

研究生教育书系
信息与电子学科

流处理器 研究与设计

张春元
梅楠义
文伍何

等著



電子工業出版社

PUBLISHING HOUSE OF ELECTRONICS INDUSTRY

<http://www.phei.com.cn>

研究生教育书系
信息与电子学科

流处理器研究与设计

张春元 文 梅 伍 楠

何 义 任 巨 管茂林 著



電子工業出版社
Publishing House of Electronics Industry
北京 · BEIJING

内 容 简 介

本书介绍了一种新型的非冯·诺依曼体系结构——流体体系结构。作者在前人的研究基础上，结合自己多年科研工作的体会，介绍了流处理的主要思想、流体系结构及其运行机制、编程模型及编译器设计，以 JPEG 和 H.264 等典型应用为例详述了应用的流化方法；并叙述了多核流体系结构设计、程序设计与编译及 VLSI 特性等多个方面的内容；最后就流体系结构的未来发展进行了讨论。本书在介绍流体系结构这一专业领域的知识和技术时，秉承实事求是的科研精神，力求做到由浅入深、文字流畅、便于阅读。

本书可作为从事处理器体系结构设计与开发的科研人员和广大爱好者的参考书，也可作为大专院校计算机相关专业本科生、研究生的教材或参考书。

未经许可，不得以任何方式复制或抄袭本书之部分或全部内容。

版权所有，侵权必究。

图书在版编目 (CIP) 数据

流处理器研究与设计 / 张春元等著。—北京：电子工业出版社，2009.4

（研究生教育书系·信息与电子学科）

ISBN 978-7-121-08487-4

I. 流… II. 张… III. 微处理器—研究生—教材 IV. TP332

中国版本图书馆 CIP 数据核字（2009）第 035098 号

责任编辑：刘 凡

印 刷：北京机工印刷厂

装 订：三河市鹏成印业有限公司

出版发行：电子工业出版社

北京市海淀区万寿路 173 信箱 邮编：100036

开 本：720×1000 1/16 印张：18 字数：363 千字

印 次：2009 年 4 月第 1 次印刷

印 数：3000 册 定价：35.00 元

凡所购买电子工业出版社图书有缺损问题，请向购买书店调换。若书店售缺，请与本社发行部联系，联系及邮购电话：(010)88254888。

质量投诉请发邮件至 zlts@phei.com.cn，盗版侵权举报请发邮件至 dbqq@phei.com.cn。

服务热线：(010)88258888。

前 言

本书的由来

基于冯·诺依曼体系结构的传统程序设计模型和体系结构模型已经很好地运行了六十多年。但是，这个模型在当前 VLSI 工艺下，可靠性、功耗、计算效率、成本等都存在问题。全世界大量从事计算机系统研究、设计和制造的科学家和工程师们一直在设计新一代计算机体系结构方面做着不懈的探索。

近二十年来，国防科技大学计算机学院有很多老师带领研究生们一直专注于新型体系结构的研究，并且取得过很多创新性成果。我们课题小组的主要成员在 20 世纪 90 年代初期攻读博士、硕士学位时，就开始关注国际上新型并行计算机体系结构的研究工作；随后的 10 年间，承担和参与了一系列国家自然科学基金项目、863 项目和国家重点工程项目，获得了多项国家和部委级科技成果奖。2000 年以后，我们课题组的老师和研究生们把学术研究的重点集中在高并行处理器技术领域。经过一段时间的学习和研究，到 2003 年，我们对该领域的研究进行了总结，课题组的研究重点逐步聚焦到以下 4 种并行结构：

- (1) 以传统处理器为基础的片内多核结构；
- (2) 以多向量处理器部件为基础的结构；
- (3) 以阵列处理器为基础的结构；
- (4) 以处理（计算）单元阵列为基础的结构。

当时，我们认为（1）和（3）两种结构面临的问题与现有的 MPP 相似，主要技术壁垒在于如何提高计算资源的使用效率，在当前程序设计模式下，要攻克这个壁垒是比较困难的。我们在结构（2）和（4）方面做了更为前沿的探索，资料表明，这种探索的结果是有吸引力的。国防科技大学计算机学院研制的银河-I 和银河-II 巨型计算机，都采用了向量处理器的结构，在向量计算机系统研究和应用程序向量

化等关键技术方面都有比较好的技术积累，同时对向量处理技术中存在的难点问题也有更加深刻的体会。我们课题组在讨论了学术新颖性和研究复杂性后，选择了结构（4）作为研究重点，而作为重点研究的方向则慎重地选择了以处理单元阵列为基础的新型计算模型——流计算模型。

美国斯坦福大学 Imagine 原型系统是以流计算模型为基础的流处理器经典代表。Imagine 在处理器的适用性、功耗、程序效率、单位成本等因素上做了很好的权衡，其公布的研究成果非常有吸引力，鼓励我们更加深入地开展流体系结构的研究。2004 年，我们申请了自然科学基金，得到 2005 年度探索项目的资助。在此期间，我们不但阅读了能够得到的所有资料，研究了流处理器的体系结构和程序设计工具链，而且启动了一个基于 FPGA 的 32 位流处理器验证系统的设计，该处理器被命名为“MASA”（Multi-dimension Adaptive Stream Architecture）。我们同时启动了 MASA 模拟器的设计和实现以及配套工具链的构造，并将课题组命名为“MASA 课题组”。到 2005 年秋，我们课题组成功地在 Altera 的 StratixII EP2S180 上实现了一个完整的 32 位 MASA 处理器演示原型 MASA-I，其清晰的结构和主要测试结果与当初的模拟结果非常接近。

MASA 课题组的研究工作得到了国防科技大学计算机学院的重视和支持。从 2004 年开始，我们加入了学院“高性能联合博导组”。2005 年 5 月，MASA 课题组主要成员又参加了学院承担的国家重点工程——“飞腾 64”（FT64）。为配合工程需要，从 2005 年 7 月下旬开始，MASA 课题组 4 名成员不舍昼夜，用两周的时间完成了共 300 余页计 40 余万字的技术研究报告，系统全面地总结了在 MASA 研究中取得的成果，对“飞腾 64”的研制起到了积极的促进作用。课题组的主要成员也因此而获得了科技进步奖。

结合应用需求和学术创新的考虑，我们在前期研究的基础上，构造了多个流处理器原型，主要包括面向图像处理的 MASA-M1、面向多线程的 MASA-SMT、面向大规模扩展的 TiSA（Tiled-MASA）等。

从 2005 年开始，MASA 课题组在流体系结构研究领域累计申请了 4 项自然科学基金项目和 1 项 863 高技术计划项目；在有关学术会议和杂志上发表了一系列的学术论文并被检索，也得到国内外同行专家的认可。国内有很多关注流体系结构的老师和研究生与我们进行了大量的学术交流，也使我们对在中国开展流体系结构的研究更加充满信心。2007 年，我们有幸遇到电子工业出版社高等教育分社的张濮先生和凌毅女士，并得到他们的建议和鼓励，把研究成果形成学术专著，以供国内从事流体系结构计算机研究的同行参考。在张先生和凌女士的鼓励和指导下，MASA 课题组经过近一年的努力，终于完成了这本专著。

关于本书

本书的主要宗旨就是介绍流计算模型的主要思想和流体系结构。流（stream）计算模型是基于“生产者-消费者”的计算模型，其产生的根源是媒体处理的大量出现，新的应用需求引发了对新型计算模型的思考和完善，随着人们对媒体处理领域计算特征的认识，流计算模型也不断发展和逐渐明晰起来。

流体系结构一般采用大规模的运算阵列和多级存储，通过片内并行和专用指令，从体系结构对应用高效支持。与通用处理器体系结构相比，这种结构大幅简化了片内控制，拥有大量的可编程运算功能单元，并且支持通信信道可编程，能够较好地缓解现有通用处理器面临的计算压力，特别适应 VLSI 技术的发展。进一步的研究表明，流处理器在媒体应用、信息处理等多种计算密集型应用领域可获得与专用芯片相当的性能。新兴流处理器包括学术界研究的原型系统 Imagine、Merrimac、VIRAM、RAW、SCORE、FT64、MASA 等，工业界的则有 Cell、NVIDIA G80、AMD-ATI R580、Storm DSP、Tiles64 等，这些流处理器都得到了广泛的关注。可以预计在未来单芯片超十亿只晶体管的时代，流体系结构将成为体系结构发展的一种主流方向。

流体系结构与应用结合十分紧密，因此本书还介绍了流应用特征、流编程编译的基本方法。此外，由于流体系结构起源于媒体处理，在面向更广泛的应用领域如科学计算时，有许多方面值得深入研究和改进，所以本书还结合课题研究，阐述了我们对未来流处理器的发展和研究方向的一些观点。全书比较全面地反映了我们在流计算领域从体系结构、理论方法、模拟系统、测试程序、芯片设计、编译、编程环境到典型应用等方面取得的研究成果。

全书共分 8 章，内容安排如下：

第 1 章主要阐述了流体系结构的产生背景、研究现状及发展趋势，并且介绍了当前几种新兴的流处理器。

第 2 章介绍了流体系结构的一些基本概念和基于流处理的硬件结构模型，并以实例说明了流处理器与向量处理器的区别。

第 3 章以 MASA 流体系结构为例，剖析了一个完整的流处理器核体系结构，深入研究流处理器内部硬件结构，包括指令集实例、流水线及各功能模块的设计等，并且对提高性能的优化设计进行了讨论。

第 4 章针对流体系结构系统运行模式和协同机制进行研究探讨，重点介绍了标量处理器核和流处理器核之间的软/硬件协同控制机制，包括机制的组成及协同过程。

第 5 章介绍流体系结构的编程模型及编译器设计。首先介绍其编程模型及与之

配套的 KernelC/StreamC 语言，然后介绍其编译器原理及设计，最后介绍一些目前国际上已有的其他流编程模型及编译环境。

第 6 章先介绍了流应用的概念和领域，并以 JPEG 编码算法为例阐述了应用的基本流化方法；然后分别讨论了媒体处理领域中的 H.264 视频编码和科学计算领域中的梅森素数求解这两个应用的流化过程和流化效果。

第 7 章从流体系结构多维可扩展性出发，首先介绍了流体系结构簇内扩展、簇间扩展和多核扩展；然后重点介绍了多核流体系结构设计、多核流体系结构程序设计与编译以及多核流处理器 VLSI 特性等多个方面的内容。

第 8 章从流处理器发展的角度，重点讨论流处理器未来可研究的工作。

本书由 MASA 课题组成员合作完成，由张春元教授、文梅副教授策划和统筹，并与伍楠、何义、任巨和管茂林 4 位博士（生）共同执笔完成。第 1 章由张春元撰写，第 2 章由张春元和文梅撰写，第 3 章由张春元、伍楠撰写，第 4 章由文梅和何义撰写，第 5 章由文梅、任巨和管茂林撰写，第 6 章由文梅和任巨撰写，第 7 章由何义和伍楠撰写，第 8 章由张春元和何义撰写。杨乾明、荀长庆、吴伟、柴俊、苏华友和全巍等硕士研究生收集和整理了大量的资料，提供了良好的素材，并参与了部分章节的编写。本书写作过程中，参阅了国内外许多作者的论文及著作，特别是 Rixner 的专著，参考了其中部分材料，在此深表谢意。

流计算模型和流体系结构的研究仍在发展之中。本书试图在前人的研究基础上，结合我们多年科研工作的体会，向同行介绍流处理器和我们的研究成果。由于作者的能力和知识面有限，疏忽、不当和错误难免，恳请读者批评指正。

有关 MASA 课题组的情况，读者可参见我们课题组的主页 <http://masa.nudt.edu.cn>。

MASA 课题组 张春元
于湖南长沙国防科技大学
2009 年 2 月

目 录

第 1 章 绪论	(1)
1.1 VLSI 技术的发展对处理器体系结构的影响	(1)
1.2 应用对体系结构提出的新要求	(2)
1.3 高性能体系结构面临新的挑战	(3)
1.3.1 专用处理器	(5)
1.3.2 通用微处理器	(5)
1.3.3 DSP 与可编程的媒体处理器	(6)
1.4 新兴的流处理器	(6)
1.4.1 Imagine 和 Merrimac	(7)
1.4.2 CELL 处理器	(8)
1.4.3 基于片上存储的 VIRAM 体系结构	(9)
1.4.4 片内多处理机体系结构的代表：RAW 和 TRIPS	(10)
1.4.5 流计算模型 SCORE	(12)
1.4.6 新型流体体系结构小结	(12)
第 2 章 流处理	(17)
2.1 流处理思想	(17)
2.2 硬件结构模型	(19)
2.2.1 解耦合计算和访存	(19)
2.2.2 多级存储层次	(20)
2.3 流处理实例及与向量处理的比较	(25)
2.4 小结	(29)

第3章	流处理器微体系结构	(31)
3.1	流体系结构设计思想	(31)
3.1.1	控制子系统	(32)
3.1.2	存储子系统	(33)
3.1.3	计算子系统	(34)
3.1.4	对外接口	(34)
3.2	流处理器的指令集设计技术	(35)
3.2.1	流级指令	(35)
3.2.2	kernel 级指令	(39)
3.3	流处理器的流水线设计技术	(44)
3.3.1	核心指令执行流水线的组织结构	(45)
3.3.2	流数据访问流水线的组织结构	(46)
3.3.3	流水线的数据通路及其相关处理	(46)
3.4	流处理器计算子系统的设计	(49)
3.4.1	簇内寄存器文件系统	(52)
3.4.2	簇内交叉互连开关	(56)
3.4.3	运算单元 ALU 和 DSQ	(56)
3.4.4	计算簇内便笺存储器设计	(59)
3.4.5	计算簇间的通信单元设计	(60)
3.4.6	流 IO 单元设计	(62)
3.5	流处理器控制子系统的设计	(63)
3.5.1	核心指令控制逻辑	(63)
3.5.2	流指令控制逻辑	(70)
3.5.3	标量处理器核的控制逻辑	(75)
3.6	流处理器存储子系统的设计	(76)
3.6.1	流寄存器文件	(76)
3.6.2	Cache	(81)
3.6.3	DRAM 存储器	(81)
3.7	流处理器核对外接口的设计	(86)
3.7.1	与标量处理器核的接口	(86)
3.7.2	多片流处理器核互连接口	(87)
3.7.3	片上总线接口	(92)
3.8	提高性能的优化设计讨论	(92)

3.8.1 条件流	(92)
3.8.2 片上索引流	(99)
3.8.3 软件流水	(101)
3.9 流处理器实例	(103)
3.9.1 用于数字媒体处理的流处理器 Storm	(103)
3.9.2 用于科学计算的流处理器 FT64	(104)
3.10 小结	(105)
第 4 章 流处理器协同机制	(107)
4.1 概述	(107)
4.2 软件协同控制机制	(108)
4.2.1 软件协同的工作原理	(109)
4.2.2 协处理中间件	(111)
4.3 硬件协同控制机制	(119)
4.3.1 异构核协同单元	(120)
4.3.2 协处理中间件	(121)
4.3.3 异构核协同任务的实现	(121)
4.4 异构核间的数据传输	(123)
4.4.1 分离片外存储空间	(123)
4.4.2 共享片外存储空间	(127)
4.5 小结	(129)
第 5 章 流编程模型与编译器	(131)
5.1 流编程模型	(131)
5.1.1 流级程序	(132)
5.1.2 核心程序	(134)
5.1.3 StreamC/KernelC 语言	(135)
5.1.4 层次化的流编译器结构	(141)
5.2 核心级编译器设计	(143)
5.2.1 核心级编译器流程	(143)
5.2.2 VLIW 调度	(145)
5.2.3 通信调度	(148)
5.2.4 寄存器分配	(151)
5.3 流级编译器设计	(158)
5.3.1 流级编译概述	(158)

5.3.2	流级编译的流程	(159)
5.3.3	流调度	(161)
5.4	其他流编程语言与编译环境	(164)
5.5	小结	(167)
第6章	流应用与编程	(169)
6.1	流应用概述	(169)
6.2	流应用的映射	(170)
6.2.1	程序特征分析	(170)
6.2.2	划分数据流图	(172)
6.2.3	计算核心的编写	(173)
6.2.4	流组织及其常用方法	(175)
6.3	媒体应用实例：高清 H.264 编码	(178)
6.3.1	H.264 简介	(178)
6.3.2	流化平台介绍	(180)
6.3.3	H.264 的流化实现	(183)
6.3.4	H.264 的流化结果分析	(199)
6.4	科学计算实例：梅森素数求解法	(201)
6.4.1	梅森素数及其求解方法 LUCAS	(201)
6.4.2	LUCAS-Lehmer 法的流式实现	(202)
6.4.3	LUCAS 流式算法在 FT64 上的性能分析和评测	(214)
6.5	小结	(215)
第7章	多核流处理器设计	(217)
7.1	流体系结构多维可扩展性	(217)
7.1.1	簇内扩展	(219)
7.1.2	簇间扩展	(220)
7.1.3	多核扩展	(221)
7.2	多核流体系结构	(224)
7.2.1	多核流体系结构的顶层硬件结构	(224)
7.2.2	流传输协议	(228)
7.2.3	流互连网络模块	(236)
7.3	TiSA 多核流程序的设计与编译	(239)
7.3.1	面向多核的流程序设计	(240)
7.3.2	多核流程序的静态编译	(243)

7.4	流体系结构的多维扩展代价与性能	(249)
7.4.1	VLSI 开销	(249)
7.4.2	性能效率	(257)
7.5	小结	(260)
第 8 章	未来流处理器研究	(261)
8.1	非规则流研究	(261)
8.2	编程模式研究	(263)
8.3	基于流模型的多核调度	(264)
8.4	通用领域的发展	(265)
8.5	其他的研究点	(266)
8.6	小结	(266)
参考文献	(267)
缩略语表	(273)

第 1 章

绪论

计算机系统设计、制造和应用的发展，使我们正处在计算机体系结构大变革的边缘。在过去的 50 年里，冯·诺依曼体系结构一直被不断更新。而在下一个十年里，体系结构会出现较大的变革^[1]。流体系结构（Stream Architecture）就是这样一种非冯·诺依曼体系结构，它能高效开发应用并行性，达到可与专用处理器相比的高性能；同时它还具有可编程性，因此被工业界和学术界广泛关注和采用。本章主要简述流体系结构的产生背景，并介绍当前几种新兴的流处理器。

1.1 VLSI 技术的发展对处理器体系结构的影响

摩尔定律相当准确地预言了芯片上可集成的晶体管数目的增长规律。到2008年，单芯片上已经发展到可以放置十亿只晶体管，设计者可以将大量的运算单元集成在一个芯片上。在 $0.15\text{ }\mu\text{m}$ CMOS 工艺下，一个 32 位整数加法器占用的芯片面积还不到 0.05 mm^2 ，而单芯片可以集成上百个 1 GHz 的浮点单元，总体性能超过了 100 GFLOPS/片^[2]。在线宽缩小的同时，计算功耗也在降低。例如，在 Imagine 处理器中^[3]，以 $0.18\text{ }\mu\text{m}$ 工艺制造的单精度浮点乘加单元占用了 0.486 mm^2 ，而每个乘法操作只耗能 185 pJ (0.185 mW/MHz)。计算成本相对来说越来越低。目前，约 100 GFLOPS 和超过 1 TOPS（渲染）能力的图形芯片，其价格还不到 100 美元，如 NVIDIA 的 GeForce4 处理器，其性能达到 120 GFLOPS 和 1.2 TOPS^[4]。嵌入式处理器尽管性能没有那么强大，但价格便宜，原始的 1 GFLOPS 的成本不到 1 美元。

但与此同时，片内、片外的通信延迟、带宽和功耗却与运算单元的大规模集成难以匹配。随着线宽缩小，线延迟与门延迟相当，这成为制约频率的关键因素，高负载长线的功耗也变得不可忽视。以多端口存储器的访问为例，假定每个乘法需要三次访问多端口存储器，且使用三个 5 mm 的总线（2 读 1 写）进行数据传输，每

次驱动 32 位 5 mm 的总线来传送数据平均耗费 24 pJ，通信的代价将与乘法的代价在同一个数量级。片外通信更是一种关键资源，即使采用现在最新的封装方式，芯片上最多也只能引出大约 1000 个引脚，这极大地限制了片外数据带宽。并且，片外通信也耗费了大量的能量（每 32 位的数据传送的耗费大于 1nJ）^[5]。这就是现代 VLSI 技术的一个典型特征：运算单元相对廉价，而运算单元之间的通信较昂贵。

1.2 应用对体系结构提出的新要求

随着科学和技术的进步，人们对计算性能的追求永无止境。在科研、国防、商务、娱乐等众多领域，有一类典型的应用——流应用正成为微处理器的主要负载，日益引起人们的关注。所谓流，就是指不间断的、连续的、移动的数据队列，队列长度可以是定长或不定长的，流元素的组成可以是复杂或简单的。流应用主要分为以下两类：

（1）数字信号处理应用，实时地处理信号、音频、视频、静态图像及其他密集型数据（Data-intensive），典型的应用包括图像处理、图形处理、视频编码与解码、雷达处理等。

（2）科学计算，主要用于科学模型的建立和模拟，典型的应用包括流体力学、气象、分子动力学、有限元方法问题等。

流应用具有以下几个主要特征。

（1）计算密集性：与传统的桌面应用相比，流应用对每次从内存取出的数据都要进行大量的算术运算，这称为计算密集性。计算密集性可以用计算访存比（对每次访存数据执行的计算操作个数）来量化，流应用的计算访存比一般在 20 以上，这意味着对每次从内存取出的数据都要进行大量的算术运算。典型流应用的计算访存比见表 1.1。

表 1.1 典型流应用的计算访存比

计算密集型应用	计算访存比
深度提取的卷积过滤与绝对值求和	473.3
视频的编码与解码 MPEG	57.9
矩阵的 QR 分解	155.3
Reed-Solomon 解码	129.6
二维拉格朗日和欧拉结合法求解爆轰流体力学问题 Ygx2 (IAPCM Benchmarks)	76

（2）并行性：大多数流应用中的运算可以并行，并且以数据级并行为主，同时还存在指令级和任务级并行。例如，流体力学中解恒定流场的偏微分方程时，各个通量可以并行计算，同时每个通量的各个结点也存在数据并行。

(3) 局域性：局域性可以分为数据重用局域性与生产者-消费者局域性。数据重用局域性是指在执行计算核心程序时重复使用系数与数据，如卷积过滤。生产者-消费者的局域性体现在计算流水线的不同阶段中，即生产出的数据（即一个计算核心程序写的数据）被另一个计算核心程序所消费（被读取），且不会再被生产数据的计算核心程序使用。

现有的用于媒体处理和流式科学计算领域的处理器，如通用处理器、DSP、向量微处理器、专用图形处理器和片内多处理机等，都存在各自的局限性。例如，在通用微处理器中，数据重用局域性通常被寄存器文件或小容量的一级 Cache 所捕获。而生产者-消费者局域性则不容易被传统的存储层次所捕获，因为它并不符合 LRU (Least-Recently-Used) 规则。

表 1.1 列出的 5 个流应用的计算访存比表明：深度提取在卷积过滤与绝对值求和过程中对每个访存的数据需要 473.3 个算术操作，其他应用对每个访存数据需要做 57.9~155.3 个算术操作^[1]。比较而言，在 SPECint2000 测试平台上，传统的桌面整数应用的算术运算只占指令总数的 2%~50%，而访存指令则占 15%~80%^[6]。这些差别说明：为桌面整数运算而设计的体系结构（如通用微处理器）并不适合流应用。

1.3 高性能体系结构面临新的挑战

随着 VLSI 技术的发展和应用需求的提高，体系结构研究人员面临的问题是：如何发展新的体系结构技术来有效地利用这些丰富而廉价的资源，并在芯片上提供更大规模的高效计算能力。换言之，就是如何有效地利用芯片面积来计算。关于这一点，相对于通用微处理器而言，图形处理器和 DSP 做得更好。一个现代的高端图形芯片拥有超过 64 个浮点 ALU 和 1000 个整数 ALU，运算密度达到通用微处理器的 100 倍^[4]。而通用的微处理器，如 Itanium，其运算单元只占芯片面积的 6.5%^[7]，大部分芯片面积都被用于实现 Cache、分支预测、乱序执行、指令/通信调度等功能。在运行密集型运算任务时，具有相近物理指标（主要是工艺、面积、功耗和元件数）的通用微处理器与专用芯片相比，运算单元的低占有率带来了运算速度上 2~3 个数量级的巨大差别（1 GOPS 的通用处理器和 1 TOPS 的图形芯片）。

高密度 VLSI 给设计人员也带来了新的问题：芯片设计的复杂度越来越难以控制，而且研制周期过长，芯片面积利用率也不高，过于专用的芯片体系结构难以适应灵活多变的市场需要，如何给众多的 ALU 提供足够的指令和数据，如何合理地消除带宽瓶颈，片上密度增加带来的连线延迟，低功耗需求，等等。

如图1.1^[8]所示为微处理器随的性能（每条指令的执行时间）增长示意图。从图

中可见，随着集成电路制造工艺的发展，从 2000 年开始，微处理器在工艺上仅仅可以获得大约每年 20% 的性能增长。要保持微处理器性能继续以每年 50% 甚至更高的速度增长，就必须在微处理器的并行处理体系结构上进行创新。先进的体系结构可以使微处理器性能更快地增长，与传统体系结构的微处理器的性能差异将逐渐扩大，到 2010 年将达到约 1000 倍。

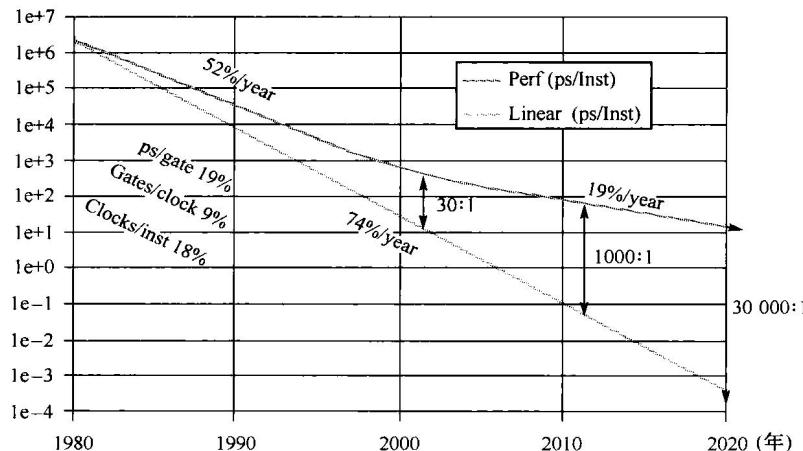


图 1.1 微处理器性能增长示意图

目前，传统微体系结构对并行能力的开发面临一系列问题。国际上体系结构研究的前沿之一，就是面向十亿只以上晶体管级别芯片的新型微体系结构研究。20 世纪 70 年代提出的因条件限制而无法有效实现和发展的脉动阵列（systolic arrays）和 20 世纪 80 年代的数据流（data flow）的核心思想得到重视和应用，20 世纪 90 年代被 RISC 挤出人们视野的向量处理等技术又焕发了新的生命力。同时，针对各种具体的应用领域，进行有限通用芯片体系结构设计的趋势日渐明显。体系结构研究人员通过分析传统高性能微处理器的特点，使新的体系结构能够达到更高的面积与能量利用率。

面积利用率高是指用相同的面积集成更多的功能单元以实现更高的性能，能量利用率高则意味着用相等的能量能够支撑更多的计算。高的面积利用率使得成本低，工艺上更方便实现；高的能量利用率使得执行一个相同任务，耗费的能量较少，这两个指标比用时钟频率衡量的性能更有意义。表 1.2^[5]列出了各种专用处理器与可编程处理器的能量利用率，所有处理器都使用 0.13 μm 制造工艺、1.2 V 工作电压。能量利用率以每个算术运算所消耗的能量为单位，并假设处理器处在最高性能且功耗最大的情况。虽然大多数处理器往往达不到最高性能，但它们的实际性能数据不易稳定地测量，所以表中只给出峰值性能。

表 1.2 传统高性能微处理器的性能效率 (0.13 μm 工艺, 1.2 V 工作电压)

微处理器	数据类型	峰值性能	功耗	能耗/操作
NVIDIA GeForce3	8~16 b	1200 GOPS	6.7 W	5.5 pJ
MPEG4 Decode	8~16 b	2 GOPS	6.2 mW	3.2 pJ
Intel Pentium 4 (3.08 GHz)	FP 16 b	12 GFLOPS 24 GOPS	51.2 W 51.2 W	4266 pJ 2133 pJ
SB-1250 (800 MHz)	FP 64 b 16 b	12.8 GFLOPS 6.4 GOPS 12.8 GOPS	8.7 W 8.7 W 8.7 W	677 pJ 1354 pJ 677 pJ
TI C67x (225 MHz)	FP	1.35 GFLOPS	1.2 W	889 pJ
TI C64x (600 MHz)	16 b	4.8 GOPS	720 mW	150 pJ

1.3.1 专用处理器

专用处理器也称为固定功能的处理器，它直接把一个应用的数据流图映射到硬件上，符合流应用的特征。专用处理器一般包含大量并行执行的运算单元，具有与流应用的计算密集性、并行性相适应的硬件组织与结构。这些运算单元通过专用线和存储器进行互连，体现出应用应有的局域性，使得面积利用率与能量利用率非常高，减少了片内信号传输的距离和对多端口大容量全局存储器的访存次数，使得大部分的基片面积与能量分配给运算单元，而不是控制与通信逻辑，可以获得相当高的面积利用率与能量利用率。

多边形着色芯片 NVIDIA GeForce3 和 MPEG4 Decode 这两种专用图形图像处理器的能量利用率可见表 1.2，每个操作的能量消耗还不到 6 pJ。表 1.2 中列出的其他的处理器都是可编程的处理器，可以看出，在可编程处理器与专用处理器之间，效率上存在着巨大的鸿沟。

1.3.2 通用微处理器

表 1.2 中还包括两个通用微处理器，一个是 3.08 GHz 的 Intel Pentium 4，另一个是 SiByte 的 SB-1250^[9]。

Pentium 4 微处理器为了获得高性能，是以牺牲能量利用率为代价，采用了超过 20 段的流水线、高频率的时钟、大深度分支预测和前瞻（speculative）硬件等。从表 1.2 中可以看出，Pentium 4 微处理器的能耗/操作是表中所有处理器中最高的。

SB-1250 微处理器通过使用低功耗技术，限制了流水线段数，提高了能量利用率，与其他低功耗处理器相比（如 XScale^[10]）不相上下，这些微处理器在 0.13 μm 的工艺下，平均每条指令耗费 500 pJ 左右，能量利用率高于 Pentium 4 微处理器。

通用处理器能量利用率较低的原因是采用了全局寄存器文件、Cache 甚至深度