建一个 3D 世界并在其中遨游。性能分析输出结果揭示了 CUDA 能够有效加速可视化应用的原因。该章的框架还可扩展并用到第 12 章的实时视频流的处理例子中。

第 10 章介绍使用 MPI（Message Passing Interface，通信传递接口）扩展第 3 章的示例来介绍可伸缩性以及性能。示例代码的另一版本证明了，在扩展至 500 个 GPGPU（单体最高具有每秒五十万亿次单精度浮点计算的能力）的过程中其性能近似于线性增长，最终获得了每秒超过三百万亿次浮点计算能力，相当于 60 000 个 x86 处理器核。

第 11 章。在 CUDA 技术浪潮中，还没有能覆盖所有方面的书籍。该章是一个综述章节，介绍了应用于不同技术的 CUDA 开源代码的项目，其中包括支持向量机（Support Vector Machine，SVM）、多维尺度分析（MDS）、交互信息、力导引图、分子模型等。对于这些项目的了解以及第 8 章讨论的其他高级语言接口，将有助于你成为一名成熟的 CUDA 开发者。

第 12 章。基于第 9 章的框架，该章给出了一个用于视频识别的实时视频流处理示例。仅需要一个廉价的摄像头或一个视频文件，你就可以完成实时的视频识别任务。本示例为学习 CUDA 而设计，因此很容易修改。经过扩展，本示例与技术可用于机器人、增强现实游戏，以及综合状态视窗（HUD）的数据融合，该章对此也进行了探讨。

学习 CUDA（以及 GPGPU）编程思想与开发技术是非常有趣的。然而，性能是使用 GPGPU 技术的终极目的，正如我的一位大学教授经常说的：“要知道布丁的滋味就要去品尝一下。”图 1 是 2011 年 7 月 12 日 NVIDIA CUDA Showcase 上报道的获得性能加速比前 100 名的应用。这些报告都来自于经权威评议的学术论文和商业实体，可以看出，对于许多不同种类的应用，相比于多核处理器，GPGPU 技术都能够获得 2 个数量级（100 倍）或以上的加速。展示的这些应用值得我们花时间进行探究，而且其中大部分都可以免费下载源代码或支持库。

GPGPU 是一种颠覆性的技术，它改变了人们对计算过程的认识。正如 NVIDIA 所说，“从超智能手机到超级计算机”。该技术的到来正赶上这样一场机遇风暴，传统的多核处理器已经无法再通过增加时钟频率来实现有效的加速。传统处理器生产商唯一

能吸引消费者升级计算机的方式，就是通过双核或四核的并行来获得 2 ~ 4 倍的加速。多核并行是具有颠覆性的，要求现有软件必须重写以实现对多核心的有效利用。在计算机与科研领域转向并行硬件探究之时，加入到这前沿软件应用研发中来吧！学习 CUDA，把握这一难得的机遇。

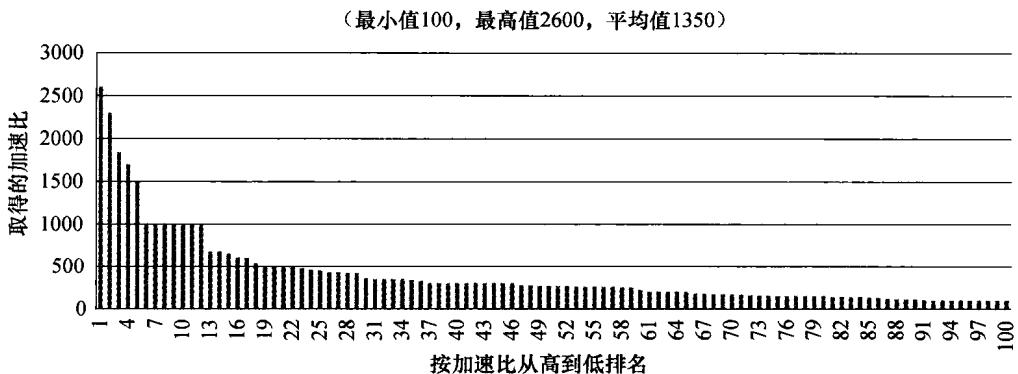


图 1 NVIDIA 展示的性能加速比前 100 名应用

## 致谢

仅以此书献给我的妻子 Margy 和我的儿子 Ryan，虽然他们没有直接帮助我写作，却在我创作本书的过程中一直深深地支持着我。尤其是我的儿子 Ryan 的成长，证明了长江后浪推前浪，感谢上帝让我陪他度过童年时光。

献给许多审阅过此书的朋友们，尤其是那些指出本书错误的朋友，对于他们付出的时间与帮助，我感激不尽。特别感谢爱尔兰高端计算中心的每个人，他们从我开始创作到完成本书手稿期间一直给予我支持。最后，感谢 NVIDIA 公司的朋友与同事，是他们将所有 CUDA 革命性的梦想变成了现实。

# 目 录

译者序

序言

前言

|   |    |
|---|----|
| 第 1 章 CUDA 入门与编程思想 .....                    | 1  |
| 1.1 源代码与维基 .....                            | 1  |
| 1.2 一个用以区别 CUDA 与传统程序开发的示例 .....            | 2  |
| 1.3 选择合适的 CUDA API .....                    | 5  |
| 1.4 CUDA 的一些基本概念 .....                      | 7  |
| 1.5 理解首个 Runtime Kernel .....               | 10 |
| 1.6 GPGPU 编程的三条法则 .....                     | 11 |
| 1.6.1 法则 1：将数据放入并始终存储于 GPU .....            | 12 |
| 1.6.2 法则 2：交给 GPGPU 足够多的任务 .....            | 12 |
| 1.6.3 法则 3：注重 GPGPU 上的数据重用，以避免带宽限制 .....    | 12 |
| 1.7 大 O 记号的思想与数据传输 .....                    | 13 |
| 1.8 CUDA 和 Amdahl 定律 .....                  | 15 |
| 1.9 数据并行与任务并行 .....                         | 15 |
| 1.10 混合执行：同时使用 CPU 和 GPU 资源 .....           | 16 |
| 1.11 回归测试与正确性 .....                         | 18 |
| 1.12 静默错误 .....                             | 19 |
| 1.13 调试简介 .....                             | 20 |
| 1.14 UNIX 调试方法 .....                        | 21 |
| 1.14.1 NVIDIA cuda-gdb 调试器 .....            | 21 |
| 1.14.2 CUDA 内存检查器 .....                     | 23 |
| 1.14.3 通过 UNIX ddd 界面使用 cuda-gdb .....      | 24 |
| 1.15 使用 Parallel Nsight 进行 Windows 调试 ..... | 25 |

|  |           |
|--|-----------|
| 1.16 本章小结 .....                                      | 27        |
| <b>第 2 章 CUDA 在机器学习与优化中的应用 .....</b>                 | <b>28</b> |
| 2.1 建模与模拟 .....                                      | 28        |
| 2.1.1 拟合参数化模型 .....                                  | 29        |
| 2.1.2 Nelder-Mead 方法 .....                           | 30        |
| 2.1.3 Levenberg-Marquardt 方法 .....                   | 30        |
| 2.1.4 算法加速 .....                                     | 31        |
| 2.2 机器学习与神经网络 .....                                  | 32        |
| 2.3 异或逻辑：一个重要的非线性机器学习问题 .....                        | 33        |
| 2.3.1 目标函数示例 .....                                   | 35        |
| 2.3.2 针对多 GPU 设备、多 CPU 处理器的完整仿函数 .....               | 35        |
| 2.3.3 完整 Nelder-Mead 优化代码的简要讨论 .....                 | 37        |
| 2.4 异或逻辑的性能结果 .....                                  | 45        |
| 2.5 性能讨论 .....                                       | 45        |
| 2.6 本章小结 .....                                       | 48        |
| 2.7 C++ NELDER-MEAD 代码模板 .....                       | 48        |
| <b>第 3 章 CUDA 工具套件：对 PCA、NLPCA 进行性能分析 .....</b>      | <b>53</b> |
| 3.1 PCA 和 NLPCA .....                                | 53        |
| 3.1.1 自编码网络 .....                                    | 55        |
| 3.1.2 用于 PCA 分析的仿函数示例 .....                          | 56        |
| 3.1.3 用于 NLPCA 分析的示例仿函数 .....                        | 58        |
| 3.2 获得基础性能分析数据 .....                                 | 60        |
| 3.3 gprof：通用 UNIX 性能分析器 .....                        | 61        |
| 3.4 NVIDIA 可视化性能分析器：computeprof .....                | 62        |
| 3.5 Microsoft Visual Studio 中的 Parallel Nsight ..... | 65        |
| 3.5.1 Nsight 时间表分析 .....                             | 66        |
| 3.5.2 NVTX 跟踪支持库 .....                               | 67        |
| 3.5.3 CUDA API 的可扩展性表现 .....                         | 68        |
| 3.6 性能调节与分析实用工具（TAU） .....                           | 70        |
| 3.7 本章小结 .....                                       | 70        |

|   |     |
|---|-----|
| 第 4 章 CUDA 执行模型 .....                   | 72  |
| 4.1 GPU 架构综述 .....                      | 72  |
| 4.1.1 线程调度：通过执行配置统筹性能与并行度 .....         | 74  |
| 4.1.2 computeprof 中 Warp 相关值 .....      | 77  |
| 4.1.3 Warp 分歧 .....                     | 77  |
| 4.1.4 关于 Warp 分歧的若干准则 .....             | 78  |
| 4.1.5 computeprof 中 Warp 分歧相关值 .....    | 79  |
| 4.2 Warp 调度与 TLP .....                  | 79  |
| 4.3 ILP：高性能低占用率 .....                   | 80  |
| 4.3.1 ILP 隐藏算术计算延迟 .....                | 81  |
| 4.3.2 ILP 隐藏数据延迟 .....                  | 84  |
| 4.3.3 ILP 的未来 .....                     | 84  |
| 4.3.4 computeprof 中指令速率相关值 .....        | 85  |
| 4.4 Little 法则 .....                     | 86  |
| 4.5 检测限制因素的 CUDA 工具 .....               | 87  |
| 4.5.1 nvcc 编译器 .....                    | 88  |
| 4.5.2 启动约束 .....                        | 90  |
| 4.5.3 反汇编器 .....                        | 90  |
| 4.5.4 PTX Kernel 函数 .....               | 92  |
| 4.5.5 GPU 模拟器 .....                     | 92  |
| 4.6 本章小结 .....                          | 93  |
| 第 5 章 CUDA 存储器 .....                    | 94  |
| 5.1 CUDA 存储器层次结构 .....                  | 94  |
| 5.2 GPU 存储器 .....                       | 95  |
| 5.3 L2 缓存 .....                         | 98  |
| 5.4 L1 缓存 .....                         | 99  |
| 5.5 CUDA 内存类型 .....                     | 100 |
| 5.5.1 寄存器 .....                         | 101 |
| 5.5.2 局域内存 .....                        | 101 |
| 5.5.3 和局域内存相关的 computeprof 性能分析参数 ..... | 102 |

|   |            |
|---|------------|
| 5.5.4 共享内存 .....                        | 102        |
| 5.5.5 和共享内存相关的 computeprof 性能分析参数 ..... | 105        |
| 5.5.6 常量内存 .....                        | 105        |
| 5.5.7 纹理内存 .....                        | 106        |
| 5.5.8 和纹理内存相关的 computeprof 性能分析参数 ..... | 108        |
| 5.6 全局内存 .....                          | 109        |
| 5.6.1 常见的整合内存示例 .....                   | 110        |
| 5.6.2 全局内存的申请 .....                     | 111        |
| 5.6.3 全局内存设计中的限制因素 .....                | 113        |
| 5.6.4 和全局内存相关的 computeprof 性能分析参数 ..... | 114        |
| 5.7 本章小结 .....                          | 115        |
| <b>第 6 章 高效使用 CUDA 存储器 .....</b>        | <b>116</b> |
| 6.1 归约 .....                            | 116        |
| 6.1.1 归约模板 .....                        | 117        |
| 6.1.2 functionReduce.h 的测试程序 .....      | 122        |
| 6.1.3 测试结果 .....                        | 126        |
| 6.2 使用非规则数据结构 .....                     | 127        |
| 6.3 稀疏矩阵和 CUSP 支持库 .....                | 131        |
| 6.4 图论算法 .....                          | 132        |
| 6.5 SoA、AoS 以及其他数据结构 .....              | 134        |
| 6.6 分片和分块 .....                         | 135        |
| 6.7 本章小结 .....                          | 136        |
| <b>第 7 章 提高并行度的技巧 .....</b>             | <b>137</b> |
| 7.1 CUDA 上下文环境对并行度的扩展 .....             | 137        |
| 7.2 流与上下文环境 .....                       | 138        |
| 7.2.1 多 GPU 的使用 .....                   | 139        |
| 7.2.2 显式同步 .....                        | 139        |
| 7.2.3 隐式同步 .....                        | 141        |
| 7.2.4 统一虚拟地址空间 .....                    | 141        |
| 7.2.5 一个简单的示例 .....                     | 142        |

|  |     |
|--|-----|
| 7.2.6 分析结果 .....                       | 144 |
| 7.3 使用多个流乱序执行 .....                    | 144 |
| 7.3.1 在同一 GPU 内并发执行 Kernel 函数的建议 ..... | 147 |
| 7.3.2 隐式并行 Kernel 的原子操作 .....          | 147 |
| 7.4 将数据捆绑计算 .....                      | 149 |
| 7.4.1 手动分割数据 .....                     | 150 |
| 7.4.2 映射内存 .....                       | 150 |
| 7.4.3 映射内存的工作机制 .....                  | 152 |
| 7.5 本章小结 .....                         | 153 |
| 第 8 章 CUDA 在所有 GPU 与 CPU 程序中的应用 .....  | 154 |
| 8.1 从 CUDA 到多种硬件后端的途径 .....            | 155 |
| 8.1.1 PGI CUDA x86 编译器 .....           | 155 |
| 8.1.2 PGI CUDA x86 编译器 .....           | 157 |
| 8.1.3 将 x86 处理器核心用作流多处理器 .....         | 159 |
| 8.1.4 NVIDIA NVCC 编译器 .....            | 160 |
| 8.1.5 Ocelot .....                     | 160 |
| 8.1.6 Swan .....                       | 161 |
| 8.1.7 MCUDA .....                      | 162 |
| 8.2 从其他语言访问 CUDA .....                 | 162 |
| 8.2.1 SWIG .....                       | 162 |
| 8.2.2 Copperhead .....                 | 163 |
| 8.2.3 EXCEL .....                      | 164 |
| 8.2.4 MATLAB .....                     | 164 |
| 8.3 支持库 .....                          | 164 |
| 8.3.1 CUBLAS .....                     | 164 |
| 8.3.2 CUFFT .....                      | 165 |
| 8.3.3 MAGMA .....                      | 174 |
| 8.3.4 phiGEMM 支持库 .....                | 175 |
| 8.3.5 CURAND .....                     | 176 |
| 8.4 本章小结 .....                         | 177 |

|  |     |
|--|-----|
| 第 9 章 CUDA 与图形渲染混合编程 .....                 | 178 |
| 9.1 OpenGL .....                           | 178 |
| 9.1.1 GLUT .....                           | 179 |
| 9.1.2 通过 OpenGL 映射 GPU 内存 .....            | 179 |
| 9.1.3 使用基元重启提升 3D 处理性能 .....               | 181 |
| 9.2 框架内各文件的介绍 .....                        | 183 |
| 9.2.1 Kernel 与 Perlin Kernel 演示的示例代码 ..... | 184 |
| 9.2.2 simpleGLmain.cpp 文件 .....            | 192 |
| 9.2.3 simpleVBO.cpp 文件 .....               | 196 |
| 9.2.4 callbacksVBO.cpp 文件 .....            | 199 |
| 9.3 本章小结 .....                             | 204 |
| 第 10 章 在云计算和集群环境中使用 CUDA .....             | 205 |
| 10.1 消息传递接口 .....                          | 205 |
| 10.1.1 MPI 编程模型 .....                      | 206 |
| 10.1.2 MPI 通信器 .....                       | 206 |
| 10.1.3 MPI 进程号 .....                       | 206 |
| 10.1.4 主从模式 .....                          | 208 |
| 10.1.5 点对点模式基础 .....                       | 208 |
| 10.2 MPI 通信机制 .....                        | 209 |
| 10.3 带宽 .....                              | 211 |
| 10.4 平衡率 .....                             | 212 |
| 10.5 运行大型 MPI 程序需要考虑的因素 .....              | 214 |
| 10.5.1 初始数据加载的可扩展性 .....                   | 214 |
| 10.5.2 使用 MPI 进行计算 .....                   | 215 |
| 10.5.3 可扩展性检查 .....                        | 216 |
| 10.6 云计算 .....                             | 217 |
| 10.7 代码示例 .....                            | 218 |
| 10.7.1 数据的产生 .....                         | 218 |
| 10.7.2 主体代码部分 .....                        | 220 |
| 10.8 本章小结 .....                            | 225 |

|   |     |
|---|-----|
| 第 11 章 CUDA 在现实问题中的应用 .....               | 227 |
| 11.1 高维数据的处理 .....                        | 228 |
| 11.1.1 PCA/NLPCA .....                    | 228 |
| 11.1.2 多维尺度分析 .....                       | 229 |
| 11.1.3 K 均值聚类算法 .....                     | 229 |
| 11.1.4 期望最大化 .....                        | 229 |
| 11.1.5 支持向量机 .....                        | 230 |
| 11.1.6 Bayesian 网络 .....                  | 230 |
| 11.1.7 互信息 .....                          | 231 |
| 11.2 力导向图 .....                           | 232 |
| 11.3 Monte Carlo 方法 .....                 | 232 |
| 11.4 分子建模 .....                           | 233 |
| 11.5 量子化学 .....                           | 234 |
| 11.6 交互式工作流 .....                         | 234 |
| 11.7 其他众多的项目 .....                        | 235 |
| 11.8 本章小结 .....                           | 235 |
| 第 12 章 针对现场实况视频流的应用程序 .....               | 236 |
| 12.1 机器视觉话题 .....                         | 236 |
| 12.1.1 3D 效果 .....                        | 237 |
| 12.1.2 肤色区域分割 .....                       | 238 |
| 12.1.3 边缘检测 .....                         | 238 |
| 12.2 FFmpeg .....                         | 239 |
| 12.3 TCP 服务器 .....                        | 241 |
| 12.4 实况视频流应用程序 .....                      | 244 |
| 12.4.1 kernelWave(): 动画 Kernel 函数 .....   | 244 |
| 12.4.2 kernelFlat(): 在平面渲染图像 .....        | 245 |
| 12.4.3 kernelSkin(): 仅保留肤色区域 .....        | 245 |
| 12.4.4 kernelSobel(): Sobel 边缘检测过滤器 ..... | 246 |
| 12.4.5 launch_kernel() 方法 .....           | 247 |
| 12.5 simpleVBO.cpp 文件 .....               | 248 |