

# 作者简介

冬瓜哥，现任某半导体公司高级资深架构师

《大话存储》系列图书作者。大话存储 公众号作者

4项专利唯一发明人



大话计算机



大话存储

(12) <b>United States Patent</b> <b>Zhang</b>		(10) <b>Patent No.:</b> <b>US 9,257,144 B1</b>
		(45) <b>Date of Patent:</b> <b>Feb. 9, 2016</b>
(54) <b>SHINGLED MAGNETIC RECORD HARD DISK DRIVE AND METHOD FOR CREATING A LOGICAL DISK FROM PHYSICAL TRACKS</b>		USPC ..... 360/79, 78.14, 78.07, 77.06, 77.07, 360/77.08, 74.05, 55, 31, 72.1, 72.3 See application file for complete search history.
(71) Applicant: <b>PMC-SIERRA US, INC.</b> , Sunnyvale, CA (US)		(56) <b>References Cited</b>
(72) Inventor: <b>Dong Zhang</b> , Beijing (CN)		U.S. PATENT DOCUMENTS
(73) Assignee: <b>PMC-Sierra US, Inc.</b> , Sunnyvale, CA (US)		8,786,969 B2 * 7/2014 Kasiraj et al. .... G11B 5/746 360/39
(*) Notice: Subject to any disclaimer, the term of this patent is extended or adjusted under 35 U.S.C. 154(b) by 0 days.		2011/0075292 A1 3/2011 New et al.
(21) Appl. No.: <b>14/688,696</b>		FOREIGN PATENT DOCUMENTS
		EP 1521261 4/2005
		* cited by examiner
		Primary Examiner — Nabil Hindi
		(74) Attorney, Agent, or Firm — Dennis R. Haszko



¥94.10

大话存储2：存储系统架构与底层原理极

1300+条评价

清华大学出版社



¥141.60

大话存储：存储系统底层架构原理极限剖

7000+条评价

清华大学出版社



¥84.60

大话存储后传：次世代数据存储思维与技

7000+条评价

清华大学出版社

# 新书简介

**新书书名：**大话计算机——计算机系统底层架构原理极限剖析

**页数：**12章，1500页，图片相当多，全彩印刷

**定位：**计算机专业科普书、百科全书。从零到最顶端全技术覆盖深度阐述。IT领域从业者全覆盖。

**特点：**从零开始，覆盖全面，而且通俗，接地气，代入感强烈，引发读者思考，然后解决问题，将思考过程和解决办法描述下来。符合人类原生态认知、求知思路。

**出版日期：**2019年4月

## 本书特点

- 01 绝对从初学者角度出发，看了**不迷茫、不撕书、不会骂作者**（我也不想被人骂）。
- 02 介绍事物绝对庖丁解牛，轻易不留“坑”，不得已留了则必填。
- 03 带着思考来写作，**促发读者思考**；问题导向，带着解决问题过程来写。
- 04 事物之间带有**前因后果关联**，而不是孤立地介绍，整本书从第一页到最后一页有一条清晰的**因果脉络**。
- 05 全局框架和局部细节兼顾，**大而全，深而细**，就像一部精美的游戏，宏观场景震撼，局部细节惊艳！
- 06 覆盖面极广，涵盖多个领域关键知识：**数字电路、模拟电路、计算机体系结构、计算机组成原理、操作系统原理、计算机图形学、高性能计算机集群/超级计算机、信号与系统、存储系统、网络/通信系统、机器学习与人工智能**等。
- 07 随便点选任何一页的任何一段，都是**精彩和拥有丰富细节的内容**。
- 08 注重大框架的建立，让读者阅后**成竹在胸**。
- 09 **便于自学**，看文字就像是在聆听作者当面讲授。
- 10 图片细节丰富，带有**事物运行的流程**，其过程是动态的而不是静态的。

## 业界大牛点评

打开来读，他的文字带我走过一个个章节，一直翻到了最后一章，感觉就像自己温习了一遍大学本科的课程，还有从业二十年来看过的许多专业书籍、科研文献、设计文档。纵观市面上所有计算机图书，能在一本书中说清楚计算机工程(COMPUTER ENGINEERING)的方方面面的关键要点的，目前只有冬瓜哥的这本《大话计算机》了。

**廖恒**

海思半导体 首席科学家、FELLOW



# 读者反响强烈



Aaron

@冬瓜哥  
刚大概看了您发的pdf,想必您一定花了不少心血,书中的每一幅图都简洁明了,每一个问题的解释通俗易懂又恰到好处,几乎做到了极致,堪称完美。



【高手】缘来&似梦

冬瓜哥的图确实画的好



【入门】同学甲

冬瓜哥的图确实画的好



夕阳如斯

小白不敢点菜 冬瓜上啥都行!

2018-04-23 17:49:06

冬瓜哥这本百科全书我是服气的,对于不懂计算机的人能提供强大的工具书作用



【高手】Bottle cap

你这ppt简直登峰造极



【高手】Bottle cap

感觉我对不起CAD这款绘图工具



IQYA

我之前看书的目录,发现都看不懂,汗颜啊



深山野老



QYA

不过看了刚才发的具体内容,发现很容易明白。



卧槽 牛啊



这个搞的比计算机博士生还专业啊



彩色的牛

不要,我要看完瓜哥的《大话计算机》就算是慰藉我的计算机生涯了



初学者

不过分



彩色的牛

不过分



TIM

不过分

@Nervermore 内核部分也是迄今为止

【大师】AMP

冬瓜哥这ppt技术我是服气的



Anonymous

我看到了未来的dalao



教官\*

等待你的书。快点出来



教官\*

@冬瓜哥



Anonymous

我也是,这本书我也等了好久@冬瓜哥



【专家】Captain Sparrow

体量好庞大



【管理员】不是摄影

确实漂亮 是个好ppt教程



【专家】聂永

确实牛逼

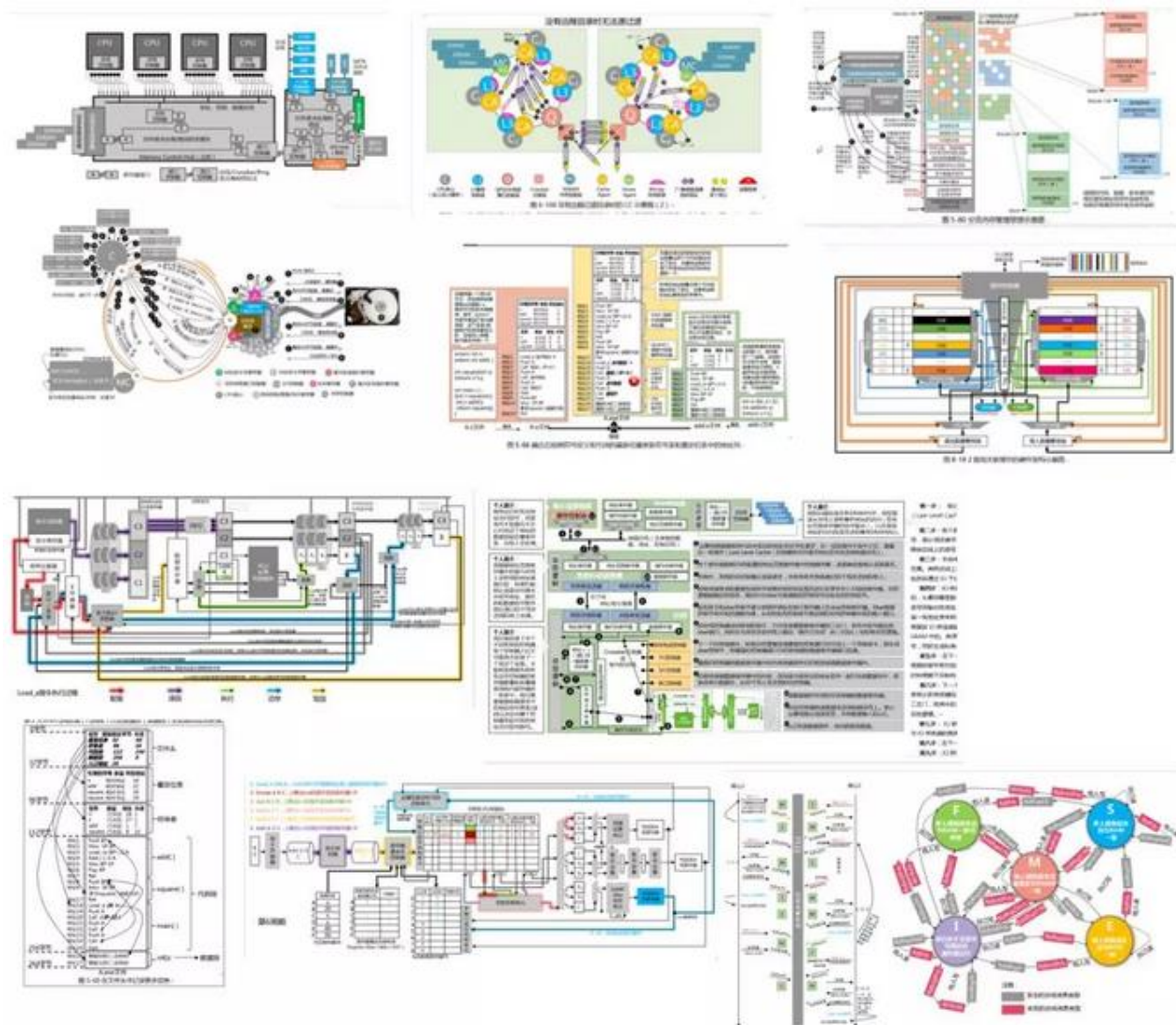


【专家】聂永

这本书



为何本书会顺手被网友称为PPT教程？（更多图片下面还有）



记得多年前有这么一道研究生入学面试题：“播放幻灯片时，按下一个空格键到屏幕显示下一页，请问这个过程计算机做了什么？”大多数参加面试的学生都答不上来。但是有一天看到冬瓜哥的《大话计算机》时，我第一反应是意识到自己错了。这本将近1400页的恢宏巨作，约500个章节，涵盖了处理器流水线、缓存、内存、并行计算、网络、声卡、GPU、操作系统，甚至包括半导体制造工艺等，每一章节都是深入浅出。不同于一般的教材，这本书采用诙谐幽默的笔法与图文并茂的形式向读者揭示计算机内部各个部件的工作原理，并穿插着技术背后的种种名人轶事，读来生动活泼、引人入胜。这像是一部小说，更像是一本百科全书，按图索骥，总能找到你想了解的知识点。

**包云岗**

中科院计算所研究员，  
先进计算机系统研究中心主任，中国科学院大学岗位教授

当年我看游戏《半条命：反恐精英》的二三百万行源代码时，有同事说，网传看完就只剩下半条命了。当我看到冬瓜哥这本《大话计算机》时，瞬间感觉我看完这本书可能也只剩下半条命了。并不是因为这本书的难度摧残大脑，相反，这完全是另一种感觉，它太通俗了，把事物的关系、流程、架构讲的太清晰了，信息量太大了，看完之后就是身体被掏空感。因为，我发现自己这么多年习得的仅有的点点计算机方面的绝招、秘密在这本书里竟然一点不落的都有，而且还通俗易懂，让我瞬间感觉之前自己耗费在学习计算机上的时间，简直成了浪费人生。

**张勇**

深圳市科力锐科技有限公司创始人之一





期待了很久很久的冬瓜哥巨著终于到手😁，这部拆分为三厚本印刷出版的巨著，图文数量和专业度质量，都直追高纳德圣经，堪称镇宅典藏珍品，终生研读！



终于收到瓜哥的大作，瓜哥付出4年的时间专注于一本书，其对技术的热爱，对专业的专注和坚持，让人敬佩！



凝结智慧与辛劳的巨著，天下贤士共赏之，莘莘学子向往之，必不脛而走，洛阳纸贵😁



【瓜哥】

扎实的分量，精致的印刷，这绝对赚了



瓜哥出品，必属精品





收到冬瓜哥呕心沥血大作，看到云岗和万青的作序，会心一笑，非常感动，这年头专心致志写一本书，花几年时间，在鎏金镀银随时随地赚大钱的计算机行业，独处书斋，屏蔽万事，放弃了无数赚钱的机会，真是可敬可佩！

平生不看大话计算机，阅尽github所有源码也枉然 😄

嗯 超值

冬瓜哥作图技能点已满

已推荐同事购买

我大致翻了下，很舒服

那些复杂如蜘蛛网的示意图和架构图可能刚开始的时候确实看不懂，但是多看一段时间之后，就会从线条和文字中隐隐约约看出四个字：瓜哥牛逼

廖田靖

是舒服

我看这本书的图画得很牛逼

@冬瓜哥 这套书，你确实费了心思 👍

下午 16:39:48

区别其他书，图和内容都是老哥自己搞明白后，自己画的

拿到书了 🙌 😄 /顶呱呱/顶呱呱

这书写的特别细 👍

@永无止境 这个盒子回收值钱

下午 12:24

不舍得回收



弱智天君

PLUS会员



一年多前就在书的序言人之一的朋友圈里看到对此书和作者的大力称赞。了解到作者坚持了四年，用心为各个层次的读者，也为中国计算机界带来了一份原创巨作，这份努力实在值得敬佩，国内各行各业都太需要这样精雕细琢的工匠精神。看到预售消息之后，果断买了首发版表示支持。书拿到手快速翻了一下，从印刷质量到排版方式，诚意十足，丝毫不掺水分。先翻看了前几页内容，暂不做过多吹捧，内容绝对够细致，能感受作者融会贯通之后的娓娓道来。第3页正文看到了个校对时未发现的typo，相信作者一定是个追求完美的人，略觉可惜。相信该书一定能成为一个传奇，向坚持奋斗不懈努力的人致敬。

举报



1



0

真的震惊，这三本书的信息量



已经收到的我表示沉迷东哥的书无法自拔

2019-04-28 11:40:24

光看第一章关于东西方人对符号的理解的思维方式就已经让我体验一次恍然大悟了

信息量太大， 我表示没个 几个月，吸收不了。

不过书的质量确实可以。 内容丰富

书很赞,已经推荐给学弟学妹啦

正在翻看，图片，图表 印刷很好，挺值的

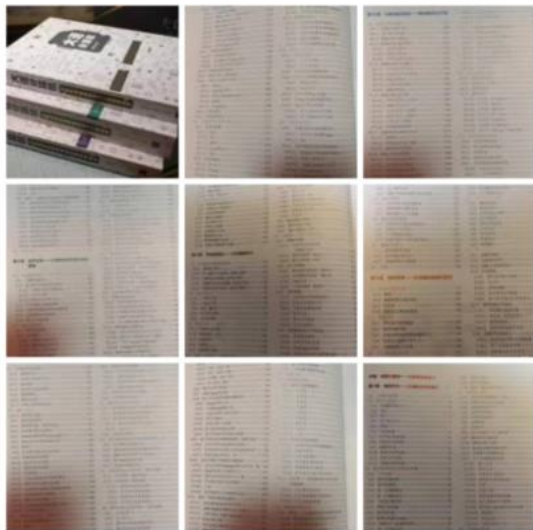
梳理知识体系 挺不错的

“2019.4.26.星期五 下午 23:08:41 查看原文  
没有一点it知识的人能看懂这个书吗

没问题的

昨晚睡前一口气看了第一本的70%

晚饭后浏览了这厚厚的大话计算机，此书堪称计算机的百科全书，通俗易懂，从初学者角度出发，看了不迷茫，不撕书，不骂作者。全书只有冬瓜哥一个作者，不禁佩服是什么驱使冬瓜哥坚持完成了此著作？唯有情怀！向东哥致敬，东哥也给自己赞一下吧😊



1小时前

大作已收到，真是学习计算机系统难得一见的好书。如此巨著以一人之力完成，深感钦佩！一定在师生中大力推荐👍



j\*\*\*y

PLUS会员



计算机书籍的经典之作啊，真的太棒了！

大话计算机 2019-04-28 13:37



k\*\*\*i



很不错，深得我心

专业必备 正版 脉络清晰 理论基础 性价比高

大话计算机 2019-04-28 13:10



2\*\*\*2



期待已久的神书，心心念念了好久。精美的插图，详尽的知识。果然冬瓜哥出品必原有的一堆书 仔细阅读它就可以了 希望自己有大的提升！

另：获得了为数不多的彩蛋版 冬瓜哥亲笔签名，激动一下

大话计算机 2019-04-28 12:24



章节导图导出图像时选的分辨率不够，导致文字略微有点糊



我以前搞过高清喷绘、数码快印🤔



书中的图片非常精美！



我要是在高校任教的话一定要选为必修教材



中国移动

下午5:21

43%

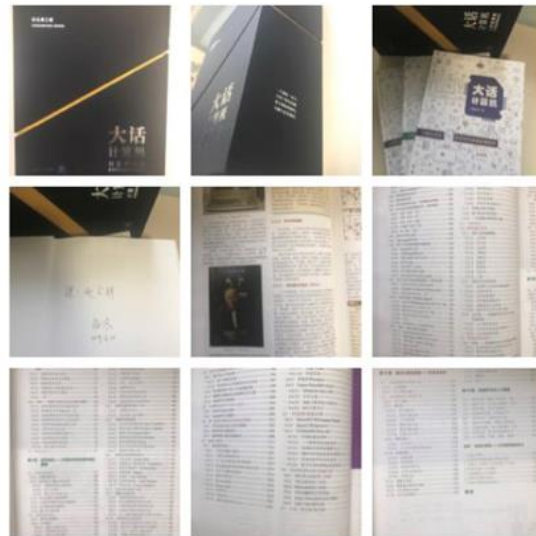


朋友圈



感谢冬瓜哥皓首穷经完成的鸿篇巨著《大话计算机》，全彩色印刷，图文并茂，堪称科技届的《资治通鉴》。覆盖了硬科技的方方面面，既有前沿的人工智能，又有厚重的科技历史。古有司马光，今有冬瓜哥。大人看完值得再传给孩子读，让他们发现科技的世界竟然这么宏大和有趣👍

收起





16:57

... 4G 54

< 详情

15:30

回复 冬瓜哥: 来借书啊 😊

15:31

回复 冬瓜哥: 有这想法, 😊

16:46

回复 冬瓜哥: 刚上市的 🤔

16:50

回复 冬瓜哥: 😂 网上评价蛮高的, 我也买了这个书!

16:54

回复 冬瓜哥: 冬瓜哥历时四年的憋尿之作, 的确值得买来学习

16:56

回复 冬瓜哥: 书真还不错, 就是有点贵

16:56

回复 冬瓜哥: 我买了两套, 一套给公司, 一套自己看

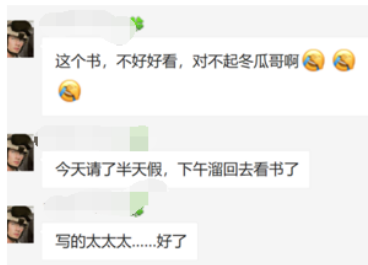
16:57

回复 冬瓜哥: 👍

评论



发送



@冬瓜哥 瓜哥, 杭州的书到了! 可惜我出差, 我妈代收的。我妈问我什么东西这么重, 我说 是知识的重量 🤔

< 朋友圈



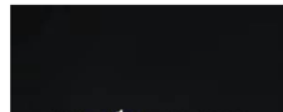
葵花宝典已到手, 从来没见过字这么! 的大部头书, 太实诚了。



4分钟前



冬瓜哥: 要是常规5号字, 就得2000页, 价格就会更高



✓ 只看当前商品

好评度100%

全部(70+)

最新

好评(60+)

中评(0)

差评(0)

视频晒单(1)

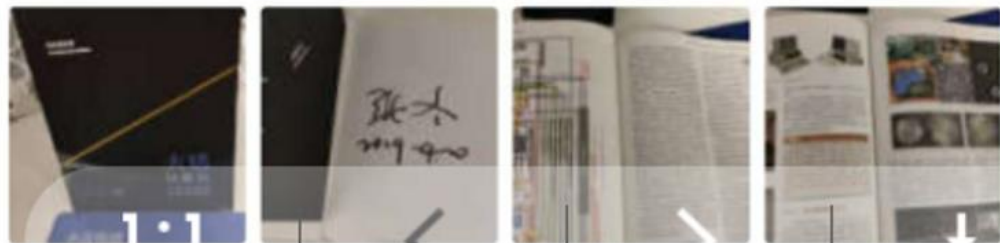
有图(9)



怎\*\*\*买 ★★★★★

2019-04-29

作为一个从事计算机it行业20年的，读书人。强烈推荐！如果你想快速跨入这个行业，本书必读！我已经把这套书作为，公司培训教材了。公司新人必读！这不是教材！人文看浪潮之巅，技术看《大话计算机》！通读过后，你的能力，会跃然提升！可以短时间积累足够的行业经验！菜鸟变专家，这套书是捷径！书贵吗？对比各种培训，简直是超级超级超级的有性价比！^



颜色：大话计算机

回复



@冬瓜哥 大话计算机的印刷真的不错，读起来特别引人入胜，随时可以沉浸进入👍👍



「冬瓜哥：@老鹰 哎，误人子弟了。。🤔」

-----  
我记得你当时和我讲 - 这本书的写作，就是一个边学习边总结的过程，其实这种写法，挺适合学习者的，特别不是科班的学习者 - 咱们自己的科班教育模式，对学习计算机其实并不好，这也是这本书的价值所在，我读的时候，这种感觉。

这两天晚上粗略翻了一下这套书，累计一个多小时吧。

人已四十，咱既不是粉，也不是吹。只是有些想法，想留一下。

首先，这套书写的真是太系统太细致了。让人不自觉的想沉浸进去，享受那种从迷失到恍然大悟的快感，似不似有点类似于吸毒？书到之前，本以为这套书应该是有联合作者的，因为这套书关联了三个领域。这是我第一次看到有这么一套从逻辑电路一点点展开自下而上延伸到计算机结构再延伸到操作系统内核的经典书籍。

那么我会怎么看这套书？首先，由于所从事的领域不在电子方面，所以前四章我估计暂时不看的。我将从第五章开始，正式起步从第五章后半段开始。我会怎么用这套书？其实迷失在一套书里，那种欲仙欲死的着魔状态非常的痛苦。为了避免，我将会打开这套书，同时也打开电脑中积累的老资料。沿着这套书的思维导图，遍历每个小节。然后让我脑中的认知、老的资料、大话计算机中的论述，让这三部分的内容在大脑中碰撞，或者是相互融合，或者是互相纠正。总之无论是互补互通还是破而后立，都会让自己有质的提高。

每拿到一本书，我总喜欢先琢磨目录结构。最后，说一下从目录结构上看，我的一个粗浅的认识，一点遗憾：这本书专门讲了Linux32位内存/地址结构，但似乎没有讲64位内存/地址空间的内容。相关的64位内容，网上的资料太杂，我一直在找到权威的书籍去澄清脑中认知。看来还要继续找下去。

与描述相符

5



1 2 3  
非常不满 不满意 一般

全部 追评 (0) 图片 (0)

好书，建议店家多大广告一下宣传这个书，作者的书是目前看到最佳的一本没有之二！



这本书真的是niubility

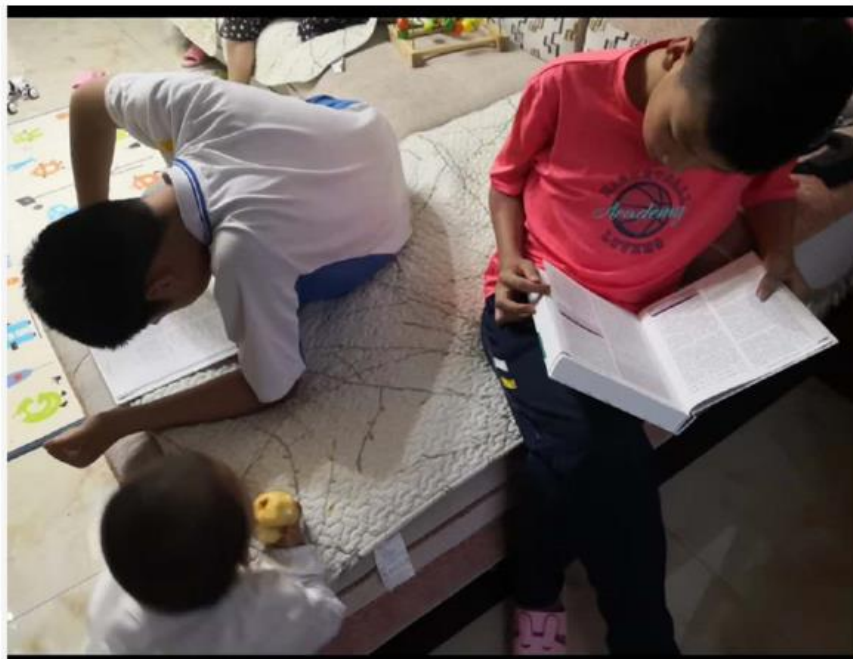
啥书

大话计算机





不过严格的来说这是一本精美的ppt展示书籍



< 详情



让两个娃看点计算机底层的书

按冬瓜哥的说法，建议16岁以上看，但我这边这几个娃娃都被我开光了，居然能够边看边哦，很多问题他们居然可以边看边找到答案，不错。



提到了我

19:15



假装你喜欢我

★★★★★

我敢保证，这确实是我看过国内外最好的同类图书，没有之一。像作者致敬，感谢他分享自己的知识，为后辈写出如此好的著作。

大话计算机

2019-05-02 13:15

举报

👍 0

💬

老\*\*\*头

PLUS会员

★★★★★

首先对这样一部倾注心血而写成的“奇书”的作者表示感谢，曾几何时自己在看到不同教材，辅导书的知识讲解时，对术解，都需要查阅更多其他资料，而缺乏系统性的整体了解，而断续的知识获取又往往在问题解决前出现新的问题，因此这本书的最终原因，希望与我有同样感受的人都能喜欢这本书。

大话计算机

2019-05-02 08:40

举报



张工 🍎🍎🍎

五一这么过 😊😊



旷世神作，此书已是顶峰，后无人可超越



共三本 还在看第一本



循序渐进 很棒



第三章 挺有意思 不愧是化学毕业 写的疏而不漏



写的太棒了



五一几天 啃完前四章 看的不想放下来



不过还是烧脑的

< 详情



支持冬瓜哥，本科化学，没有读过研究生，纯靠兴趣，对计算机底层的理解和技术的追求已经到了一种洁癖，也是目前计算机著作的一股清流👍👍🤔🤔



刺猬四年磨一剑 -- 冬瓜哥《大话...

昨天 12:58



13:01

冬瓜哥认识大约十年前，博导帮他第一本书写序，当时博士的时候博导经常拿他来鞭策我们说“你们这些计算机科班的难道还不如学化学的？一定要深入基层，理解原理”，好久没他声音了，这次四年磨一剑，终于出了重磅著作，强烈推荐😊



买了一套



内容真心翔实👍



深入浅出



四年辛苦不寻常，字字读来都是精



刚刚大致翻了一下，图文挺多的。感觉深入到计算机的底层去思考了



有种打破砂锅，深究到底的感觉。



今天已经收到书了，回家迫不及待打开箱子，一个词形容“震撼”，震撼到难以形容，我对作者的才华和决心佩服至极！如果有机会我想请教作者，是什么支撑着您学习完这么多东西然后再消化、整理再写出这本巨作？我真的很难相信这本书会是一个人写出来的，不在于字数而在于知识面之宽广和深度。佩服至极！

2019-05-03 23:15:19





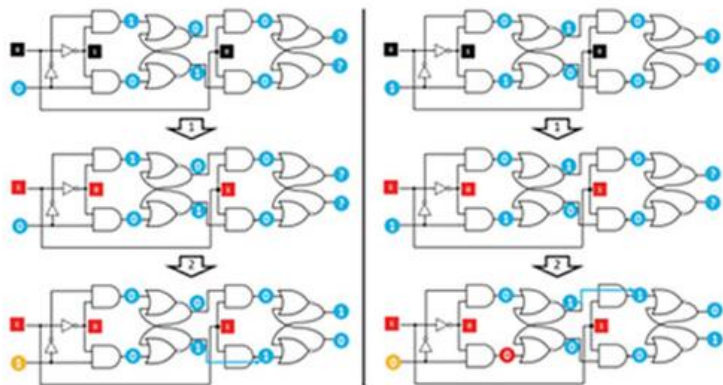


图1-27 D端为0和1时的锁存过程示意图

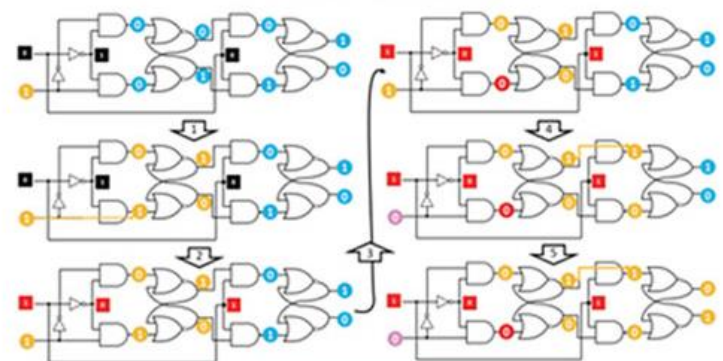


图1-28 连续输入两个新数据时依次锁定过程示意图

橙色信号传递到输出端。与此同时，紫色的新数据到达了D端，等待下一个时钟下沿将其传递到中间，然后再来一个时钟上沿将其传递到输出端。

触发器就相当于照相机镜头（输入端）、快门（触发器电路中的与门）和底片（输出端）。时钟信号控制快门开合，时钟边沿的一瞬间，快门打开，镜头投射到快门上的影像会被瞬间透光锁定到底片上感光保存，在下一个时钟下沿之前，输出端信号不受输入端信号影响。

边沿触发器更像是两道依次开合的闸门，两道闸门之间是一个临时等待区（图1-26中的A、B点），如图1-29所示，时钟下边沿会将输入端信号传递到触发器内部，但并不输出；时钟上边沿则将上一步暂存

在内部的信号输出，并在瞬间封闭输入端通路，直到下一个下边沿到来时，被堵在门口的新输入值才会进门在临时等待区等待。

综上所述，边沿型触发器电路仅当锁定信号从0变成1的瞬间，才能够锁住当前D的值并在Q端输出，后续不管是D变化，还是锁定信号从1变成0，Q值保持不变，所以被称为边沿型触发器。而之前的D触发器，称之为电平型触发器，因为其在0电平时锁闭，1电平时打开。边沿触发，这一点已经比较奇特了，但是这位神人还没有就此罢休，他调皮地将Q'端和D端连接了起来，并分析了电路的时序逻辑，发现了更奇特的结果，有时候真理只差一步，不要断然否认任何尝试，不要觉得某件事肯定没意义，发现真理的人，

个固定的位置先让程序读出来，然后程序就可以一层一层地按照指针指向的路径继续读数据。

图5-26中所示的“超级块/根入口/根指针”指的就是这个固定入口位置，超级块可以是一个扇区，或者多个连续的扇区组合起来（称为块），其中存放着一个或者多个根指针（也可以放最终内容，但是这里空间太小，放进去反而不划算，一般只放一些简要的信息比如系统版本号、最后访问时间等）。超级块内部的根指针指向另外的扇区或者块，这些被指向的扇区/块可以存放最终的数据内容，也可以继续存放指针，如果选择后者，则这些

位于根指针下一层的存放指针的扇区/块被称为**二级指针块**，指针块内也可以既包含实际文件内容又包含指向下游数据块的指针。二级指针可以继续指向三级指针块或者混存块。同理，也可以有更多级的指针块，最后一级指针块会指向纯用于存放最终数据内容的块，公示表里那些文件名、起始扇区号之类的信息，最终会被存放于此。

按照上述的思想，我们可以将文件系统改造为如图5-27所示的这番景象，这样是不是更壮观了？其实这只是最简易的设计，现实中的文件系统设计远比这复杂，当然，功能、可靠性、性能也都更强。

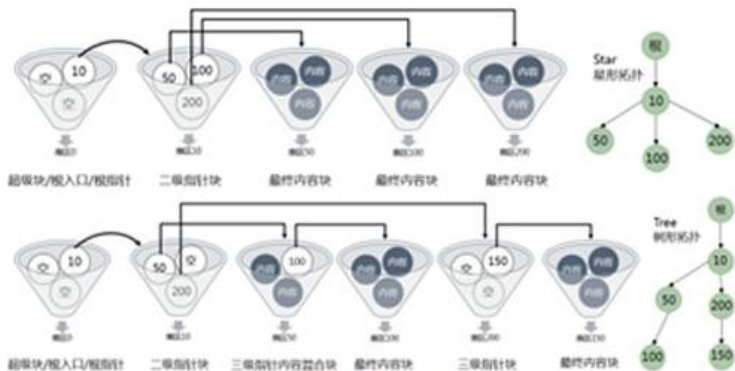


图5-26 星型和树形结构

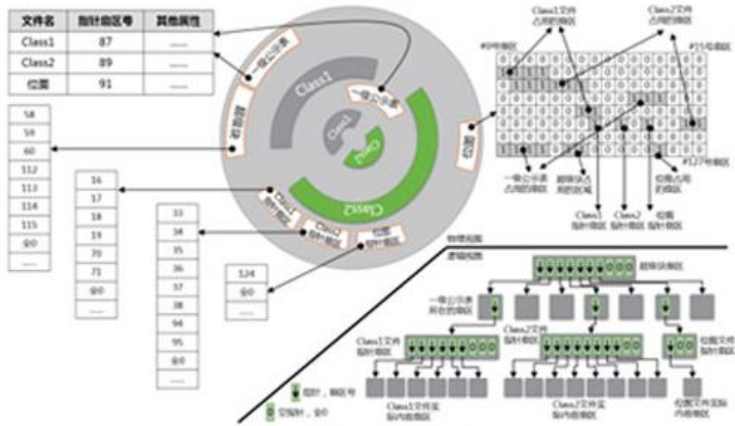


图5-27 一种改良版的简易文件系统设计



为a.out的可执行文件，/a.out一样可以执行它。a.out的意思是Assembler Output（编译输出）的意思。）

值得注意的是,假设1.0调用了2.0中的函数A,而A里面又调用了位于3.0中的函数B,B则没有再次调用外部函数。那么此时你必须将1/2/3.0这三个文件都给出,链接才能完成。这在现实中是非常普遍的现象。所以,我们一直没有回答的一个问题:为何需要多个.o文件并打包成库文件呢?此时你应该知道答案了。如果能够将有联系或者相互引用关系的多个库文件打包成一个文件,那么链接的时候只需要给出这个文件名即可,而不需要把所有相关的目标文件名都列出来。

然而从图5-67中可以看到，main()函数的第一句代码Load a 26 A是该程序的第一句代码，程序必须从这里开始执行。然而，该句代码却被放到了地址18上。我们前文中一直在讲这样放，把程序放到存储器的第一行上，第0行放一个NOP指令。我们设计的这个CPU加电之后，用NOP指令避免了初始状态的不确定性，然后从第1行开始执行代码。那么说，这个程序就无法运行了么？想一想都有哪些方法可以让这个程序继续运行了？

(1) 真是闲的没事了，为什么不把main()放到程序的第一行呢？把main()的位置定死，必须从第一行执行，可以这样做。



图 4-26



从而输出对应的电压,这就是动态功耗。如果不读取它,开关关闭,对应的baseline就不会被充放电,只维持原来的电压,并存在一定的微小漏电流,这就是静态功耗。

于是,有些设计采用了冬瓜哥又花了2个小时雕琢而成的图6-20所示的方式,将Tag部分拿出来放到







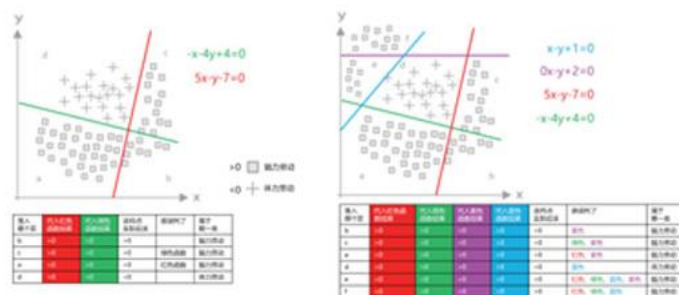


图12-7 两个分类场景下的结果总结

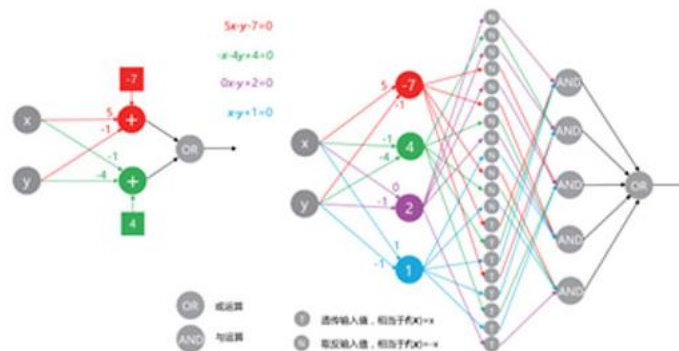


图12-8 对应图12-7中场景的逻辑表达式图

除了把事情表达成与或非的关系之外，还可以从另外的角度和方法来入手，从而可以实现更多效果和结论。

如图12-9所示，还是刚才的场景，c区被绿色函数误判，但是对于蓝色函数，c区的判定全部正确。如果以某个样点距离函数直线的绝对距离值作为该样点的分类可能性的话，那么越远离直线越可能正确或不正确。对于c区中的蓝色样点，蓝色函数说该样点是脑力劳动者，可能性假设为1，但是绿色函数说该样点是脑力劳动者的可能性为-5。到底该听谁的？如果两者都听得话，那就将两者输出的可能性值相加：1+(-5)=-4，结果显示该点不是脑力劳动者，但是开启上帝视角之后可以看到蓝色样点的确为脑力劳动者，应该听蓝色函数的，但是无奈蓝色函数判定该样点是脑力劳动者的可能性低得可怜，只有1，压不

过-5，那怎么办？

好办，我们给蓝色函数一个较高的权重，也就是直接将它乘以10，再与绿色函数相加，如图12-9中间所示。这下好，10+(-5)=5，脑力劳动者！让你明显！蓝色函数有了特权，就可以在c区无视一切其他可能性值的贡献度。对于橙色样点，其本来的可能性就已经很高了，已经压过绿色的反对声音了，被放大10倍之后，就更高了，但是回过头来看看a区，蓝色曲线本来就完全误判了a区中所有样点，认为其为脑力劳动者的可能性都<0，现在被放大了10倍，相当于把错误也放大了10倍，而绿色函数的正确声音在a区被完全压制了，a区彻底沦陷！

那好，既然蓝色函数不给力，把绿色函数放大10倍看看，同样的事情发生了，a区收复失地，但是c区

## 四大理由值得拥有



### 很专业但能看懂

作者基于多年半导体从业的思考与沉淀，以独特视角切入缤纷计算机世界，用接地气的语言、华丽精准的图表纵情解读。



### 卓越阅读体验

全彩印刷+温婉柔和的纸张+8色顶配印刷机+平铺桌面不回弹的精装书。

## 精细豪华解读

三册精装，超过1500页，超过350万字，超过500张原创图表，11张8开思维导图海报，大量二维码引入图片、音视频.....堪称计算机原理饕餮盛宴。




## 再造IT基本功

深入研磨计算机系统，你的软硬件开发思维、AI路线、算法思路，甚至信息安全、芯片制造等领域，一定会出现让你自己大吃一惊的感悟和火花，你的收获将远远超过图书定价。









冬瓜哥其人，喜钻研，擅用与科班教材截然不同的方式把计算机原理细细分解，娓娓道来。

这本书如此通俗，以至于假设当代文明毁灭之后，下一个文明从地壳中找到这本书，按照书中所讲即可迅速建立起计算机体系，极大推动下一个文明的进程……

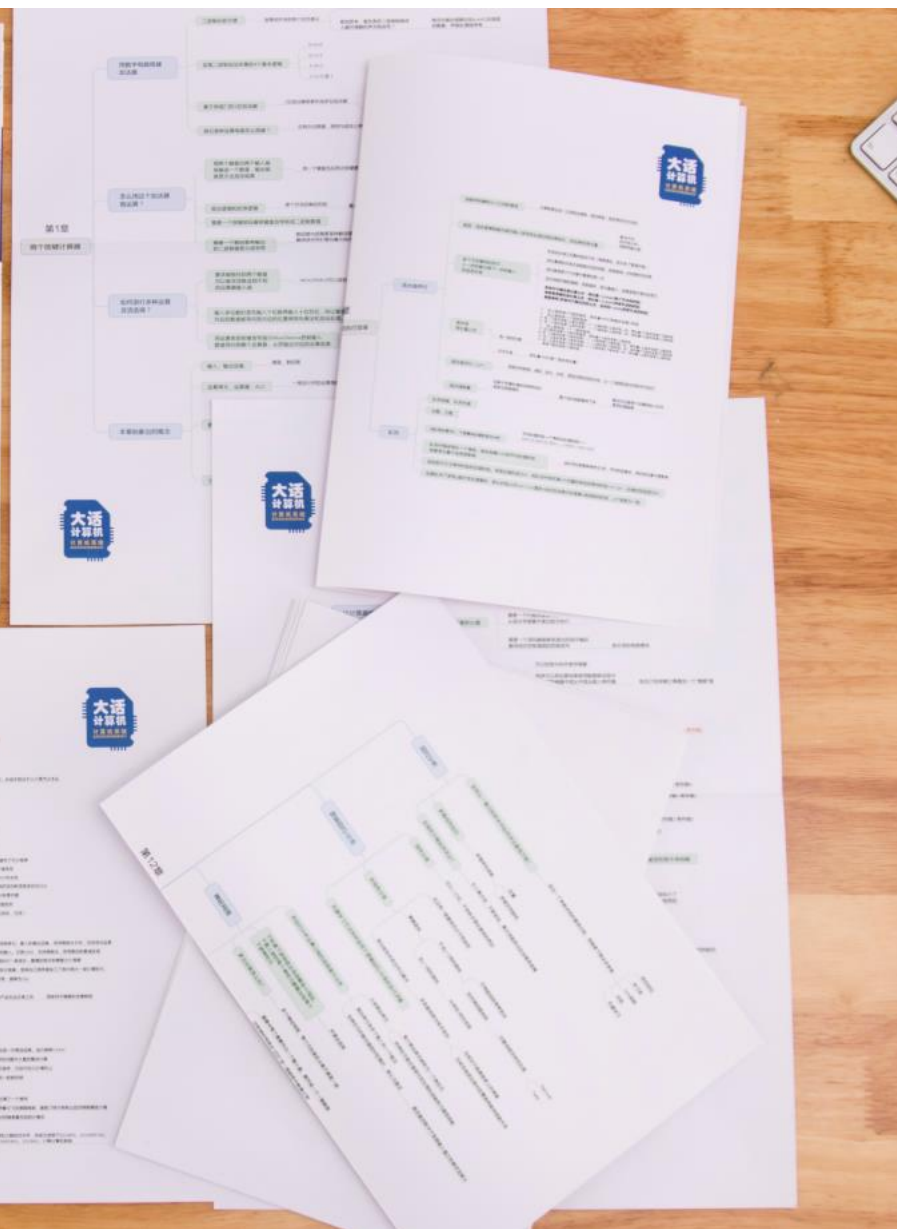












转交。Eager模式则是节点在收到消息之后就立即转发出去。同时节点会异步查询。直到网络上所有节点都收到该消息。一次之后回到读未节点。这种方式浪费了很多带宽。如果提升带宽速度。与广播无异。这两种方式都是没有过滤目录的。也就是说每个节点消息必须所有人都知道。免不了。除非某个节点判断出最新数据就在自己。没必要往下转发了(比如Lazy模式)。第三种方法称为“Oracle”(先知)。之所以未知。当然是因为每个节点都实现了精确的过滤目录。可以点对点传递消息。

6.9.10.3 增设远程目录过滤片外广播

上文中介绍了单CPU芯片内部Inclusive的LLC缓存行或者Exclusive的所有缓存行中存放过滤bitmap目录。以实现片内广播过滤。该bitmap又被称为目录。由于其只能追踪本CPU片内的缓存状况。所以暂且称之为本地目录。同时也提到了。由于不知道其他CPU中是否缓存了该行。所以一些作成的消息不得不向其他所有核心广播。如果能够存在某种专门用于追踪其他CPU片内可能缓存有哪些缓存行的bitmap的话。这个bitmap就可暂且称之为远程目录。那样就可以避免广播而使用单播。节省网络流量。提升性能。下面我们使用一些流程图来帮助大家梳理清楚本地目录和远程目录在访问操作时的过滤作用全景。本系统采用Inclusive模式的缓存设计。直接在LLC缓存行中存储本地目录bitmap。

图6-167为一个在没有远程过滤目录情况下的CC案例示意图。C<sub>0</sub>核心发出一个Stor请求。C<sub>0</sub>的CA通过查询本地目录成功地发现C<sub>3</sub>核心的私有缓存中不会命中。那就不需要向其发送WrtIvid消息了。但是C<sub>0</sub>

的CA无法判断右侧的CPU内都有哪些缓存行。所以发送了个目标地址为广播地址的WrtIvid消息到该CPU控制器。右侧CPU的OP控制收到该广播。便会产生该向本地目录(假设该访问的地址会被分配到C<sub>0</sub>核心的L3缓存分片)以判断自己的私有缓存里是否有该行。

如图6-168所示。C<sub>0</sub>和C<sub>3</sub>通过读取本地目录bitmap后发现自己对应的位为1。表明自己的私有缓存命中该行。于是去将对应行作废处理。然后返回IvidRspF消息。C<sub>0</sub>和C<sub>3</sub>通过bitmap发现自己核心里并没有缓存该地址。所以直接返回IvidRspN消息。同时。C<sub>0</sub>还从其他CPU收到了C<sub>3</sub>核心返回的IvidRspF消息。然后。右侧CPU收到了C<sub>0</sub>核心返回的IvidRspF消息。那么。该由哪个CA来更新目录呢?这是个问题。就像图中这样由这两个CA都去更新一遍。那会做无用功。降低性能。而且还会带来潜在的一致性问题。这个问题下文再来解释。

可见。上述过程发送了太多广播。我们强烈期望下。本地CPU如何知道远程CPU到底缓存了哪些行呢?一开始。所有CPU的缓存都是空的。一旦任意一个CPU发生Load操作。那么必然会发出RdPrb消息。由于发送RdPrb消息的CPU的本地目录是空的。其必将该消息广播到全网。那么本地CPU就可以收到该广播。很自然地就知道了“哪个CPU读取了这一行并缓存了”。如果是Stor指令。一旦命中缓存。也会发起RdIvidPrb消息广播到全网。这样。其他CPU也就知道“这个CPU缓存了这一行”。只要经过是够长的时

没有远程目录时无法源过滤



图6-167 没有远程过滤目录时的CC示意图(3)



转发；Eager模式则是节点在收到消息之后就立即转发出去，同时后台做异步查询，直到环路上所有节点都收到且转发一次之后回到请求节点，这种方式浪费了链路带宽，却能够提升探测速度，与广播无异。这两种方式都是没有过滤目录的，也就是说每个探测消息必须所有人都转一圈，免不了，除非某个节点判断出最新数据就在我这，没必要往下转发了（比如Lazy模式）。第三种方法称为“Oracle”（先知），之所以先知，当然是因为每个节点都实现了精确的过滤目录，可以点对点推送消息。

### 6.9.10.3 增设远程目录过滤片外广播

上文中介绍了单CPU芯片内部在Inclusive的LLC缓存行或者Exclusive的所有缓存行中存放过滤bitmap目录，以实现片内广播过滤。该bitmap又被称作目录，由于其只能追踪本CPU片内的缓存状况，所以暂且称之为本地目录。同时也提到了，由于不知道其他CPU中是否缓存了该行，所以一些作废消息不得不向片外所有核心广播。如果能够在某种专门用于追踪其他CPU片内可能缓存有哪些行的bitmap的话，这个bitmap就可暂且称之为远程目录，那样就可以避免广播而使用单播，节省网络流量，提升性能。下面我们用一些流程图来帮助大家理清清楚本地目录和远程目录在访存操作时的过滤作用全景。本系统采用Inclusive模式的缓存设计，直接在LLC缓存行中存储本地目录bitmap。

图6-167为一个在没有远程过滤目录情况下的CC案例示意图。C<sub>0</sub>核心发出一个Stor请求，C<sub>0</sub>的CA通过查询本地目录成功地发现C<sub>0</sub>核心的私有缓存中不会命中，那就不需要向其发送WrtIvid消息了。但是C<sub>0</sub>

的CA无法判断右侧的CPU内都有哪些缓存行，所以发送了个目标地址为广播地址的WrtIvid消息到QPI的控制器。右侧CPU的QPI控制器收到该广播，便会生成四条独立的消息分别发给它的4个核心，后者各自查询本地目录（假设该访问的地址会被分配到C<sub>0</sub>核心的L3缓存分片）以判断自己的私有缓存中是否有该行。

如图6-168所示，C<sub>0</sub>和C<sub>1</sub>通过读取本地目录bitmap后发现自己对应的位为1，表明自己的私有缓存命中了该行，于是去将对应行作废处理，然后返回IvidRsp消息；C<sub>2</sub>和C<sub>3</sub>通过bitmap发现自己核心里没有缓存该行地址，所以直接返回IvidRsp消息。同时，C<sub>0</sub>还从其他CPU收到了C<sub>2</sub>核心返回的IvidRsp消息，然后，右侧CPU的本地过滤目录需要被更新，因为对应的行已经被作废，不存在于该CPU内了。那么，该由哪个CA来更新目录呢？这是个问题，如果像图中这样由这两个CA都去更新一遍，那会做无用功，降低性能，而且还会带来潜在的一致性问题。这个问题下文再来解释。

可见，上述过程发送了太多广播，我们强烈需要源过滤机制。为了实现远程过滤目录，首先思考一下，本地CPU如何知道远程CPU到底缓存了哪些行了呢？一开始，所有CPU的缓存都是空的，一旦任意一个CPU发生Load操作，那么必然会发出RdPrb消息。由于发送RdPrb消息的CPU的本地目录是空的，其必定将该消息广播到全网，那么本地CPU就可以收到该广播，很自然地就知道了“哪个CPU读取了这一行并缓存了”。如果是Stor指令，一旦命中缓存，也会发起RdIvidPrb消息广播到全网，这样，其他CPU也就知道“这个CPU缓存了这一行”。只要经过足够长的时

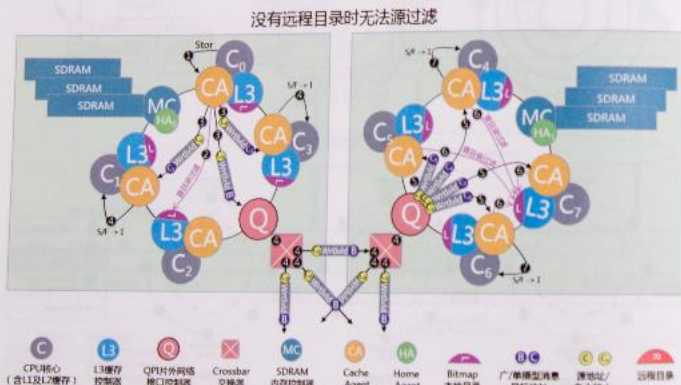


图6-167 没有远程过滤目录时的CC示意图 (1)

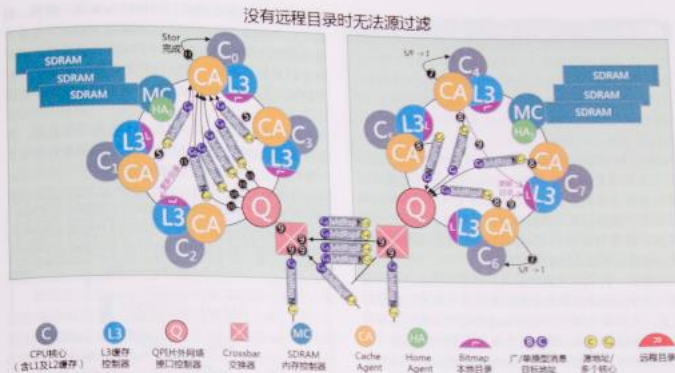


图6-168 没有远程过滤目录时的CC示意图 (2)

间，就可以积累足够多的追踪记录，从而逐渐实现过滤的目的。这个过程，称之为目录的预热过程。

那么，这个远程追踪目录需要存储什么信息呢？一个bitmap向量是必须的，bitmap中位数的数量就是系统中所有远程核心的数量（本地核心就不需要在远程目录里追踪了，有本地LLC目录）。另外，由于追踪的是远程CPU内部可能缓存的行，对应的行可能在本地并没有缓存在LLC中，所以欲将bitmap存在本地LLC对应行中就不可能了，必须单独开辟一块高速存储空间，存储对应行的地址、对应行的bitmap向量。如果要更加精细的记录，比如不仅仅要记录对应行是否存在于某个核心缓存里，还想知道对应行目前的MESIF状态，那就还必须记录一个状态字段，比如：只要收到对方的WrtIvid消息，那么本地就可以判断对方缓存中的该行一定处于M态了。

缓存行地址需要几十位，比如64位地址，缓存行大小64字节（需要6位来编址），那么缓存行地址长度就为64×6=58位，再加上3位的MESIF状态位，以及bitmap向量的16位（假设系统为4个4核心CPU），则这一条追踪记录就得至少78位大小。如果要追踪所有的缓存行，也就是2<sup>78</sup>条记录，再乘以78位，其容量将会非常大。所以，该追踪目录势必做出类似缓存和RAM之间无法全包含时的处理，比如直接相连方式，能省掉一些地址位，这样多条追踪记录会相互挤占同一个位置，还需要再次对比缓存行的物理地址Tag才能判断是否该记录有效，如果发现不匹配，则认为未过滤成功，对应的广播必须发出。

图6-169为一条记录项的示意图。思考一下，某个变量并不是任何时候都会被系统中所有核心缓存

的。假设系统共有4个核心，那么，某个变量被缓存的场景共可分为：只被1个核心缓存（共4种组合）、被2个核心缓存（共6种情况）、被3个核心缓存（共4种情况）、被全部4个核心缓存（共1种情况），总共15种状态。系统运行的时候，根据不同程序场景而定，在这15种情况中，出现比例最大的可能只有少数几种，我们假设为8种，分别为：被核心1缓存、被核心2缓存、被核心3缓存、被核心4缓存、被核心2/4缓存、被核心1/3缓存、被核心1/2/3缓存、被核心2/3/4缓存。那么，我们只需要存储3位就可以描述这8种状态。每次CA根据接收到的CC消息需要更新该条目时，电路先将当前存储的3位根据上述8种规定译码展开之后输入到一个4位逻辑电路，然后再把本次需要修改的位合入到这4位中，再通过收缩译码器将这4位译码成上述8种情况中的一种，如果本次修改的位恰好没有匹配上述8种场景中的任何一种，比如本次合入之后结果是“被核心3/4缓存”，我们称之为溢出。此时，译码器被设计为输出一个溢出信号，然后，电路可以决定，是将该条目直接作废（因为它已经无法反映真实情况了），还是动态在线地变更译码策略。后者这种做法更加智能，但是其需要该译码器是可配置的，可以根据寄存器中的控制字决定如何译码，比如在上述情况中，可以通过变更控制器寄存器内的控制位，让译码器改为可以将“被核心3/4缓存”（对应4位展开bitmap值0011）翻译成比如“010”，从而替代上述规定中的一种。下次如果又变化了，那就再替代一种。

Tag	Station	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100	101	102	103	104	105	106	107	108	109	110	111	112	113	114	115	116	117	118	119	120	121	122	123	124	125	126	127	128	129	130	131	132	133	134	135	136	137	138	139	140	141	142	143	144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159	160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175	176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191	192	193	194	195	196	197	198	199	200	201	202	203	204	205	206	207	208	209	210	211	212	213	214	215	216	217	218	219	220	221	222	223	224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239	240	241	242	243	244	245	246	247	248	249	250	251	252	253	254	255	256	257	258	259	260	261	262	263	264	265	266	267	268	269	270	271	272	273	274	275	276	277	278	279	280	281	282	283	284	285	286	287	288	289	290	291	292	293	294	295	296	297	298	299	300	301	302	303	304	305	306	307	308	309	310	311	312	313	314	315	316	317	318	319	320	321	322	323	324	325	326	327	328	329	330	331	332	333	334	335	336	337	338	339	340	341	342	343	344	345	346	347	348	349	350	351	352	353	354	355	356	357	358	359	360	361	362	363	364	365	366	367	368	369	370	371	372	373	374	375	376	377	378	379	380	381	382	383	384	385	386	387	388	389	390	391	392	393	394	395	396	397	398	399	400	401	402	403	404	405	406	407	408	409	410	411	412	413	414	415	416	417	418	419	420	421	422	423	424	425	426	427	428	429	430	431	432	433	434	435	436	437	438	439	440	441	442	443	444	445	446	447	448	449	450	451	452	453	454	455	456	457	458	459	460	461	462	463	464	465	466	467	468	469	470	471	472	473	474	475	476	477	478	479	480	481	482	483	484	485	486	487	488	489	490	491	492	493	494	495	496	497	498	499	500	501	502	503	504	505	506	507	508	509	510	511	512	513	514	515	516	517	518	519	520	521	522	523	524	525	526	527	528	529	530	531	532	533	534	535	536	537	538	539	540	541	542	543	544	545	546	547	548	549	550	551	552	553	554	555	556	557	558	559	560	561	562	563	564	565	566	567	568	569	570	571	572	573	574	575	576	577	578	579	580	581	582	583	584	585	586	587	588	589	590	591	592	593	594	595	596	597	598	599	600	601	602	603	604	605	606	607	608	609	610	611	612	613	614	615	616	617	618	619	620	621	622	623	624	625	626	627	628	629	630	631	632	633	634	635	636	637	638	639	640	641	642	643	644	645	646	647	648	649	650	651	652	653	654	655	656	657	658	659	660	661	662	663	664	665	666	667	668	669	670	671	672	673	674	675	676	677	678	679	680	681	682	683	684	685	686	687	688	689	690	691	692	693	694	695	696	697	698	699	700	701	702	703	704	705	706	707	708	709	710	711	712	713	714	715	716	717	718	719	720	721	722	723	724	725	726	727	728	729	730	731	732	733	734	735	736	737	738	739	740	741	742	743	744	745	746	747	748	749	750	751	752	753	754	755	756	757	758	759	760	761	762	763	764	765	766	767	768	769	770	771	772	773	774	775	776	777	778	779	780	781	782	783	784	785	786	787	788	789	790	791	792	793	794	795	796	797	798	799	800	801	802	803	804	805	806	807	808	809	810	811	812	813	814	815	816	817	818	819	820	821	822	823	824	825	826	827	828	829	830	831	832	833	834	835	836	837	838	839	840	841	842	843	844	845	846	847	848	849	850	851	852	853	854	855	856	857	858	859	860	861	862	863	864	865	866	867	868	869	870	871	872	873	874	875	876	877	878	879	880	881	882	883	884	885	886	887	888	889	890	891	892	893	894	895	896	897	898	899	900	901	902	903	904	905	906	907	908	909	910	911	912	913	914	915	916	917	918	919	920	921	922	923	924	925	926	927	928	929	930	931	932	933	934	935	936	937	938	939	940	941	942	943	944	945	946	947	948	949	950	951	952	953	954	955	956	957	958	959	960	961	962	963	964	965	966	967	968	969	970	971	972	973	974	975	976	977	978	979	980	981	982	983	984	985	986	987	988	989	990	991	992	993	994	995	996	997	998	999	1000	1001	1002	1003	1004	1005	1006	1007	1008	1009	1010	1011	1012	1013	1014	1015	1016	1017	1018	1019	1020	1021	1022	1023	1024	1025	1026	1027	1028	1029	1030	1031	1032	1033	1034	1035	1036	1037	1038	1039	1040	1041	1042	1043	1044	1045	1046	1047	1048	1049	1050	1051	1052	1053	1054	1055	1056	1057	1058	1059	1060	1061	1062	1063	1064	1065	1066	1067	1068	1069	1070	1071	1072	1073	1074	1075	1076	1077	1078	1079	1080	1081	1082	1083	1084	1085	1086	1087	1088	1089	1090	1091	1092	1093	1094	1095	1096	1097	1098	1099	1100	1101	1102	1103	1104	1105	1106	1107	1108	1109	1110	1111	1112	1113	1114	1115	1116	1117	1118	1119	1120	1121	1122	1123	1124	1125	1126	1127	1128	1129	1130	1131	1132	1133	1134	1135	1136	1137	1138	1139	1140	1141	1142	1143	1144	1145	1146	1147	1148	1149	1150	1151	1152	1153	1154	1155	1156	1157	1158	1159	1160	1161	1162	1163	1164	1165	1166	1167	1168	1169	1170	1171	1172	1173	1174	1175	1176	1177	1178	1179	1180	1181	1182	1183	1184	1185	1186	1187	1188	1189	1190	1191	1192	1193	1194	1195	1196	1197	1198	1199	1200	1201	1202	1203	1204	1205	1206	1207	1208	1209	1210	12
-----	---------	---	---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	----



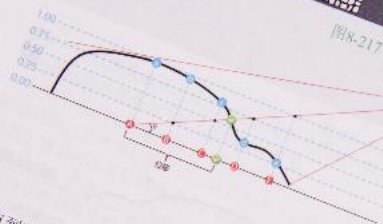
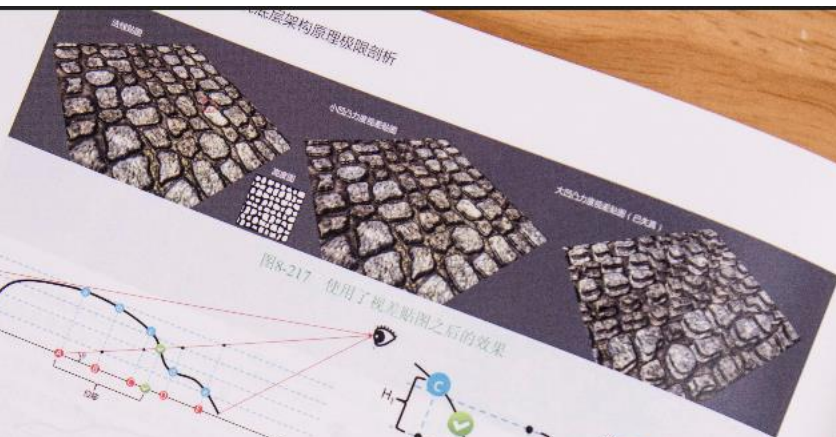


图8-218 透视图

将A点到观察者的光线分别与每一层的等高线进行求交，由于等高线为直线，所以很容易求得A、B、C、D、E各点的纹理坐标，然后从高度图中对应纹理坐标处读出这些纹理对应的假表面高度值。显而易见的一个规律是，光线穿越表面交点两边的点，也就是本例中的C和E，它们的表面高度值和层级高度值是相反的，也就是说，C点的表面高度值大于其对应的层级高度值，而E点的表面高度值小于其对应的层级高度值。交点一定位于被这样两个点夹住的表面区间内。所以，具体计算时并不需要将所有等高线对应的交点全求出，而是可以先求最高层级的点，考察其对应的表面高度值与层级高度值的关系，如果表面高度值小于层级高度值，则其继续与下一个层级高度值交并做相同比较，直到找到发生反转（表面高度值大于层级高度值）的那个点，也就可以确定两个关键点了。或者采用二分法，先从中间层级高度开始求交计

算,如果发现表面高度大于层级高度,则从再高一层  
的层级求交并计算、比较。  
然后,利用这两个点各自的表面高度与层级高度  
之间的差值,用两者的比例去插值(①和②两点间的距  
离,最终得到③点的纹理坐标值。①点与理想值并不  
重合,只是更加接近了,如图8-218右侧所示。  
这种利用寻找两个相邻且发生高度比  
交点,并在其之间水平距离  
交点的视觉贴图方

面交点的视差贴图技术, 被称为**陡峭视差 (Steep Parallax Mapping)**。如图8-219所示为SPM差贴图效果, 可以看到其凹凸得更加真实, 失真度小, 而且其精度已经高到足以实现近处像素遮挡远处的像素的真实效果了。看上去很难想象这个模型的表面原本是平坦的, 不过, 可以观察其平坦的边缘, 表明其一定是障眼法, 而没有使用置换贴图。

要提升陡峭视差贴图的近似度

提升陡峭视差贴图近似度

层级切分得  
个优化是，  
是越发正对  
以可以动态降  
在陡峭视觉贴图  
的视觉遮蔽贴图  
及迄今为止精度最高  
(mapping)。其中Oc  
障遮挡效果，其实最  
图已经可以实现遮挡。  
两个反折点之前的，  
这两个点之间寻找最  
确的算法，篇幅所  
了解。


直接把3D游戏的画面抬升了一个档次。其中就使用了视差贴图技术出现的年代：凹凸贴图（1978年）、法线贴图（1996年）、浮雕贴图（2005年）。

### 8.2.7.4 物体投影 (Shadow)

我们前文中介绍的光照效果，都是致力  
型的表面凹凸细节明暗区域，虽然表面可能  
设不出影子，但是通过法线贴图处理依然可以

上的某个氨基酸残基上的电化学生基团产生电力，从而驱动蛋白质分子上的具有钾离子选择性的分子马达的ATP结合点（油门）暴露，ATP一拥而上将其旋转起来，从一侧吸引钾离子，转到另一面，释放钾离子。从模型量底层并不是无限连续的，其最小单元，一种可能的模型是：细蛋白质分子在

但是模模糊糊。那么，  
是一个正弦波形状。那么，  
一侧的钾离子浓度如果较高，那么该蛋白质  
一侧的某些残基上的能够与钾离子相互产生物理电作  
用或者化学力作用的基团每次受力都会对整个蛋白质  
分子深处的各个原子进行牵拉，当钾离子浓度较高的  
时候，这种被牵拉的频率就会越高，当牵拉频率达到  
某个值的时候，引起分子的构象产生一个比较大的跃  
变，将分子马达的ATP结合点暴露从而转运离子。当  
钾离子浓度不够的时候，阈值频率未达到，本侧浓度持  
停止工作。随着钾离子被运输到另一侧，本侧浓度持  
续降低，最后形成一个闭环的负反馈控制系统。  
蛋白质分子就是利用原子来搭建的精密机械，  
是造物者的杰作。分子马达犹如热机一样，也分  
个冲程，比如结合ATP分子时，构象改变，导致蛋  
质旋转一定角度，ADP被释放后，转子再次旋转。  
转一定角度，ADP被释放后，转子再次旋转。  
转一定角度，ADP被释放后，转子再次旋转。



转一定角度，就像热机的吸气、压缩、膨胀这四个冲程不断循环一样。而蛋白质分子马达纯粹依靠由ATP分子形成的分子间作用力来拉动整个转子旋转。扫描二维码观看蛋白质分子马达的冲程动画。

图5所示其实是一个ATP生成器。通电可以让转子转动，对于蛋白质分子马达产生电流。对于蛋白质分子马达，通电可以产生电流。对于蛋白质分子马达，通电可以产生电流。

处是，其几乎无所不能，既可以作为离子通道，还可以作为离子泵，如图4所示。比如钠钾泵，可以将细胞外的钾离子运到胞外。那么，钠钾泵是怎么来的呢？

图5所示其实是一个ATP生成器。对于电动机来  
讲，通电可以让转子转动，但是如果让转子转动，则  
可以反过来产生电流。对于蛋白质分子马达也是一样  
的，ATP生成器也是一个分子马达，其在细胞膜内的  
区域是一个转子，转子表面包含多个质子结合点，可  
以利用高酸性环境驱动转子转动，转子转动导致胞外  
部分构象形变，从而将ADP+Pi合成为ATP，这相当于  
发电机。

对于图4左侧，诗意一些的表达则是：我要送你  
一棵生命之树，象征着你我共同的目标。树根之下  
是能量的源泉。质子驱动着马达，  
这是能量之源。质子把ADP和Pi强行结合成  
ATP。好吧，编

，诗意一些的表达则是：我要送你  
象征着你我共同的目标。树根之下  
这是能量的源泉。质子驱动着马达  
努力把ADP和P<sub>i</sub>强行结合成













## 大话 计算机

准备开脑洞了，  
请读者准备好耐心和勇气，  
与作者一同畅游  
缤纷计算机世界。

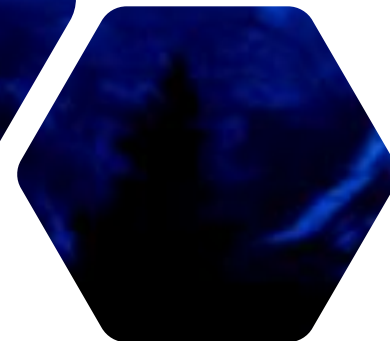
Enjoy!

清华大学出版社

清华大学出版社

## 大话 计算机

计算机系统  
原理架构原理及应用








-  01-《大话计算机》开篇-苦想计算机-2017-02-14.docx
-  02-《大话计算机》第1章-数字电路-2016-10-20-102.docx
-  03-《大话计算机》第2章-简易CPU-2016-11-12-53.docx
-  04-《大话计算机》第3章-计算机硬件史-2016-11-14-136.docx
-  05-《大话计算机》第4章-流水线与执行优化-2016-12-21-54.docx
-  06-《大话计算机》第5章-程序世界-2017-04-06-106.docx
-  07-《大话计算机》第6章-处理器微体系结构2017-07-15-182.docx
-  08-《大话计算机》第7章-计算机IO系统-2017-08-5-221 (待加入sa
-  09-《大话计算机》第8章-声音图象-2017-10-10-197 .docx


---

-  10-《大话计算机》第9章-并行计算2016-11-19-71.docx

---

-  13-《大话计算机》后记-写作经历-2016-11-28.docx
-  11-《大话计算机》第10章-现代计算机形态与生态-2016-11-19.docx
-  11-《大话计算机》第10章-现代计算机形态与生态-2016-11-19.pptx

---

-  12-《大话计算机》尾声-狂想计算机-2016-11-19.docx

---

## 第1章

### 1.1 十余年的迷惑

#### ▲ 1.2 从 $1+1=2$ 说起

- 1.2.1 用电路实现 $1+1=2$
- 1.2.2 或门/OR门
- 1.2.3 与门/AND门
- 1.2.4 非门/NOT门和与非门
- 1.2.5 异或门/XOR门
- 1.2.6 一位加法器
- 1.2.7 全手动一位加法机
- 1.2.8 实现多位加法器
- 1.2.9 电路的时延
- 1.2.10 新世界的新规律
- 1.2.11 先行/并行进位
- 1.2.12 电路化简和变换

#### ▲ 1.3 我们需要真正可用的计算器

- 1.3.1 产生记忆
- 1.3.2 解决按键问题
- 1.3.3 数学的懵懂
- 1.3.4 第一次理解数学
- 1.3.5 第一次理解语义
- 1.3.6 七段显示数码管
- 1.3.7 野路子乘法器
- 1.3.8 科班乘法器
- 1.3.9 数据交换器Crossbar
- 1.3.10 多媒体声光按键转码器
- 1.3.11 第一次驾驭时间

#### ▲ 1.4 信息与信号

- 1.4.1 录制和回放
- 1.4.2 振动和信号
- 1.4.3 低通滤波
- 1.4.4 高通滤波
- 1.4.5 带通滤波
- 1.4.6 带阻滤波
- 1.4.7 傅里叶变换
- 1.4.8 波动与电磁波
- 1.4.9 载波、调制与频分复用

#### ▲ 1.5 完整的计算器

- 1.5.1 用时序控制增强用户体验
- 1.5.2 用MUX来实现Crossbar
- 1.5.3 奇妙的FIFO队列
- 1.5.4 同步/异步FIFO
- 1.5.5 全局共享FIFO
- 1.5.6 多路仲裁
- 1.5.7 交换矩阵
- 1.5.8 时序问题的产生与触发器
- 1.5.9 擒纵机构与触发器
- 1.5.10 擒纵机构与晶振
- 1.5.11 Serdes与MUX/DEM...
- 1.5.12 计算离不开数据传递
- 1.5.13 几个专业概念的由来

#### ▲ 1.6 多功能计算器

- 1.6.1 算术逻辑单元/ALU

## 第2章

### ▲ 2.1 从累积计算说起

#### 2.1.1 实现累积计算

### ▲ 2.2 自动执行

#### 2.2.1 将操作方式的描述转化为指令

#### 2.2.2 实现那只章鱼——控制通路及部件

#### 2.2.3 动起来吧！——时序通路及部件

#### 2.2.4 半自动执行！——你得推着它跑

#### 2.2.5 全自动受控执行！——不用扬鞭自奋蹄！

#### 2.2.6 NOOP指令

#### 2.2.7 利用边沿型触发器搭建电路

#### 2.2.8 判断和跳转

#### 2.2.9 再见，Mr.章鱼！

### ▲ 2.3 更高效的执行程序

#### 2.3.1 利用循环缩减程序尺寸

#### 2.3.2 实现更多方便的指令

#### 2.3.3 多时钟周期指令

#### 2.3.4 微指令和微码

#### 2.3.5 全局地址空间

#### 2.3.6 多端口存储器

#### 2.3.7 多级缓存与CPU

#### 2.3.8 数据遍布各处

#### 2.3.9 降低数据操作粒度

#### 2.3.10 取指令/数据缓冲加速



第3章	3.3 制造工艺革命——集成电路	3.3.8 半导体工艺的瓶颈	3.4.2 电子存储器
3.1 从薄铁片到机械计算机	3.3.1 从泥活字看量产晶体管	3.3.8.1 寄生电容	3.4.2.1 静态随机存储器SRAM
3.1.1 算盘和计算尺	3.3.2 跟冬瓜哥学做P/N结蛋糕	3.3.8.2 静态/动态功耗	3.4.2.2 动态随机存储器DRAM
3.1.2 不可编程手动机械十进制计算机	3.3.3 提升集成度	3.3.8.3 栅氧厚度和High-K材料	3.4.2.3 Flash闪存
3.1.3 可编程自动机械十进制计算机	3.3.4 芯片内的深邃世界	3.3.8.4 导线连接和Low-K材料	3.4.2.4 只读存储器ROM
3.1.4 可编程自动电动机械二进制计算机	3.3.5 cMOS集成电路工艺概述	3.3.8.5 驱动能力及时延	3.4.3 光存储器
3.1.5 可编程自动全电动二进制计算机	3.3.6 cMOS工艺步骤概述	3.3.8.6 时钟树	3.4.3.1 光盘是如何存储数据的
3.2 电子管时代	3.3.7 cMOS工艺详细步骤	3.3.9 集成电路计算机	3.4.3.2 压盘与刻盘的区别
3.2.1 二极管	3.3.7.1 热氧化	3.3.10 微处理器计算机	3.4.3.3 光盘表面微观结构
3.2.2 三极管	3.3.7.2 氮化硅积淀	3.3.11 暴力拆解奔三CPU	3.4.3.4 多层记录
3.2.3 AM广播革命	3.3.7.3 浅槽隔离蚀刻	3.4 存储器-不得不说的故事	3.4.3.5 激光头的秘密
3.2.4 电子管计算机	3.3.7.4 pMOS和nMOS生成	3.4.1 机械存储器	3.4.3.6 蓝光光盘简介
3.2.5 石头会唱歌	3.3.7.5 触点电极的生成	3.4.1.1 声波/扭力波延迟线 ( Delay Line )	3.4.4 不同器件担任不同角色
3.3 固态革命——晶体管	3.3.7.6 通孔和支撑柱 ( via ) 的生成	3.4.1.2 磁鼓存储器 ( Drum )	3.4.4.1 寄存器和缓存
3.3.1 P/N结与晶体管	3.3.7.7 第一层导线连接	3.4.1.3 磁芯存储器 ( Core )	3.4.4.2 主运行内存/主存
3.3.2 场效应管 ( FET )	3.3.7.8 第二层导线连接		3.4.4.3 Scratchpad RAM
3.3.3 MOSFET	3.3.7.9 表面钝化		3.4.2.4 内容寻址内存CAM/TCAM
3.3.4 cMOS			3.4.2.5 外存
3.3.4 晶体管计算机			

## ▲ 4.1 大话流水线

### 4.1.1 不高兴的译码器

## ▲ 4.1.2 思索流水线

### 4.1.2.1 流水线的本质是并发

### 4.1.2.2 不同时延的步骤混杂

### 4.1.2.3 大话队列

### 4.1.2.4 流水线的应用及优化

## ▲ 4.2 优化流水线

### 4.2.1 拆分慢速步骤

### 4.2.2 放置多份慢速部分

### 4.2.3 加入缓冲队列

### 4.2.4 图解五级流水线指令执行过程

## ▲ 4.3 流水线冒险

### 4.3.1 访问冲突与流水线阻塞

### 4.3.2 数据依赖与数据前递

### 4.3.3 跳转冒险与分支预测

## ▲ 4.4 指令的动态调度

### 4.4.1 结构相关与寄存器重命名

### 4.4.2 保留站与乱序执行

### 4.4.3 分步图解乱序执行

### 4.4.4 重排序缓冲与指令顺序提交

## ▲ 4.5 物理并行执行

### 4.5.1 超标量和多发射

### 4.5.2 VLIW超长指令字

### 4.5.3 SIMD单指令多数据

## ▲ 第5章

### ▲ 5.1 基本的数据结构

#### 5.1.1 数组

#### 5.1.2 数据类型与ASCII码

#### 5.1.3 结构体

#### 5.1.4 数据怎么摆很重要

### ▲ 5.2 高级语言

#### 5.2.1 简单的声明和赋值

#### 5.2.2 编译和编译器

#### 5.2.3 向编译器描述数据的编排方式

#### 5.2.4 高级语言编程小试牛刀

#### 5.2.5 人脑编译忆苦思甜

### ▲ 5.3 浮点数及浮点运算

#### 5.3.1 数值范围和精度

#### 5.3.2 浮点数的用处和表示方法

#### ▲ 5.3.3 浮点数的二进制表示

##### 5.3.3.1 二进制浮点数转十进制小数

##### 5.3.3.2 十进制小数转二进制浮点数

##### 5.3.3.3 负指数和0的表示

##### 5.3.3.4 无穷与非规格化数的表示

#### 5.3.4 浮点数运算挺费劲

#### 5.3.5 浮点数的C语言声明

#### 5.3.6 十六进制表示法

### ▲ 5.4 程控多媒体计算机

#### 5.4.1 键盘是前提

#### 5.4.2 搜索并显示

#### 5.4.3 实现简易计算器

#### 5.4.4 录入和保存

#### 5.4.5 简易文件系统

#### 5.4.6 计时/定时

#### 5.4.7 发声控制

#### 5.4.8 图像显示

#### 5.4.9 网络聊天

### ▲ 5.5 程序社会

#### 5.5.1 函数和调用

#### 5.5.2 设备驱动程序

#### 5.5.3 函数之间的联络站

#### ▲ 5.5.4 库和链接

##### 5.5.4.1 静态库和静态链接

##### 5.5.4.2 头文件

##### 5.5.4.3 API和SDK

##### 5.5.4.4 动态库和动态链接

##### 5.5.4.5 库文件/可执行文件的格式

### ▲ 5.5.5 程序的执行和退出

#### 5.5.5.1 初步解决地址问题

#### 5.5.5.3 更好的人机交互方式

#### 5.5.5.4 程序的退出

#### 5.5.5.5 使用外部设备和内存

### ▲ 5.5.6 多程序并发执行

#### 5.5.6.1 利用时钟中断来切换线程

#### 5.5.6.2 更广泛的使用中断

#### 5.5.6.3 虚拟地址空间与分页

#### 5.5.6.4 虚拟与现实的边界——系统调用

#### 5.5.7 呼唤操作系统

### ▲ 5.6 计算机操作系统

#### ▲ 5.6.1 内存管理

##### 5.6.1.1 分区式内存管理

##### 5.6.1.2 分段+分区式内存管理

##### 5.6.1.3 Intel CPU的分段支持机制

##### 5.6.1.4 分页+分段式内存管理

##### 5.6.1.5 DOS和Linux操作系统的内存管理

#### 5.6.2 中断响应及处理

#### 5.6.3 驱动程序与设备管理

#### 5.6.7 文件系统和I/O协议栈

#### 5.6.5 系统调用

#### 5.6.4 线程调度与管理

#### 5.6.6 进程间通信

#### 5.6.8 人机交互界面



第6章	6.2.15 Load/Stor Queue与Stream Buffer	6.5 QPI片间互连网络简介	6.8.6 解决提前执行测错乱问题
6.1 从超线程到多核心	6.2.16 非阻塞Cache与MSHR	6.5.1 QPI物理层与同步异步通信原理	6.8.7 解决乱序执行错乱问题
6.1.1 超线程并行	6.2.17 缓存行替换策略	6.5.2 QPI链路层网络层和消息层	6.8.8 小结
6.1.2 多核心/多CPU并行	6.2.18 I_Cache/D_Cache/TLB_Cache	6.5.3 QPI的初始化与系统启动	6.9 解决多核心访存空间一致性问题
6.1.3 idle线程	6.2.19 对齐和伪共享	6.5.3.1 链路初始化和拓扑发现	6.9.1 基于总线监听的缓存一致性实现
6.1.4 乱序执行还是SMT ?	6.3 连起来，为了一致性	6.5.3.2 系统启动	6.9.1.1 Snarfin/Write Sync方式
6.1.5 逆超线程 ?	6.3.1 Crossbar交换矩阵	6.5.4 QPI的扩展性	6.9.1.2 Write Invalidate方式
6.1.6 线程与进程	6.3.1.1 点阵式Crossbar	6.6 基于QPI互联的高端服务器架构一览	6.9.2 推导MESIF状态机
6.1.7 多核心访存基本拓扑	6.3.1.2 基于复用器的Crossbar	6.6.1 某32路CPU高端主机	6.9.3 MOESI状态机
6.2 缓存十九式	6.3.1.3 Crossbar的级联	6.6.2 IBM x3850/3950 X5/X6主机	6.9.4 结合MESIF协议进一步理解锁和屏障
6.2.1 缓存是分级的	6.3.2 Ring环网	6.6.3 HP Superdome2主机	6.9.5 结合MESIF深刻理解时序一致性模型
6.2.2 缓存是透明的	6.3.3 NoC片上网络	6.6.4 Fujitsu PQ2K主机	6.9.5.1 终极一致性 ( UC )
6.2.3 缓存的容量、频率和延迟	6.3.4 众核Many-Core	6.7 理解多核心访存时空一致性问题	6.9.5.2 严格一致性 ( SC )
6.2.4 私有缓存和共享缓存	6.3.5 多核心程序执行过程回顾	6.7.1 访存空间一致性问题	6.9.5.3 顺序一致性 ( SEC )
6.2.5 Inclusive和Exclusive模式	6.3.6 在众核心上执行程序	6.7.2 访存时间一致性问题	6.9.5.4 处理器一致性 ( PC )
6.2.6 Dirty和Valid标记位	6.4 存储器在网络中的分布	6.7.2.1 延迟到达导致的错乱	6.9.5.5 弱一致性 ( WC )
6.2.7 缓存行Cache Line	6.4.1 CPU片内访存网络与存储器分布	6.7.2.2 访问冲突导致的错乱	6.9.6 缓存行并发写优化
6.2.8 全关联/直接关联/组关联	6.4.2 CPU片外访存网络	6.7.2.3 提前执行导致的错乱	6.9.7 Cache Agent的位置
6.2.9 Virtual Cache	6.4.2.1 全总线拓扑及南桥与北桥	6.7.2.4 乱序执行导致的错乱	6.9.8 基于共享总线的嗅探过滤机制
6.2.10 缓存Homonym问题	6.4.2.2 AMD Athlon北桥	6.8 解决多核心访存时间一致性问题	6.9.8.1 Bitmap粗略过滤
6.2.11 缓存Alias问题	6.4.2.3 常用网络拓扑及UMA/NUMA	6.8.1 有你没我，有我没你	6.9.8.2 Vector Bitmap精确过滤
6.2.12 Page Coloring页面着色	6.4.2.4 AMD Opteron北桥	6.8.2 让子弹飞，等响声来	6.9.8.3 Bloom filter布隆过滤器与哈希采样
6.2.13 小结及商用CPU的缓存模式	6.4.3 参悟全局共享内存架构	6.8.3 硬件原生保证的基本时序	6.9.8.4 JETTY filter
6.2.14 缓存对写入操作的处理	6.4.4 访存网络的硬分区	6.8.4 解决延迟到达错乱问题	6.9.8.5 Stream Register Filter
		6.8.5 解决访问冲突错乱问题	6.9.8.6 Counter Stream Register Filter

## ▲ 第7章

### ▲ 7.1 计算机I/O的基本套路

7.1.1 Programmed IO+Polling模式

7.1.2 DMA+中断模式

7.1.3 DMA与缓存一致性

7.1.4 Scatter/Gather List ( SGL )

7.1.5 使用队列提升I/O性能

### ▲ 7.1.6 固件/Firmware

7.1.6.1 固件与OS的区别与联系

7.1.6.2 固件的层次

7.1.6.3 固件的格式

7.1.6.4 固件存在哪

7.1.6.5 固件如何加载运行

7.1.7 网络I/O基本套路

7.1.8 接入更多外部设备

7.1.9 一台完整计算机的全貌

7.2 中断处理

### ▲ 7.3 网络通信系统

#### ▲ 7.3.1 OSI七层标准模型

7.3.1.1 应用层

7.3.1.2 表示层

7.3.1.3 会话层

7.3.1.4 传输层

7.3.1.5 网络层

7.3.1.6 链路层

7.3.1.7 物理层

7.3.1.8 传送层

7.3.1.9 小结

#### ▲ 7.3.2 底层信号处理系统

7.3.2.1 AC耦合电容及N/Mb编码

7.3.2.2 加扰的作用

7.3.2.3 各种线路编码

7.3.2.4 各种模拟调制技术

7.3.2.5 频谱宽度与比特率

7.3.2.6 数字信号处理与数字滤波

#### ▲ 7.3.3 以太网——高速通用非访存式后端外部网络

7.3.3.1 以太网的网络层

7.3.3.2 以太网的链路层和物理层

7.3.3.3 以太网I/O控制器

### ▲ 7.4 典型I/O网络简介

#### ▲ 7.4.1 PCIE——高速通用访存式前端I/O网络

7.4.1.1 PCI网络拓扑及数据收发过程

7.4.1.2 PCI设备的配置空间

7.4.1.3 PCI设备的枚举和配置

7.4.1.4 PCI设备寄存器的物理地址分配和路由

7.4.1.5 中期小结

7.4.1.6 PCIE网络拓扑及数据收发过程

7.4.1.7 PCIE网络的层次模型

7.4.1.8 NTB非透明桥

7.4.1.9 PCIE Switch内部

7.4.1.10 在PCIE网络中传递消息

7.4.1.11 在PCI网络中传递中断信号

7.4.1.12 在PCIE网络中传递中断信号

7.4.1.13 MSI/MIS-X中断方式

7.4.1.14 PCIE体系中的驱动程序层次

7.4.1.15 小结

#### ▲ 7.4.2 USB——中速通用非访存式后端I/O网络

7.4.2.1 USB网络的基本拓扑

7.4.2.2 USB设备的枚举和配置

7.4.2.3 USB网络协议栈

7.4.2.4 USB网络上的数据包传送

7.4.2.5 USB网络的层次模型

7.4.2.6 小结

#### ▲ 7.4.3 SAS——高速专用非访存式后端I/O网络

7.4.3.1 SAS网络拓扑及设备编号规则

7.4.3.2 SAS网络中的Order Set一览

7.4.3.3 SAS的链路初始化和速率协商

7.4.3.4 SAS网络的初始化与设备枚举

7.4.3.5 SAS和SCSI的Host端协议栈

7.4.3.6 形形色色的登记表

7.4.3.7 SAS网络的数据传输方式

7.4.3.8 SAS网络的层次模型

7.4.3.9 SAS控制器内部架构

7.4.3.10 SAS Expander内部架构

7.4.3.11 基于SAS体系的存储系统形态

## 第8章

### 8.1 声音处理系统

#### 8.1.1 让蜂鸣器说话

#### 8.1.2 音乐是可以被勾兑出来的

##### 8.1.2.1 可编程音符生成器PSG

##### 8.1.2.2 音乐合成器

##### 8.1.2.3 FM合成及波表合成

#### 8.1.3 声卡发展史及架构简析

#### 8.1.4 与发声控制相关的Host端角色

#### 8.1.5 让计算机成为演奏家

#### 8.1.6 独立声卡的没落

### 8.2 图形处理系统

#### 8.2.1 用声音来画图

#### 8.2.2 文字模式

##### 8.2.2.1 向量文本模式显示

##### 8.2.2.2 用ROM存放字形库

##### 8.2.2.3 点阵文字显示模式

##### 8.2.2.4 Monochrome Display Adapter ( ...

##### 8.2.2.5 点阵作图与ASCII Art

#### 8.2.3 图形模式

##### 8.2.3.1 Color Graphics Adapter ( CGA )

##### 8.2.3.2 Enhanced Graphics Adapter ( EG...

##### 8.2.3.3 Video BIOS ROM的引入

##### 8.2.3.4 Video Graphics Array ( VGA )

##### 8.2.3.5 VGA的后续

##### 8.2.3.6 当代显卡的图形和文字模式

### 8.2.4 2D图形及其渲染流程

#### 8.2.4.1 2D图形加速卡PGC

#### 8.2.4.2 2D图形模型的准备

#### 8.2.4.3 对模型进行渲染

#### 8.2.4.4 矢量图和bitmap

#### 8.2.4.5 顶点、索引和图元

#### 8.2.4.6 2D图形动画

#### 8.2.4.7 坐标变换及矩阵运算

#### 8.2.4.8 2D图形渲染流程小结

#### 8.2.4.9 2D绘图库以及渲染加速

### 8.2.5 3D图形模型和表示方法

#### 8.2.5.1 3D模型的表示

#### 8.2.5.2 顶点的4个基本属性

### 8.2.6 3D图形渲染流程

#### 8.2.6.1 顶点坐标变换阶段/Vertex Transform

#### 8.2.6.2 顶点光照计算阶段/Vertex Lighting

#### 8.2.6.3 栅格化阶段/Rasterization

#### 8.2.6.4 像素着色阶段/Pixel Shading

#### 8.2.6.5 遮挡判断阶段/Testing

#### 8.2.6.6 混合及后处理阶段/Blending

#### 8.2.6.7 3D渲染流程小结

### 8.2.7 典型的3D渲染特效简介

#### 8.2.7.1 法线贴图 ( Normal Map )

#### 8.2.7.2 曲面细分与置换贴图 ( Tessellation

#### 8.2.7.3 视差/位移贴图 ( Parallax Map )

#### 8.2.7.4 物体投影 ( Shadow )

#### 8.2.7.5 抗锯齿 ( Anti-Aliasing )

#### 8.2.7.6 光照控制纹理 ( Light Mapping )

#### 8.2.7.7 纹理动画 ( Texture Animation )

### 8.2.8 当代3D游戏制作过程

### 8.2.9 3D图形加速渲染

#### 8.2.9.1 3D图形渲染管线回顾

#### 8.2.9.2 固定渲染管线3D图形加速

#### 8.2.9.3 可编程渲染管线3D图形加速

#### 8.2.9.4 Unified可编程3D图形加速

#### 8.2.9.5 深入AMD R600 GPU内部执行流程

### 8.2.10 3D绘图API 及软件栈

#### 8.2.10.1 GPU内核态驱动及命令的下发

#### 8.2.10.2 GPU用户态驱动及命令的翻译

#### 8.2.10.3 久违了OpenGL与Direct3D

#### 8.2.10.4 Windows下图形软件栈

### 8.2.11 3D图形加速卡的辉煌时代

#### 8.2.11.1 街机/家用机/手机上的GPU

#### 8.2.11.2 SGI Onyx超级图形加速工作站

#### 8.2.11.3 S3 ViRGE时代

#### 8.2.11.4 3dfx Voodoo时代

#### 8.2.11.5 NVidia和ATI时代

## 8.3 结语和期盼



- 第9章
  - 9.1 科学计算到底在算些什么
    - 9.1.1 蛋白质分子的故事
      - 9.1.1.1 氧气运输的故事
      - 9.1.1.2 更复杂的生化逻辑是如何完成的
    - 9.1.2 如何模拟蛋白质分子自折叠过程
    - 9.1.3 如何将模拟过程映射为多线程并行计算
    - 9.1.4 其他科学计算场景
  - 9.2 大规模系统共享内存之向往
    - 9.2.1 UMA/NUMA/MPP
    - 9.2.2 OpenMP并行编程
  - 9.3 基于消息传递的非共享内存系统
    - 9.3.1 采用消息传递方式同步数据
    - 9.3.2 MPI库基本函数简介
    - 9.3.3 MPI库聚合通信函数简介
  - 9.4 超级计算机
    - 9.4.1 IBM蓝色基因
      - 9.4.1.1 中央处理器CPU
      - 9.4.1.2 计算节点和I/O节点
      - 9.4.1.3 三个独立网络同时传递数据
      - 9.4.1.4 蓝色基因Q的网络控制和路由实现
      - 9.4.1.5 节点卡及整机架布局
      - 9.4.1.6 整体系统拓扑
      - 9.4.1.7 Service Node
      - 9.4.1.8 Service Card
      - 9.4.1.9 时钟同步
      - 9.4.1.10 系统启动
      - 9.4.1.11 软件安装与用户认证

- 9.4.1.12 状态监控
- 9.4.1.13 计算任务的执行
- 9.4.1.14 操作系统
- 9.4.1.15 Login Node
- 9.4.1.16 存储系统
- 9.4.2 圣地亚哥Gordon
- 9.4.3 Fujitsu PrimeHPC FX10
  - 9.4.3.1 SPARC64 IXfx CPU
  - 9.4.3.2 水冷主板
  - 9.4.3.3 Tofu六维网络互联拓扑
  - 9.4.3.4 ICC互联芯片
- 9.5 利用GPU加速计算
  - 9.5.1 Direct3D中的Compute Shader
  - 9.5.2 OpenCL和OpenACC
  - 9.5.3 NVidia的CUDA API
    - 9.5.3.1 CUDA基本架构
    - 9.5.3.2 一个极简的CUDA程序
    - 9.5.3.3 CUDA对线程的编排方式
    - 9.5.3.4 GPU对CUDA线程的调度方式
    - 9.5.3.5 CUDA程序的内存架构
    - 9.5.3.6 基于CUDA的PhysX库效果
- 9.6 利用PLD和ASIC加速计算
  - 9.6.1 PAL/PLA是如何工作的
  - 9.6.2 CPLD是如何工作的
  - 9.6.3 FPGA是如何工作的
  - 9.6.4 FPGA编程及应用形态
  - 9.6.4 Xilinx FPGA架构及相关概念
  - 9.6.6 ASIC与PLD的关系
- 9.7 小结：软归软 硬归硬

▲ 第10章 计算机操作系统——舞台幕后的工作者	10.1.6.4 高位内存区 (HMA)	10.2.2.2 Linux的任务软切换机制	10.4.2 用户态和内核态抢占
▲ 10.1 内存布局与管理	10.1.6.5 扩展内存 (XMS)	10.2.2.3 进程0的创建和运行	10.4.3 中期小结
10.1.1 实模式与保护模式	10.1.6.6 用XMS顶替EMS	10.2.2.4 进程1和2的创建和运行	10.4.4 实时与非实时内核
10.1.2 分区式内存管理	▲ 10.1.7 后DOS时代x86内存布局	10.2.2.5 在用户态创建和运行任务	▲ 10.4.5 任务调度基本数据结构
▲ 10.1.3 8086分段+实模式	10.1.7.1 E820表	10.2.2.6 fork( )自测题及深入思考	10.4.5.1 任务优先级描述
10.1.3.1 线性/逻辑/物理地址	10.1.7.2 物理地址扩展PAE	10.2.2.7 用户空间线程/协程	10.4.5.2 三大子调度器
▲ 10.1.4 80286分段+保护模式	10.1.7.3 x86物理内存布局	10.2.2.8 任务状态	10.4.5.3 运行队列的组织
10.1.4.1 GDT表	▲ 10.1.8 Linux下的内存管理	▲ 10.3 任务间通信与同步	▲ 10.5 任务调度核心方法
10.1.4.2 实现权限检查	10.1.8.1 32位Linux内存布局	10.3.1 信号极其处理	10.5.1 简单粗暴的实时任务调度
10.1.4.3 LDT表	10.1.8.2 相关模块数据结构	10.3.2 等待队列与唤醒	10.5.2 左右为难的普通任务调度
▲ 10.1.5 80386分段+分页模式	10.1.8.3 brk和mmap系统调用	10.3.3 进程间通信	10.5.3 2.4内核中的O(n)调度器
10.1.5.1 页目录/页表/页面	10.1.8.4 malloc/calloc/realloc函数	▲ 10.3.4 锁和同步	10.5.4 2.5内核中的O(1)调度器
10.1.5.2 比较分页和分段机制	10.1.8.5 buddy和slab算法	10.3.4.1 信号量 (Semaphore)	10.5.5 未被接纳的RSDL普通任务
10.1.5.3 Flat分段模式	▲ 10.2 任务创建与管理	10.3.4.2 互斥量 (Mutex)	▲ 10.5.6 沿用至今的CFS普通任务
10.1.5.4 分页的控制参数	▲ 10.2.1 32位x86处理器任务管理支持	10.3.4.3 自旋锁 (Spinlock)	10.5.6.1 指挥棒变为运行时
10.1.5.5 MMU和TLB	10.2.1.1 用户栈与内核栈	10.3.4.4 快速互斥量 (Futex)	10.5.6.2 weight/period/vr
▲ 10.1.6 DOS下的内存管理	10.2.1.2 线程和中断上下文	10.3.4.5 条件量 (Condition)	10.5.7 多处理器任务负载均衡
10.1.6.1 常规内存和上位内存	10.2.1.3 任务切换机制	10.3.4.6 完成量 (Completion)	10.5.8 任务的Affinity
10.1.6.2 EMS内存扩充卡	10.2.1.4 任务嵌套/任务链	10.3.4.7 读写锁 (RWlock) 和RCU锁	▲ 10.6 中断响应及处理
10.1.6.3 上位内存块 (UMB)	10.2.1.5 小结	▲ 10.4 任务调度基本框架	▲ 10.6.1 中断相关基本知识
10.1.6.4 高位内存区 (HMA)	▲ 10.2.2 32位Linux的任务创建与管理	10.4.1 任务的调度时机	10.6.1.1 Local和I/O APIC
	10.2.2.1 PCB/task_struct{ }		
			10.6.1.2 8259A (PIC) 中断控制器
			10.6.1.3 MSI/MSI-X底层实现
			10.6.1.4 IPI处理器间中断
			10.6.1.5 可屏蔽/不可屏蔽中断
			10.6.1.6 中断的共享和/嵌套
			10.6.1.7 中断内部/外部优先级
			10.6.1.8 中断Affinity及均衡
			▲ 10.6.2 中断相关数据结构
			10.6.2.1 中断描述符表IDT
			10.6.2.2 irq_desc[ ]和vector_irq[ ]
			10.6.2.3 相关数据结构的初始化
			10.6.3 中断基本处理流程
			10.6.4 80h号中断 (系统调用)
			▲ 10.6.5 中断上半部和下半部
			10.6.5.1 softirq
			10.6.5.2 ksoftirqd线程
			10.6.5.3 softirq与preempt_count
			10.6.5.4 tasklet
			10.6.5.5 workqueue
			10.6.6 中断线程化
			10.6.7 系统的驱动力
			▲ 10.7 时间管理与时钟中断

#### ▲ 10.7.1 表哥的收藏

10.7.1.1 RTC

10.7.1.2 PIT

10.7.1.3 HPET

10.7.1.4 Local Timer

10.7.1.5 TSC

#### ▲ 10.7.2 表哥的烦恼

10.7.2.1 软计时

10.7.2.2 软Timer

10.7.2.3 软Tick

10.7.2.4 单调时钟源

10.7.2.5 中断广播唤醒

10.7.2.6 强制周期性中断广播

#### ▲ 10.7.3 表哥的记忆

10.7.3.1 Clocksource Device

10.7.3.2 Clockevent Device

10.7.3.3 Local/Global Device

10.7.3.4 HZ/Jiffy/NOHZ

10.7.3.5 各种时间种类

10.7.3.6 低精度定时器时间轮

10.7.3.7 高精度定时器红黑树

#### ▲ 10.7.4 表哥的思维

10.7.4.1 tick\_init()

10.7.4.2 init\_timers()

10.7.4.3 hrtimers\_init()

10.7.4.4 timekeeping\_init()

10.7.4.5 time\_init()/late\_time\_init()

10.7.4.6 APIC\_init\_uniprocessor()

10.7.4.7 do\_basic\_setup()/do\_initcalls()

10.7.4.8 初始化流程全局图

#### ▲ 10.7.5 表哥的行动

10.7.5.1 初始的低精度+HZ模式

10.7.5.2 切换到低精度+NOHZ模式

10.7.5.3 切换到高精度+NOHZ模式

10.7.5.4 idle与NOHZ

10.7.5.5 切换高精度模式流程图

#### ▲ 10.8 VFS与本地FS

##### ▲ 10.8.1 VFS目录层

10.8.1.1 目录与VFS

10.8.1.2 目录承载者

10.8.2 本地FS相关数据结构

##### ▲ 10.8.3 VFS相关数据结构及初始化

##### ▲ 10.8.3 VFS相关数据结构及初始化

10.8.3.1 Mount流程

10.8.3.2 Open流程

10.8.4 从read到Page Cache

10.8.5 从Page Cache到通用块层

##### ▲ 10.8.6 Linux下的异步I/O

10.8.6.1 基于glibc的异步I/O

10.8.6.2 基于libaio的异步I/O

#### ▲ 10.9 块I/O协议栈

##### ▲ 10.9.1 从通用块层到I/O调度层

10.9.1.1 块设备与buffer page

10.9.1.2 bio

##### ▲ 10.9.2 从I/O调度层到块设备驱动

10.9.2.1 Request与Request Queue

10.9.2.2 堵盖儿和掀盖儿

10.9.2.3 \_\_make\_request主流程

10.9.2.4 I/O Scheduler

##### ▲ 10.9.3 相关数据结构的初始化

10.9.3.1 request\_queue初始化

10.9.3.2 gendisk/scsi\_disk/block\_device...

10.9.4 从块设备驱动到SCSI中间层

10.9.5 从SCSI中间层到通道控制器驱动

10.9.5 从SCSI中间层到通道控制器驱动

#### ▲ 10.10 网络I/O协议栈

10.10.1 socket的初始化

10.10.2 socket的创建和绑定

10.10.3 发起TCP连接

10.13 小结



## 第11章 现代计算机系统形态与生态

### ▲ 11.1 工业级相关计算机产品

11.1.1 工业控制

11.1.2 军工和航空航天

### ▲ 11.2 企业级相关计算机产品

11.2.1 芯片与板卡

#### ▲ 11.2.2 服务器

11.2.2.1 塔式服务器

11.2.2.2 机架式服务器

11.2.2.3 刀片服务器

11.2.2.4 模块化服务器

11.2.2.5 整机柜服务器

11.2.2.6 关键应用主机

#### ▲ 11.2.3 网络系统

11.2.3.1 以太网卡

11.2.3.2 以太网交换机和路由器

#### ▲ 11.2.4 存储系统

11.2.4.1 机械磁盘

11.2.4.2 固态硬盘

11.2.4.3 SAN存储系统

11.2.4.4 分布式存储系统

11.2.4.5 数据恢复

11.2.5 超融合系统

11.2.6 数据备份和容灾系统

11.2.7 云计算和云存储

11.2.8 自主可控系统

### ▲ 11.3 消费级相关计算机产品

11.3.1 智能手机

11.3.2 电视盒/智能电视

11.3.3 摄像机

11.3.4 玩具

## ▲ 第12章：机器学习与人工智能

12.1 回归分析：愚者千虑必有一得

12.2 逻辑分类：不是什么都能一刀切

12.3 神经网络：竟可万能拟合

12.4 深度神经网络：四两拨千斤

12.5 对象检测：先抠图后识别

12.6 卷积神经网络：图像识别利器

12.7 可视化展现：盲人真的摸出了象

12.8 具体实现：搭台唱戏和硬功夫

12.9 人工智能：本能、智能、超能

## 第1章

### 1.1 十余年的迷惑

#### ▲ 1.2 从 $1+1=2$ 说起

- 1.2.1 用电路实现 $1+1=2$
- 1.2.2 或门/OR门
- 1.2.3 与门/AND门
- 1.2.4 非门/NOT门和与非门
- 1.2.5 异或门/XOR门
- 1.2.6 一位加法器
- 1.2.7 全手动一位加法机
- 1.2.8 实现多位加法器
- 1.2.9 电路的时延
- 1.2.10 新世界的新规律
- 1.2.11 先行/并行进位
- 1.2.12 电路化简和变换

#### ▲ 1.3 我们需要真正可用的计算器

- 1.3.1 产生记忆
- 1.3.2 解决按键问题
- 1.3.3 数学的懵懂
- 1.3.4 第一次理解数学
- 1.3.5 第一次理解语义
- 1.3.6 七段显示数码管
- 1.3.7 野路子乘法器
- 1.3.8 科班乘法器
- 1.3.9 数据交换器Crossbar
- 1.3.10 多媒体声光按键转码器
- 1.3.11 第一次驾驭时间

#### ▲ 1.4 信息与信号

- 1.4.1 录制和回放
- 1.4.2 振动和信号
- 1.4.3 低通滤波
- 1.4.4 高通滤波
- 1.4.5 带通滤波
- 1.4.6 带阻滤波
- 1.4.7 傅里叶变换
- 1.4.8 波动与电磁波
- 1.4.9 载波、调制与频分复用

#### ▲ 1.5 完整的计算器

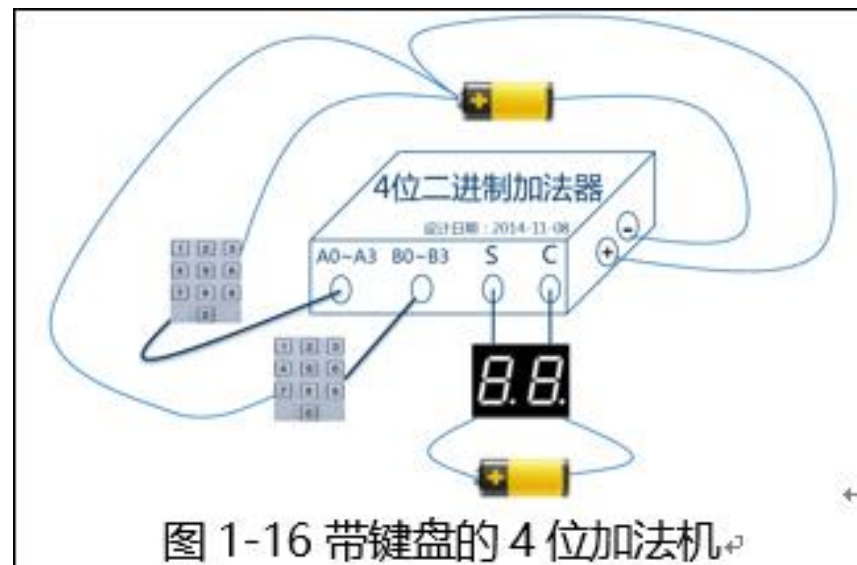
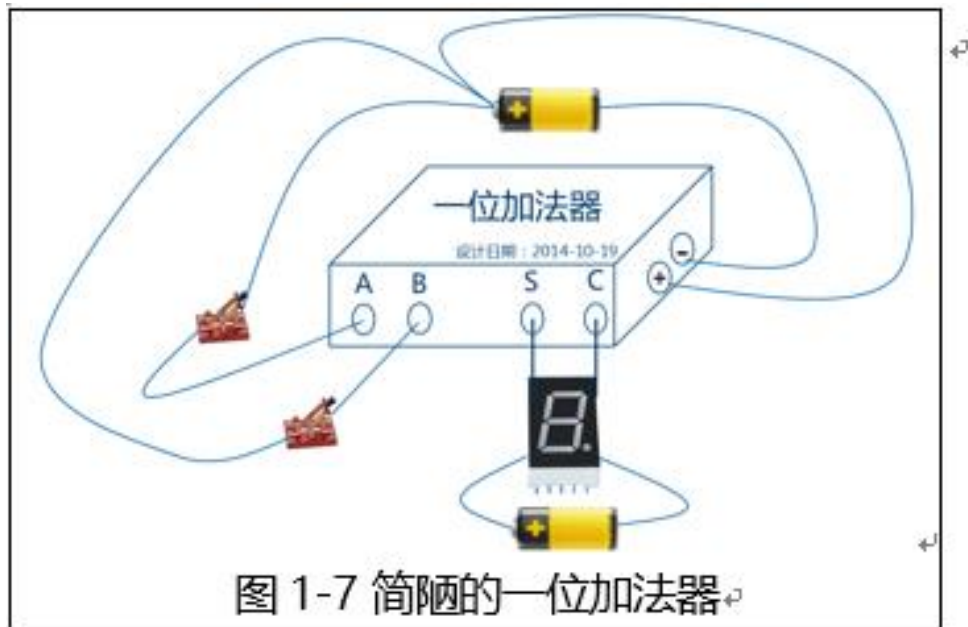
- 1.5.1 用时序控制增强用户体验
- 1.5.2 用MUX来实现Crossbar
- 1.5.3 奇妙的FIFO队列
- 1.5.4 同步/异步FIFO
- 1.5.5 全局共享FIFO
- 1.5.6 多路仲裁
- 1.5.7 交换矩阵
- 1.5.8 时序问题的产生与触发器
- 1.5.9 擒纵机构与触发器
- 1.5.10 擒纵机构与晶振
- 1.5.11 Serdes与MUX/DEM...
- 1.5.12 计算离不开数据传递
- 1.5.13 几个专业概念的由来

#### ▲ 1.6 多功能计算器

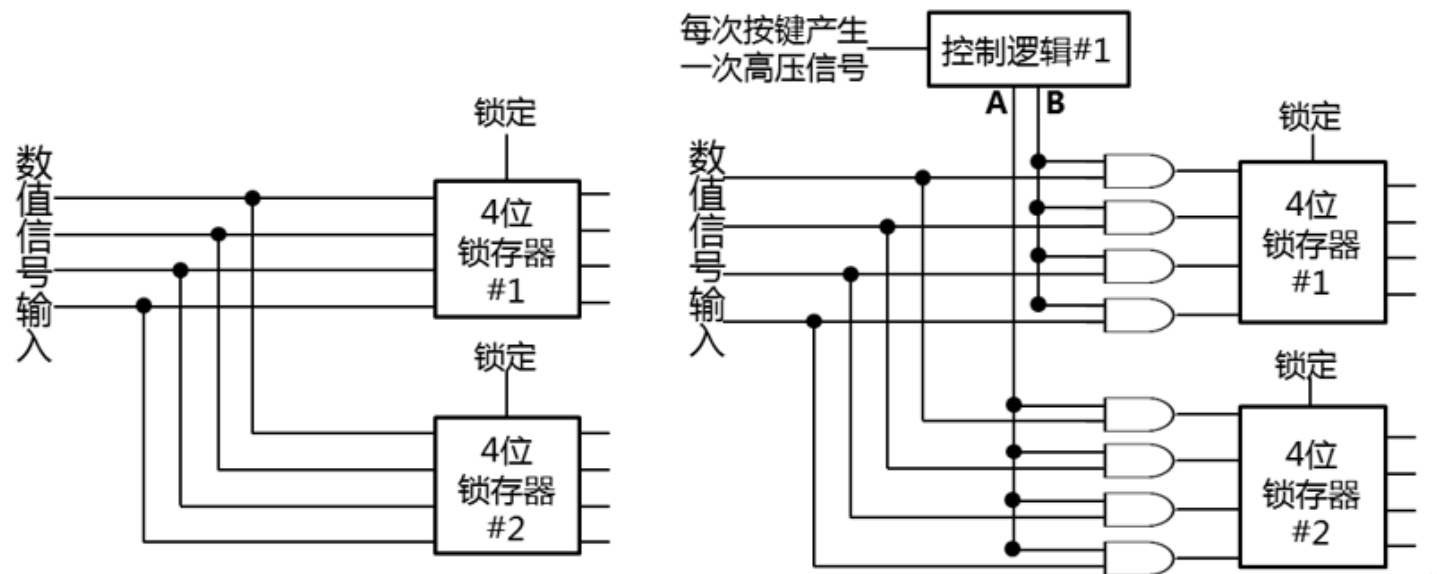
- 1.6.1 算术逻辑单元/ALU

写作思路：先提出问题，计算机到底是怎么算 $1+1=2$ 的。给出思考过程，给出答案，引出基本的组合逻辑电路，然后引出加法器。再提出问题，如何设计一个按键式计算机来算加法，给出思考过程，遇到了各种工程问题，一步步解决，最后引出译码器、触发器等逻辑电路。最后逐渐搭建出一个可以算多位数加减乘除的按键式计算器。最后引出ALU的概念。





与门和或门的组合，能实现任何逻辑功能，主要元件不用控制。



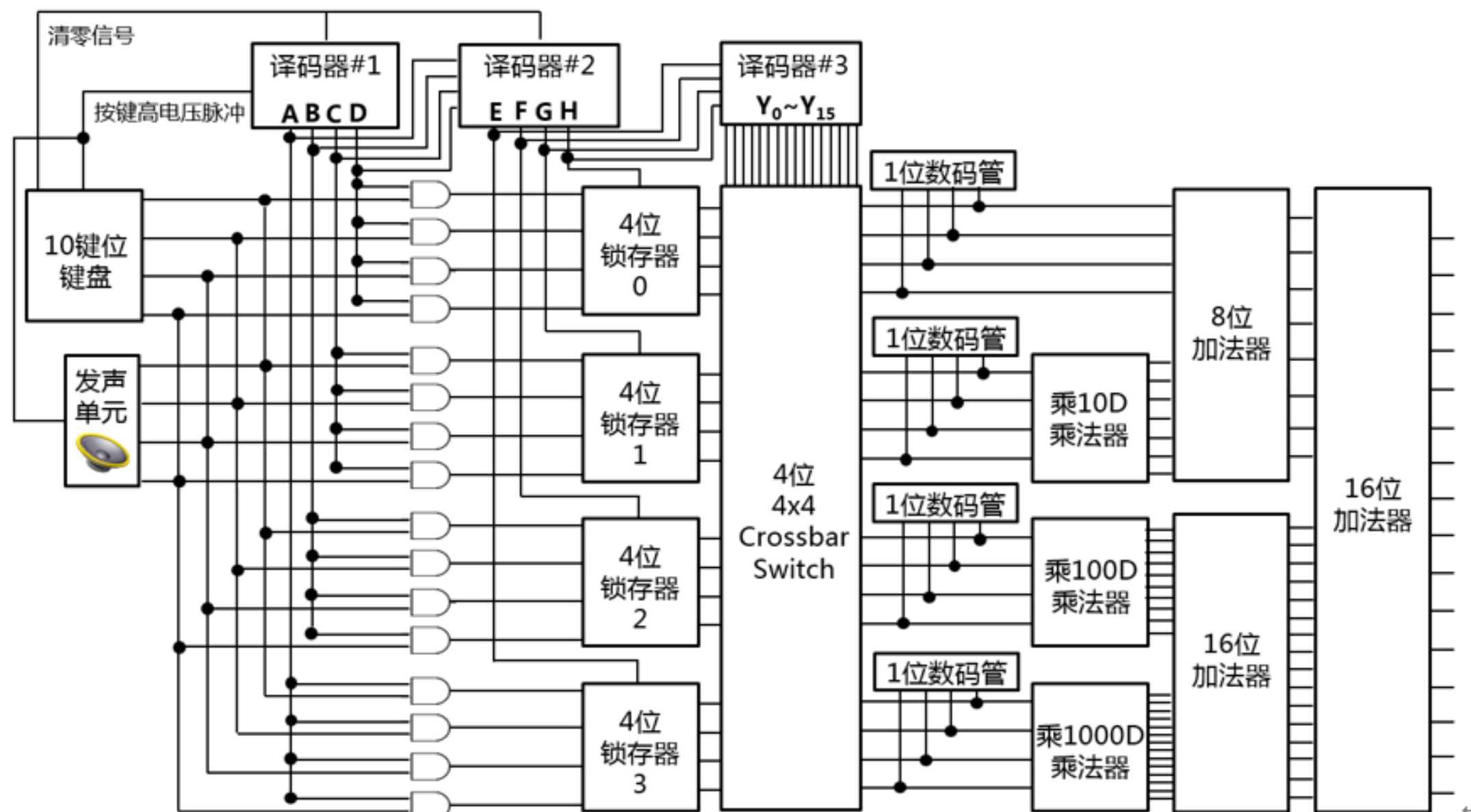


图 1-97 可发声按键转码器架构图

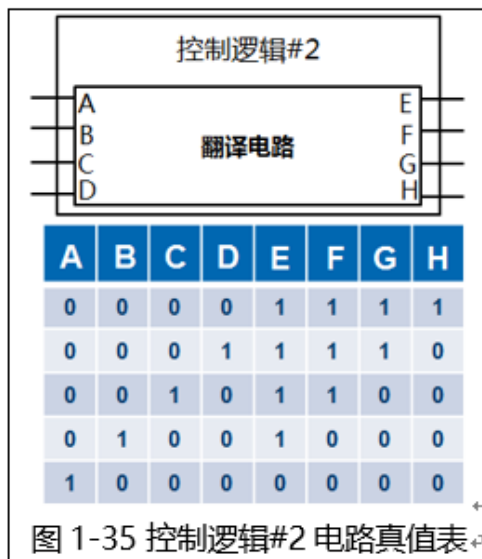


图 1-35 控制逻辑#2 电路真值表

同理，我们使用上文中给出的套路，写出 EFGH 这四个输出的表达式，结果遇到了问题。对于 E 这个输出，五行里有 4 行都是 1，如果按照第一行来给出表达式，则为  $E = A' \cdot B' \cdot C' \cdot D'$ ，而将第二行的值代入这个表达式，会得出与现实不一样的结果，这就证明，我们之前使用的套路，要么出了问题，要么并不能完整的描述所有场景而只是特例。这可怎么办，空欢喜一场。但是至少  $E = A' \cdot B' \cdot C' \cdot D'$  可以描述第一行，那么第二行既然不能用这个式子描述，那肯定也有其对应的表达式，专门为第二行写出来，就是  $E = A' \cdot B' \cdot C' \cdot D$ ，第三行则是  $E = A' \cdot B' \cdot C \cdot D'$ ，第四行  $E = A' \cdot B \cdot C' \cdot D'$ ，第五行由于  $E = 0$ ，不用写表达式。这四行都对，又都不对，那到底是个什么状态？任何一种输入组合，要么匹配第一行，要么第二行，总之总得匹配某一行，这思路就来了，既然是“要么”，那么“或”这个逻辑，是否可以用来表达这个规律？我们把这五个式子或一下， $E =$

$A' \cdot B' \cdot C' \cdot D' + A' \cdot B' \cdot C' \cdot D + A' \cdot B' \cdot C \cdot D' + A' \cdot B \cdot C' \cdot D'$ ，这个式子表达了这样一个意思，任何一组输入组合，同时输入到这 4 个子等式里，一旦某个子等式不匹配，则这个子等式的输出一定是 0，但是这 4 个子等式中总有一个匹配，那么其输出为 1，而因为这 4 个子等式的输出是相或在一起的，所以整个等式的输出就为 1，正好匹配了  $E = 1$  的现实。

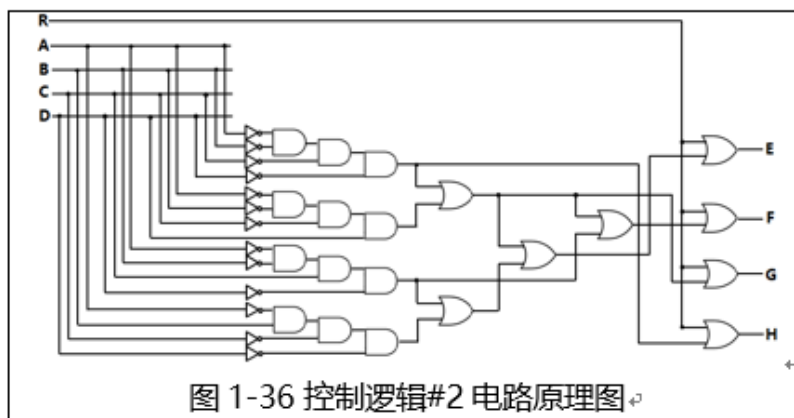


图 1-36 控制逻辑#2 电路原理图

代入验证发现，这个等式的确是正确的，其的确可以完整的描述一个真值表逻辑。同理，我们写出 EFGH 四个输出的各自的表达式：

$$E = A' \cdot B' \cdot C' \cdot D' + A' \cdot B' \cdot C' \cdot D + A' \cdot B' \cdot C \cdot D' + A' \cdot B \cdot C' \cdot D'$$

$$F = A' \cdot B' \cdot C' \cdot D' + A' \cdot B' \cdot C' \cdot D + A' \cdot B' \cdot C \cdot D'$$

$$G = A' \cdot B' \cdot C' \cdot D' + A' \cdot B' \cdot C' \cdot D$$

$$H = A' \cdot B' \cdot C' \cdot D'$$

根据表达式画出对应的电路如图 1-36 所示。这里清零逻辑使用了或门，因为对于控制逻辑#2，清零之后必将 EFGH 全设置为 1，打开所有锁存器使其处于透传状态。

所以，根据真值表生成逻辑电路的基本规律是：忽略输出值为 0 的行，找出输出值为 1 的行，然后观察该行的所有输入信号的值，若为 0，则对信号取反，然后将该行所有输入信号相与，再将所有行相或，即可得出该行输出值。每一行输入信号的正值或者反值相与形成一个乘积项 (Product-Term)，多个乘积项相或后形成一个输出值。



## 第2章

### ▲ 2.1 从累积计算说起

#### 2.1.1 实现累积计算

### ▲ 2.2 自动执行

#### 2.2.1 将操作方式的描述转化为指令

#### 2.2.2 实现那只章鱼——控制通路及部件

#### 2.2.3 动起来吧！——时序通路及部件

#### 2.2.4 半自动执行！——你得推着它跑

#### 2.2.5 全自动受控执行！——不用扬鞭自奋蹄！

#### 2.2.6 NOOP指令

#### 2.2.7 利用边沿型触发器搭建电路

#### 2.2.8 判断和跳转

#### 2.2.9 再见，Mr.章鱼！

### ▲ 2.3 更高效的执行程序

#### 2.3.1 利用循环缩减程序尺寸

#### 2.3.2 实现更多方便的指令

#### 2.3.3 多时钟周期指令

#### 2.3.4 微指令和微码

#### 2.3.5 全局地址空间

#### 2.3.6 多端口存储器

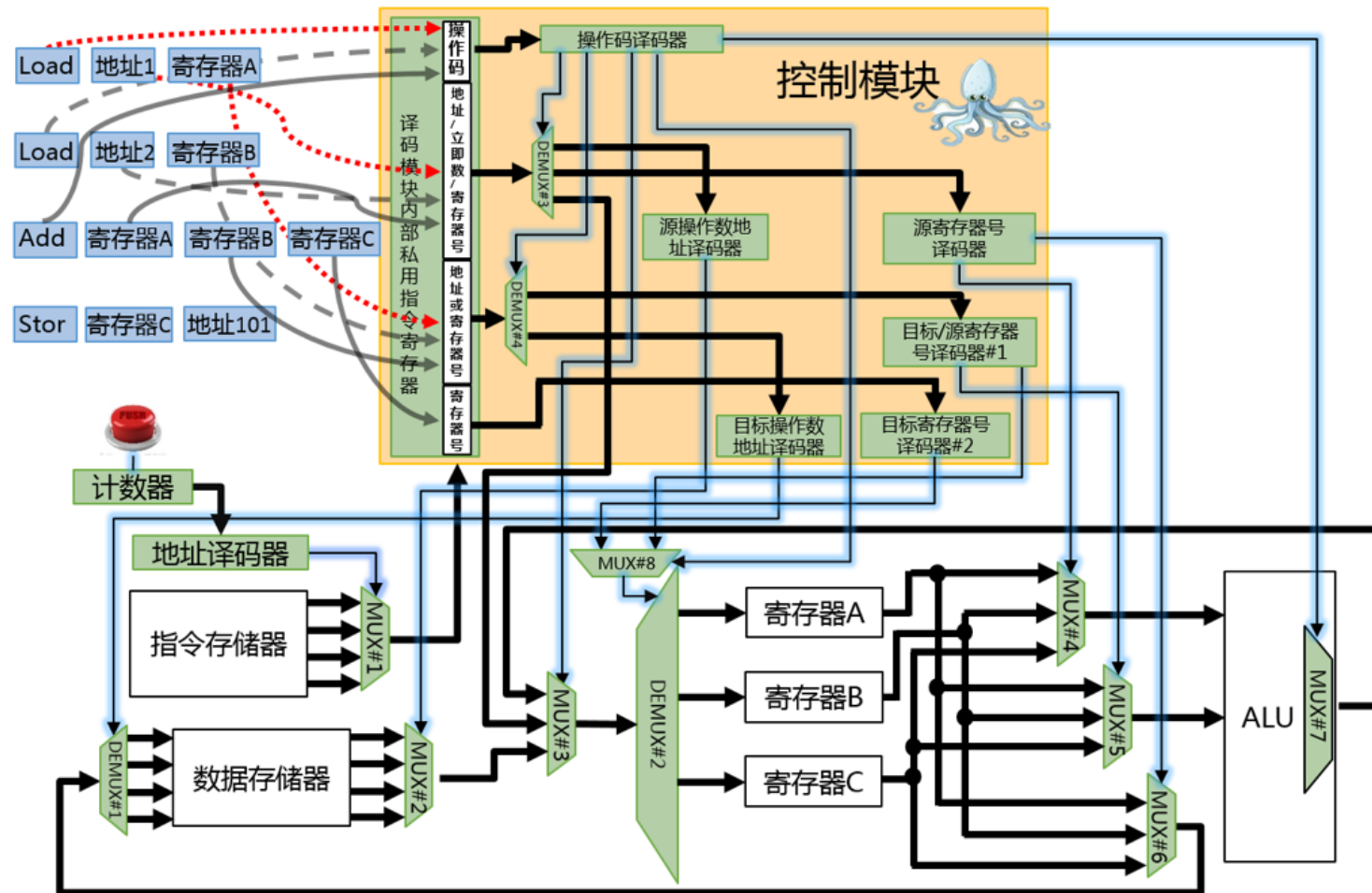
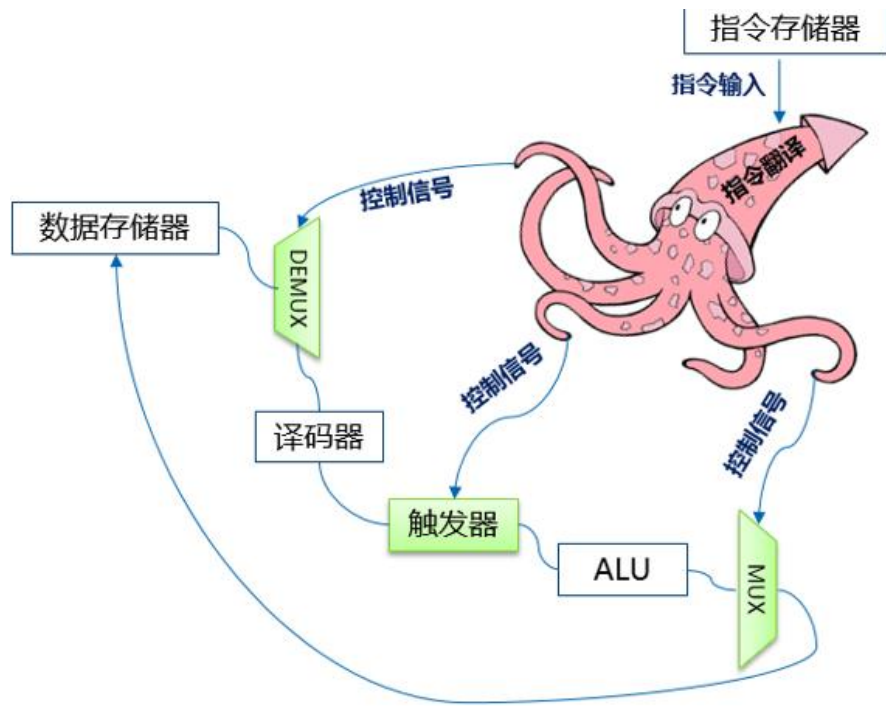
#### 2.3.7 多级缓存与CPU

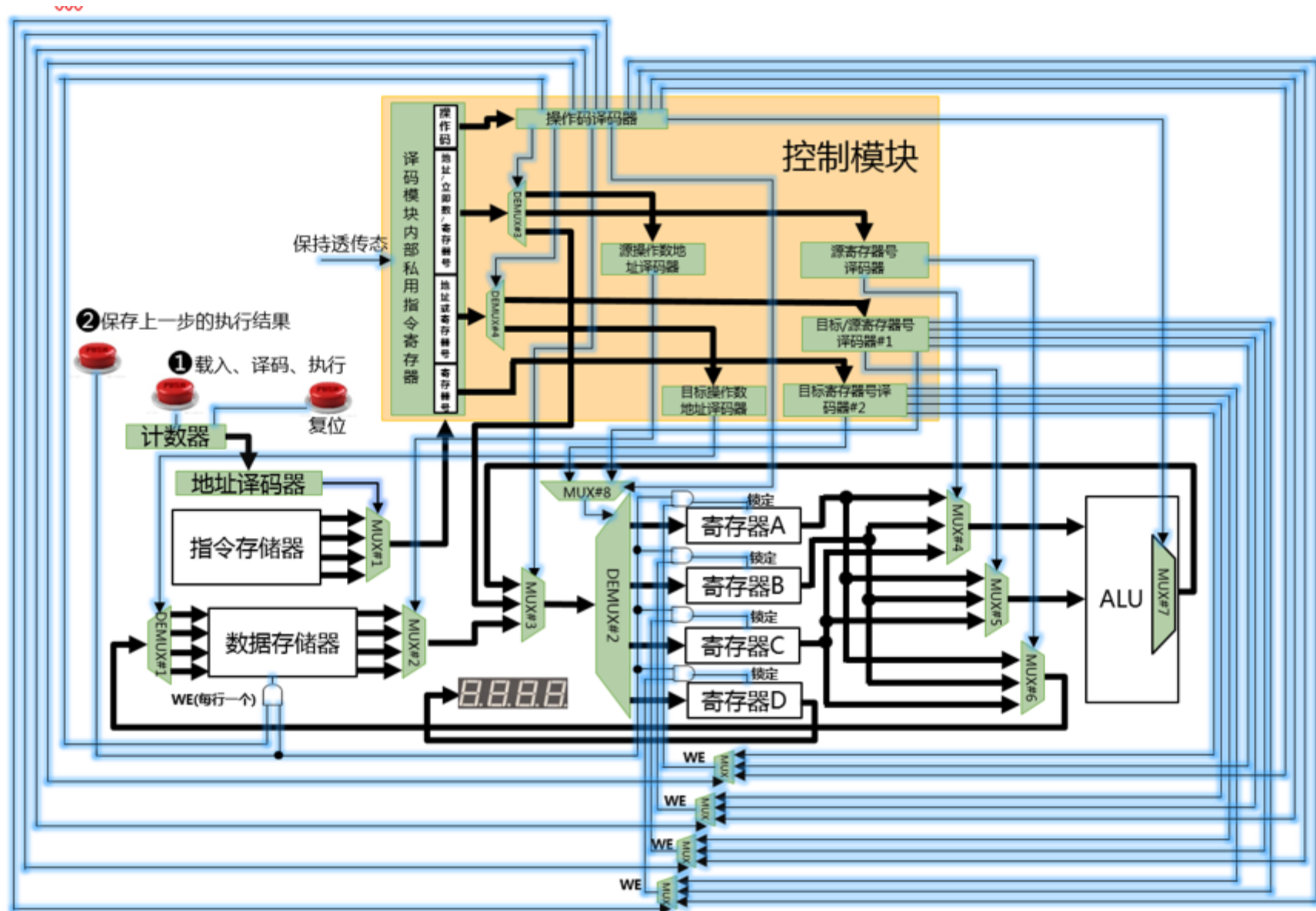
#### 2.3.8 数据遍布各处

#### 2.3.9 降低数据操作粒度

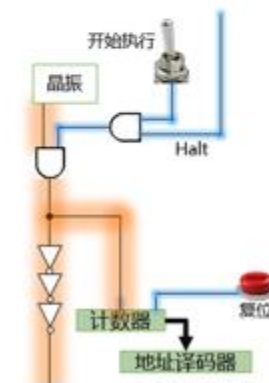
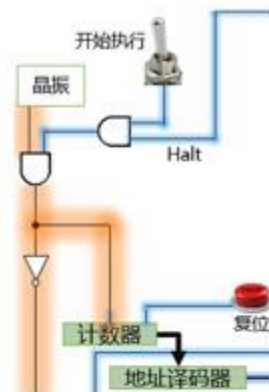
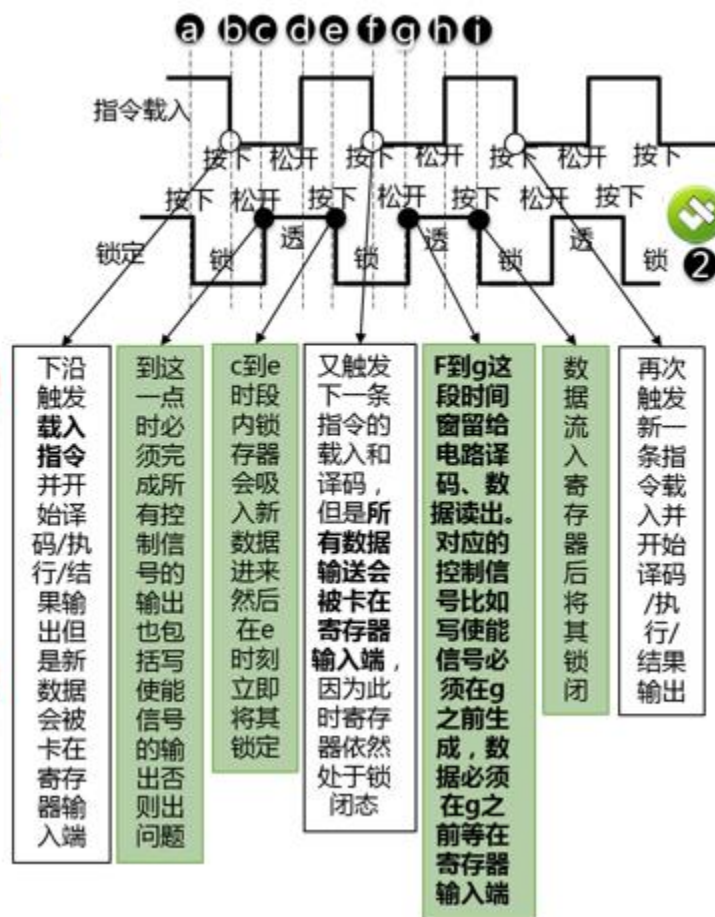
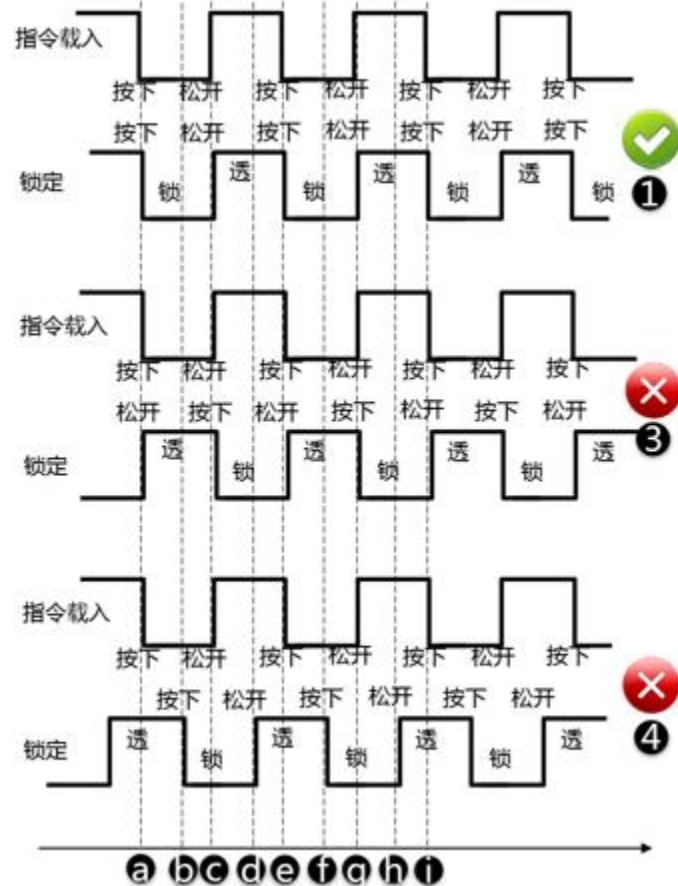
#### 2.3.10 取指令/数据缓冲加速

写作思路：在之前的简易计算器基础上，逐步实现自动根据指令来执行计算。比如算全班同学平均成绩，将800多次按键转化成指令，从而引出指令译码、控制逻辑、形成一个简易的CPU。









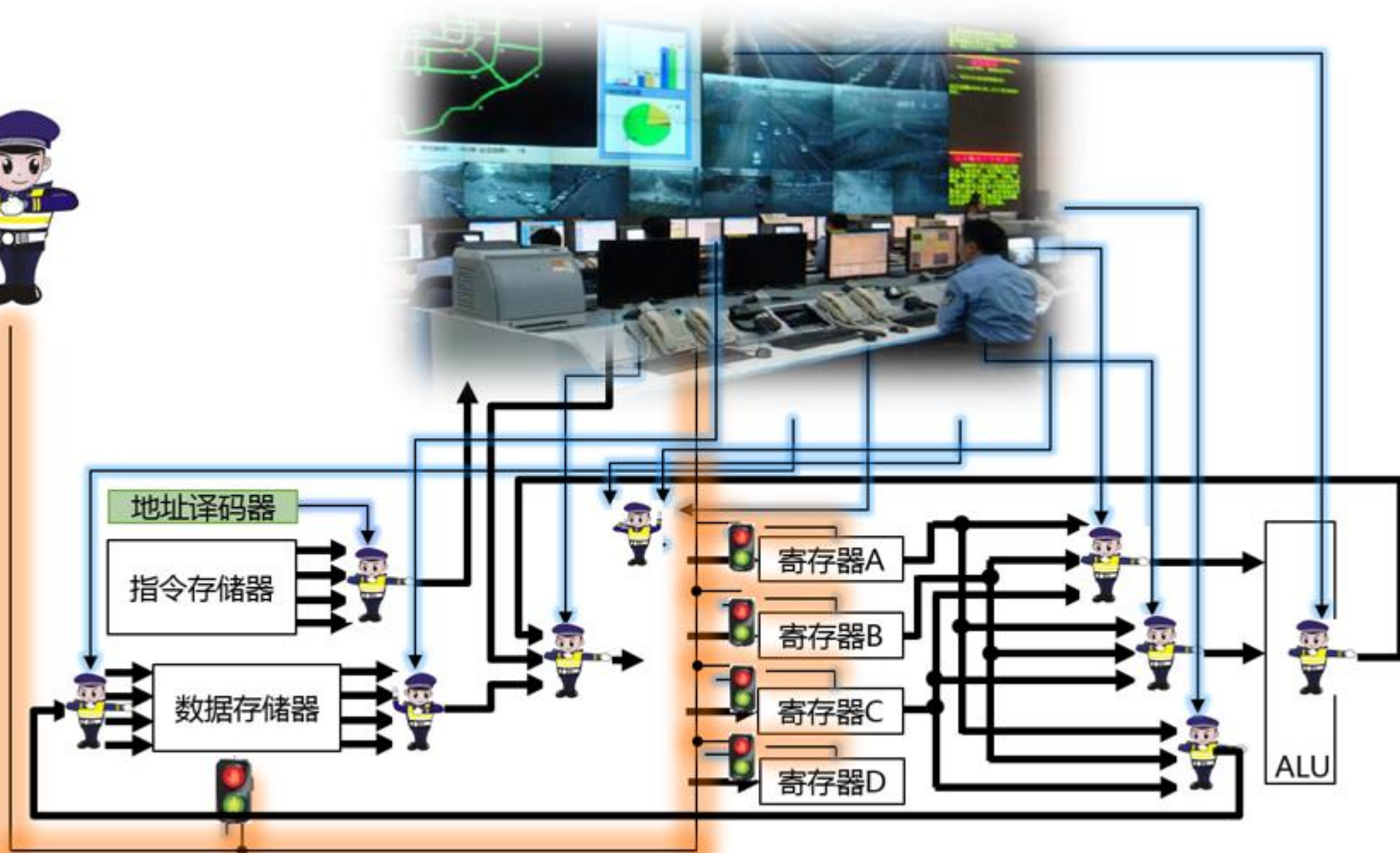


图 2-16 各控制部件只是在根据指令来指挥交通

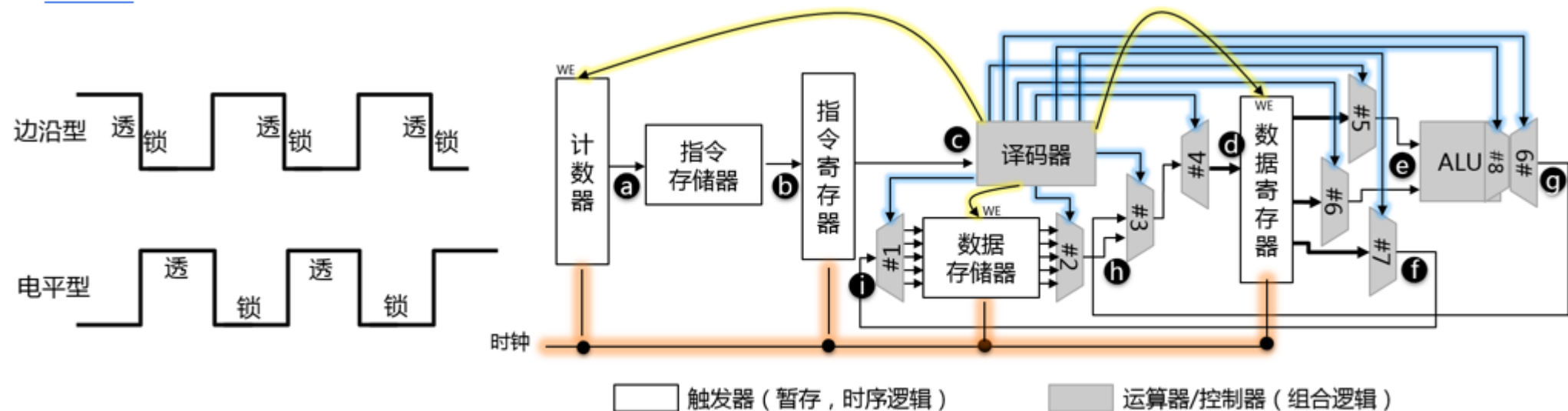


图 2-17 采用边沿型触发器时示意图

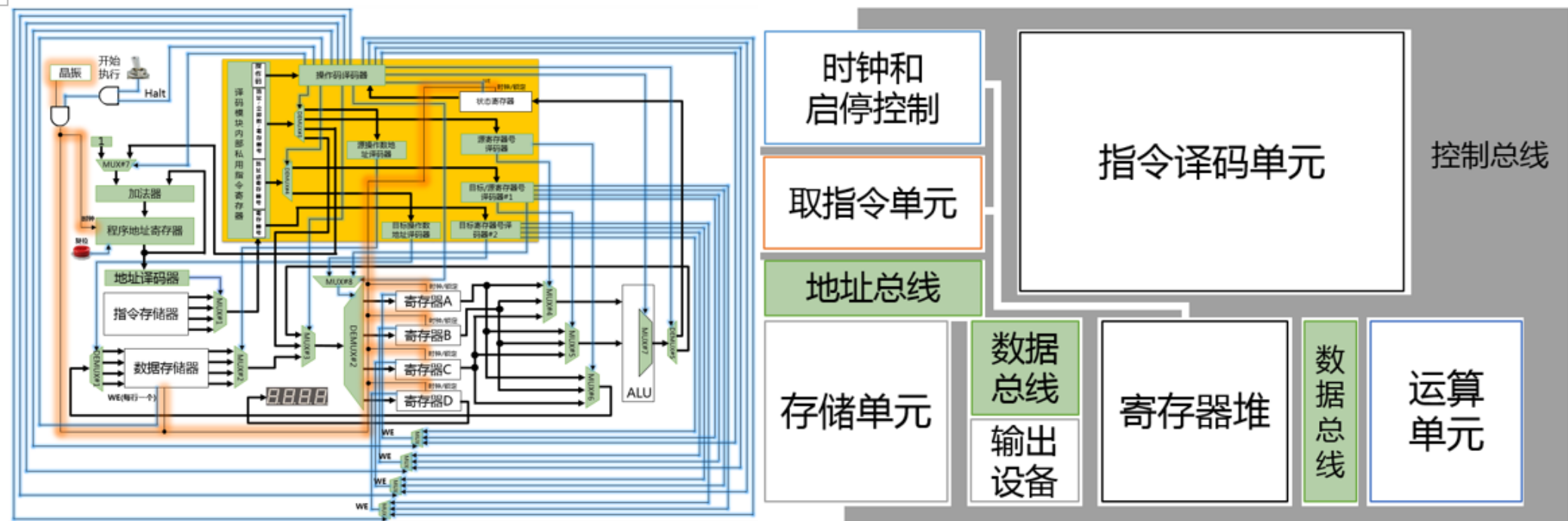


图 2-22 为主要模块起个名字



第3章	3.3 制造工艺革命——集成电路	3.3.8 半导体工艺的瓶颈	3.4.2 电子存储器
3.1 从薄铁片到机械计算机	3.3.1 从泥活字看量产晶体管	3.3.8.1 寄生电容	3.4.2.1 静态随机存储器SRAM
3.1.1 算盘和计算尺	3.3.2 跟冬瓜哥学做P/N结蛋糕	3.3.8.2 静态/动态功耗	3.4.2.2 动态随机存储器DRAM
3.1.2 不可编程手动机械十进制计算机	3.3.3 提升集成度	3.3.8.3 栅氧厚度和High-K材料	3.4.2.3 Flash闪存
3.1.3 可编程自动机械十进制计算机	3.3.4 芯片内的深邃世界	3.3.8.4 导线连接和Low-K材料	3.4.2.4 只读存储器ROM
3.1.4 可编程自动电动机械二进制计算机	3.3.5 cMOS集成电路工艺概述	3.3.8.5 驱动能力及延	3.4.3 光存储器
3.1.5 可编程自动全电动二进制计算机	3.3.6 cMOS工艺步骤概述	3.3.8.6 时钟树	3.4.3.1 光盘是如何存储数据的
3.2 电子管时代	3.3.7 cMOS工艺详细步骤	3.3.9 集成电路计算机	3.4.3.2 压盘与刻盘的区别
3.2.1 二极电子管	3.3.7.1 热氧化	3.3.10 微处理器计算机	3.4.3.3 光盘表面微观结构
3.2.2 三极电子管	3.3.7.2 氮化硅积淀	3.3.11 暴力拆解奔三CPU	3.4.3.4 多层记录
3.2.3 AM广播革命	3.3.7.3 浅槽隔离蚀刻	3.4 存储器-不得不说的故事	3.4.3.5 激光头的秘密
3.2.4 电子管计算机	3.3.7.4 pMOS和nMOS生成	3.4.1 机械存储器	3.4.3.6 蓝光光盘简介
3.2.5 石头会唱歌	3.3.7.5 触点电极的生成	3.4.1.1 声波/扭力波延迟线 ( Delay Line	3.4.4 不同器件担任不同角色
3.3 固态革命——晶体管	3.3.7.6 通孔和支撑柱 ( via ) 的生成	3.4.1.2 磁鼓存储器 ( Drum )	3.4.4.1 寄存器和缓存
3.3.1 P/N结与晶体管	3.3.7.7 第一层导线连接	3.4.1.3 磁芯存储器 ( Core )	3.4.4.2 主运行内存/主存
3.3.2 场效应管 ( FET )	3.3.7.8 第二层导线连接		3.4.4.3 Scratchpad RAM
3.3.3 MOSFET	3.3.7.9 表面钝化		3.4.2.4 内容寻址内存CAM/TCAM
3.3.4 cMOS			3.4.2.5 外存
3.3.4 晶体管计算机			

写作思路：前两章读者已经充分了解了数字电路是如何进行计算的。本章介绍人们是如何用电子开关搭建出电路的，电子开关的发展，一直到最后的芯片是怎么做出来的。

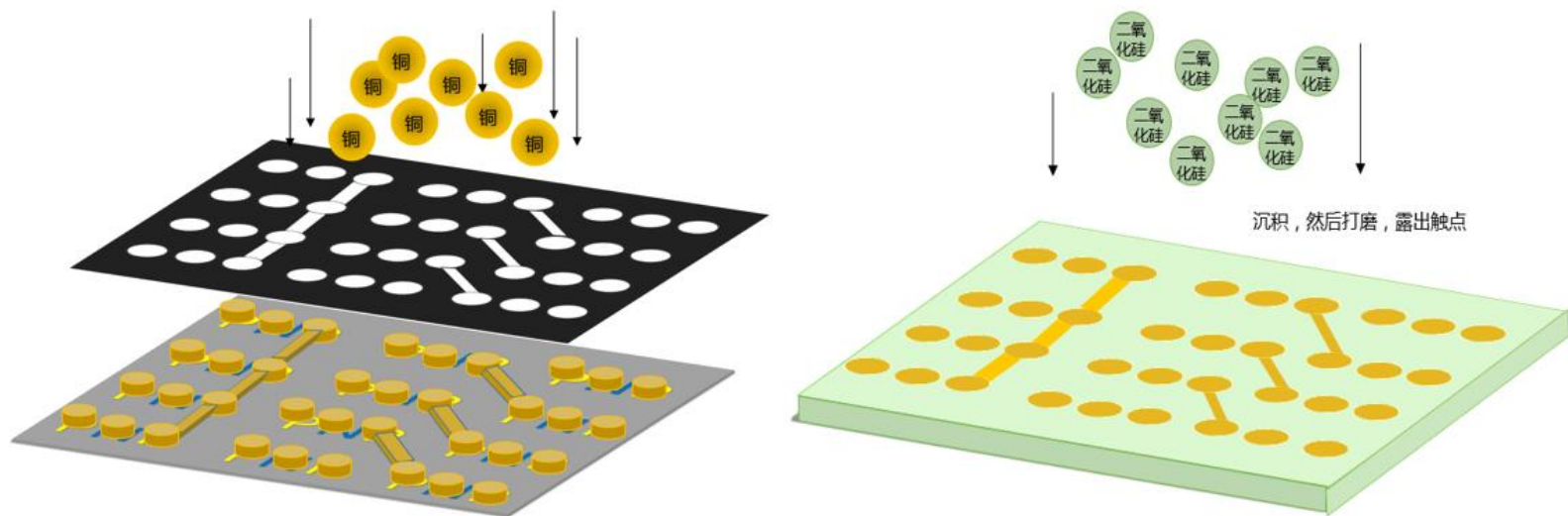


图 3-101 沉积导线以及覆盖一层二氧化硅绝缘层

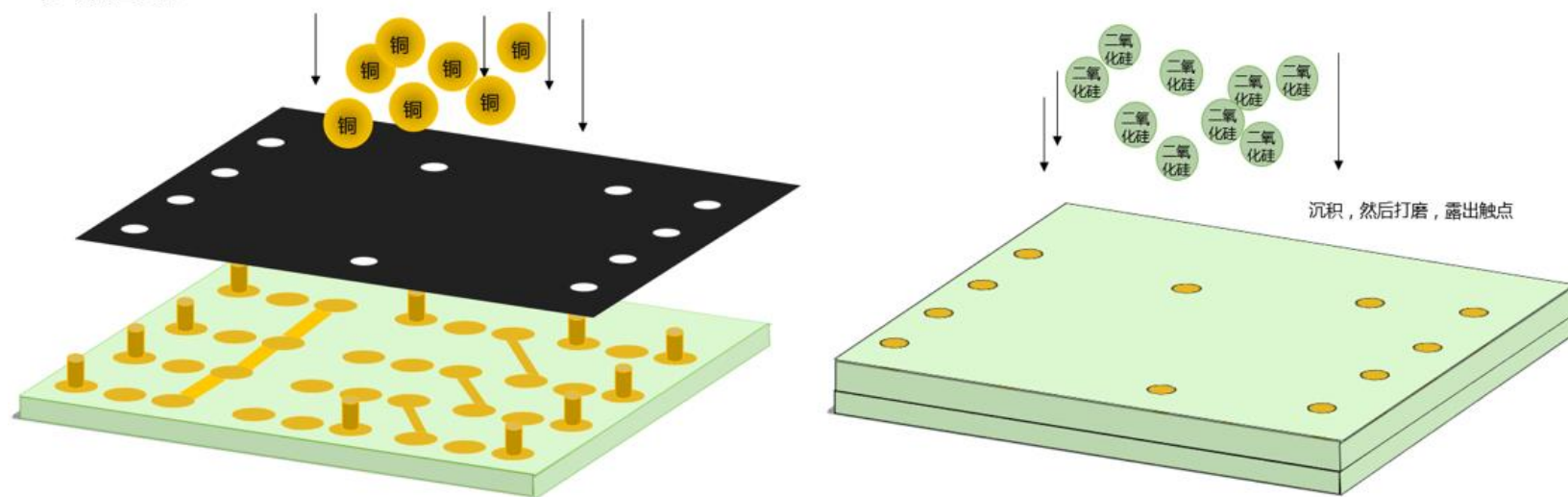
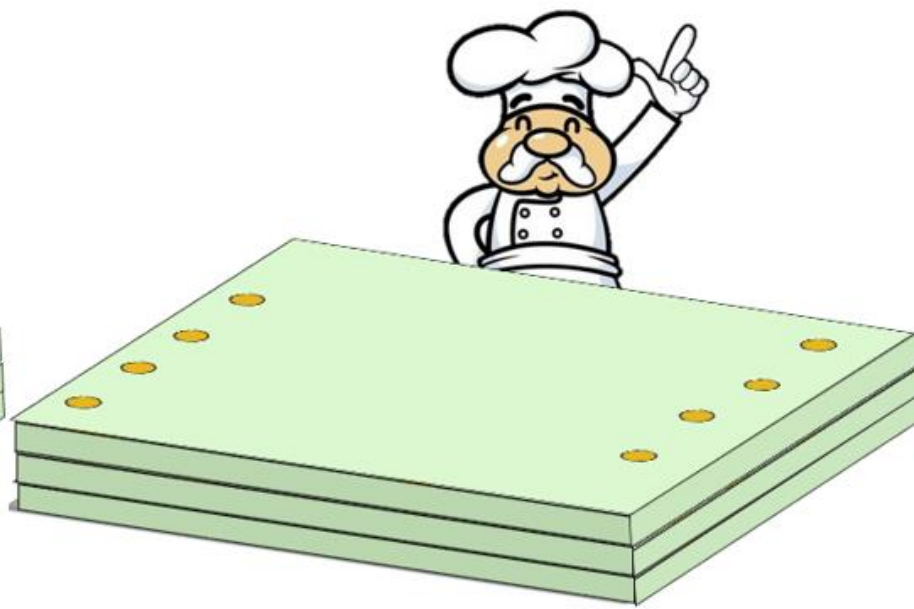
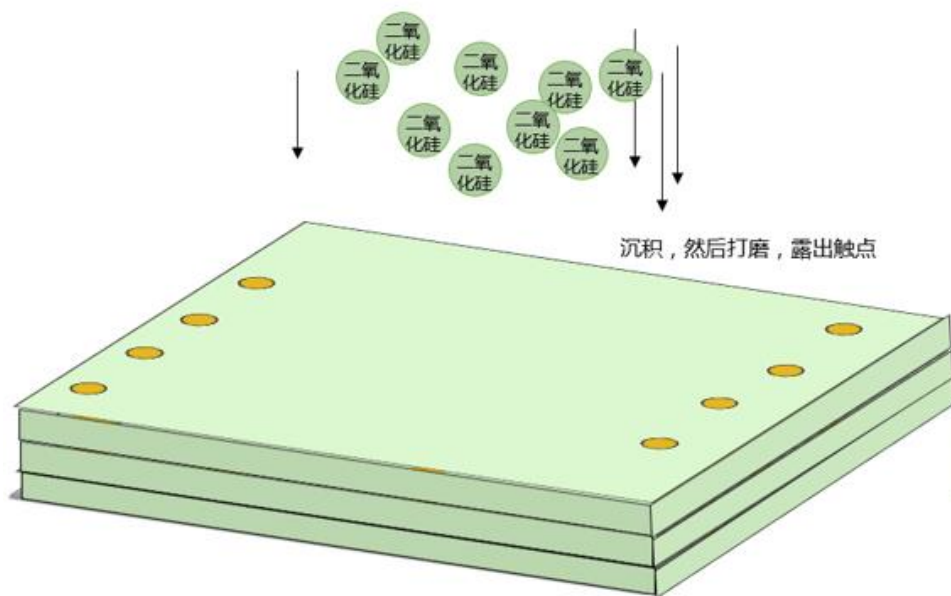
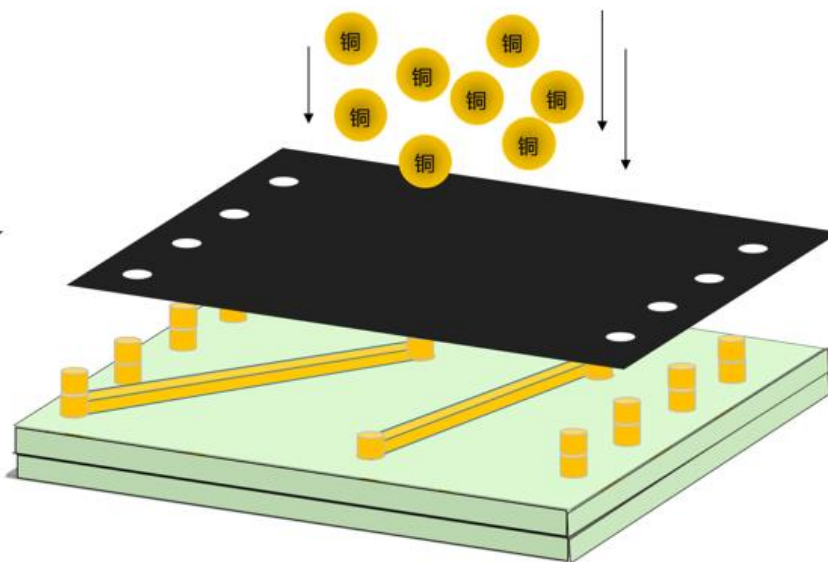
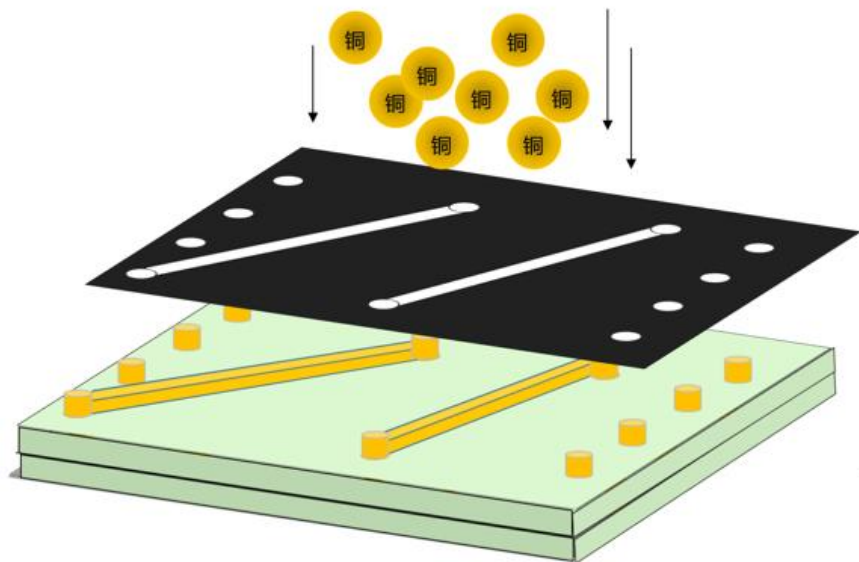


图 3-102 将需要高架的导线两端升起来并再次覆盖绝缘层





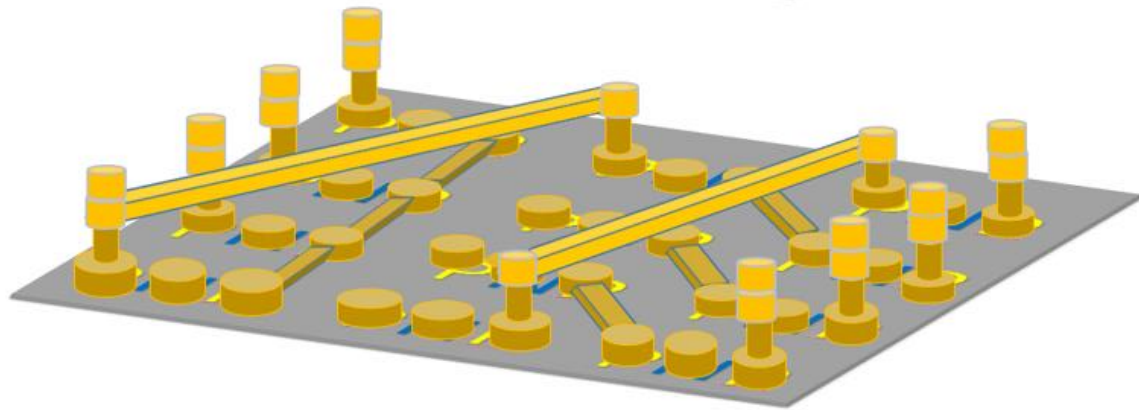
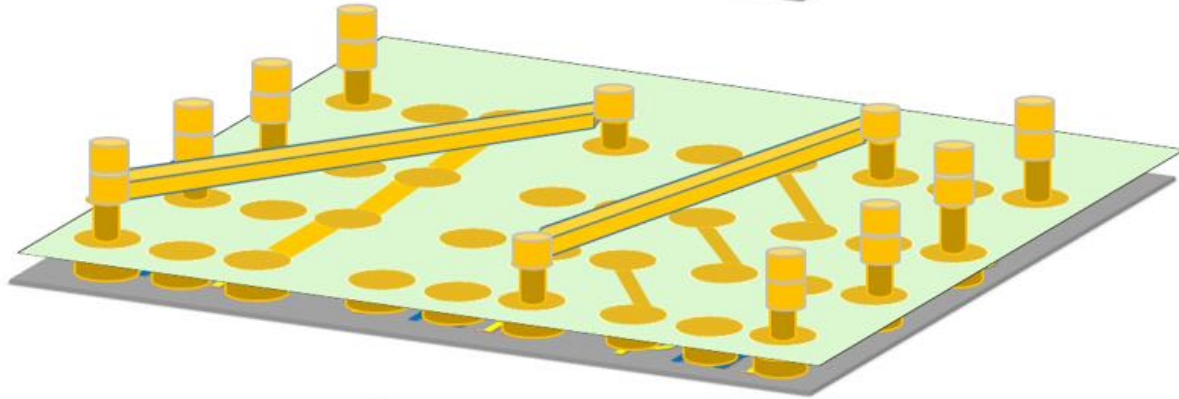
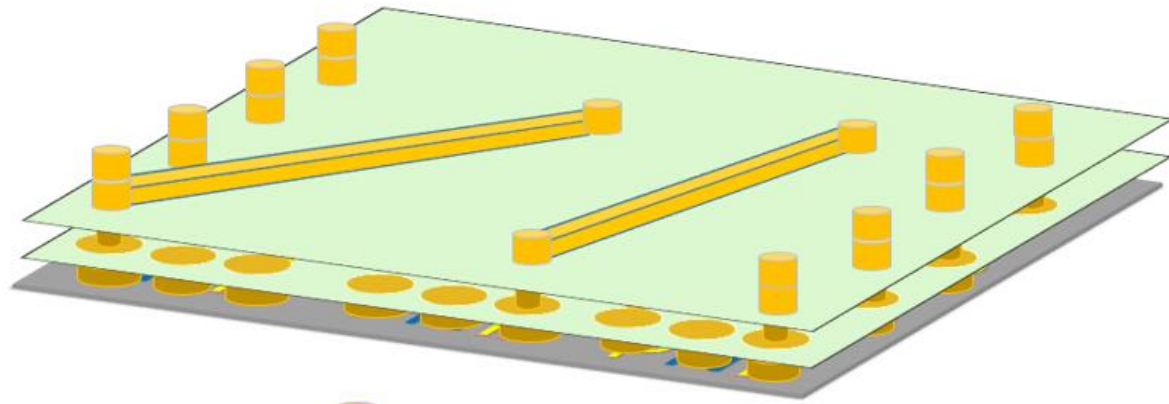




图 3-112 60 年代的人们正在红色胶片上手工雕刻遮罩

## ▲ 4.1 大话流水线

### 4.1.1 不高兴的译码器

## ▲ 4.1.2 思索流水线

### 4.1.2.1 流水线的本质是并发

### 4.1.2.2 不同时延的步骤混杂

### 4.1.2.3 大话队列

### 4.1.2.4 流水线的应用及优化

## ▲ 4.2 优化流水线

### 4.2.1 拆分慢速步骤

### 4.2.2 放置多份慢速部分

### 4.2.3 加入缓冲队列

### 4.2.4 图解五级流水线指令执行过程

## ▲ 4.3 流水线冒险

### 4.3.1 访问冲突与流水线阻塞

### 4.3.2 数据依赖与数据前递

### 4.3.3 跳转冒险与分支预测

## ▲ 4.4 指令的动态调度

### 4.4.1 结构相关与寄存器重命名

### 4.4.2 保留站与乱序执行

### 4.4.3 分步图解乱序执行

### 4.4.4 重排序缓冲与指令顺序提交

## ▲ 4.5 物理并行执行

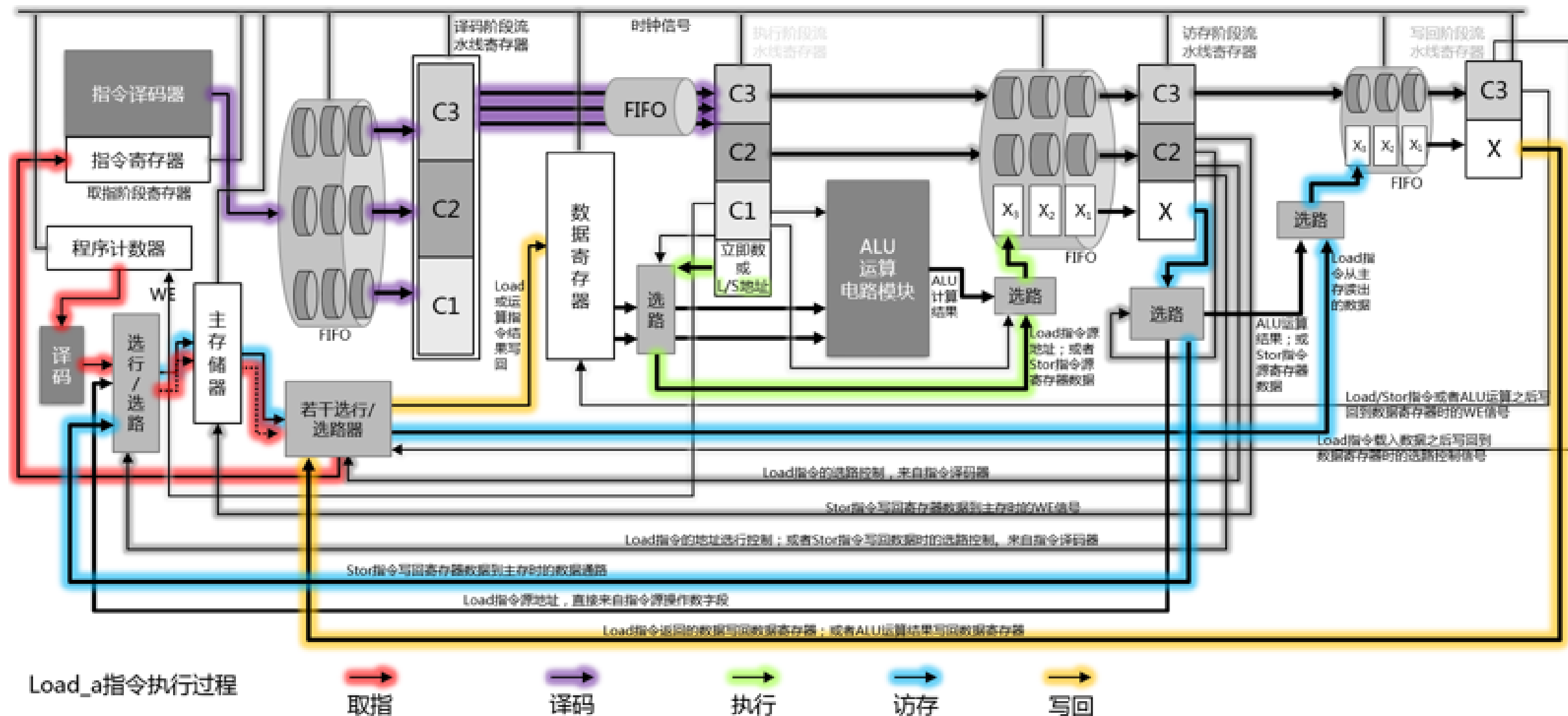
### 4.5.1 超标量和多发射

### 4.5.2 VLIW超长指令字

### 4.5.3 SIMD单指令多数据

写作思路：前三章读者已经充分了解CPU基本原理，以及人们如何制造CPU，通过芯片技术可以集成入大量复杂的电路逻辑。有了这个认知基础，开始介绍CPU内部更高级的控制方法，流水线、多发射、乱序执行、分支预测等等体系结构方面的概念。以“不高兴的译码器”为标题党，强烈将读者思维代入作者设下的套中无法自拔，并最终跟随作者的思考思路完成本章的学习。





从这 8bit 记录中，可以观察到该指令以往 8 次跳转的**关联模式/范式**，比如范式 01010101 表示隔一次跳一次；而范式 00110011 表示隔两次跳两次，同理，01100111 也是一种范式，但是看上去没什么明显规律，也没法用简洁的语言表达，但是如果 01100111 这种范式在程序执行期间重复出现的话，也就能从更久远的历史中看出之前根本无法分辨的规律。至于为什么体现出这种规律，那是更底层的计算机程序科学家所研究的问题了。对于分支预测，很显然，如果上一次该代码处于某个范式时，它跳了，那么当它执行了若干次后，如果恰好这 8bit 历史记录又轮回到了同一个范式，那么你说，该跳还是不跳？当然是跳啊！为什么？因为上一次这样的時候它跳了啊！就这么简单。用这种方式可以深入到事物的更深层次的表象中挖掘出更深层的规律。相比之下，如果只是简单的比较 1 和 0 谁多就听谁的，就显得简陋和粗暴了。如图 4-44 中所示的 BHT 表中，其中红、黄、蓝、绿四条指令的跳转范式就出现了轮回，注意同一行中的相同颜色字体表示范式轮回。

好了，看来这个思路已经非常明朗了，现在需要解决如何记录“当某个范式到来时该指令跳了没有”这件事了。请注意，之前的思路是为每一条跳转指令记录一个 2bit 的跳转记录，而现在需要为每一条指令的多个跳转范式分别各自记录一个 2bit 的跳转记录，有多少个可能的范式，就要记录多少个 2bit 的跳转记录，8bit 可以组成 256 种范式，所以需要记录 256 个 2bit。怎么记录？能否直接在 BHT 表的右侧加一列跳转记录列呢？不行，因为我们需要对单个指令的每一可能的范式都记录对应的跳转记录，8bit 可以组成最大 256 个范式，那就意味着，对于单条指令，我们需要建立一张 256 行 2 列的表，第一列 8bit 记录该指令所有可能的 256 种范式，第二列 2bit（每次遇到该范式时跳了就置 1，不跳就置 0，也用 2bit 来记录，增加防抖动缓冲）记录该范式对应的跳转记录，表格共计  $256 \times 10 \text{bit} = 320 \text{Byte}$ 。这还是单条指令，如果记录 256 条指令的话，那么就需要  $256 \times 320 \text{Byte} = 80 \text{KB}$  的 CAM 存储器，会占用较大电路面积。

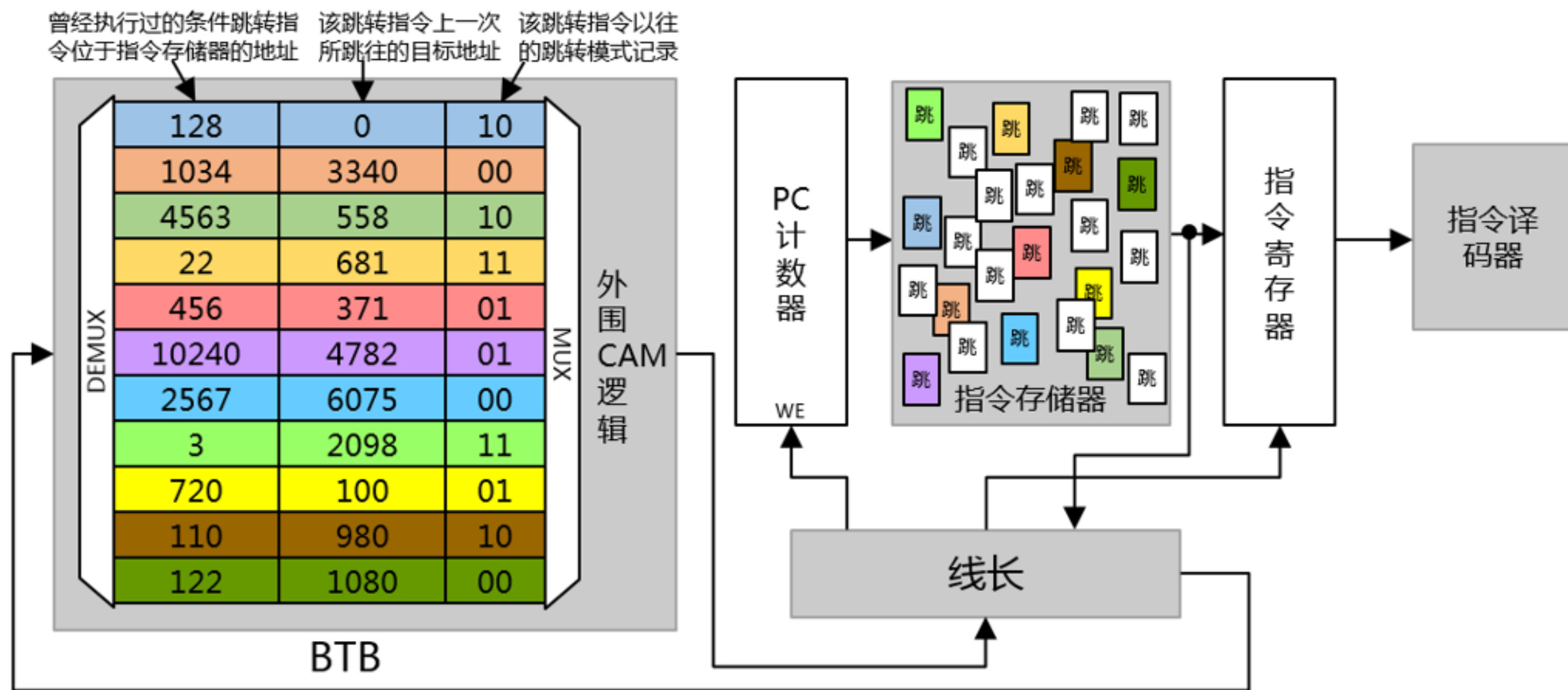


图 4-43 加入跳转目标地址列的 PHT 被称为 BTB

6 Add A D C : //算出C+E的值并存储到寄存器E中





▲ 第5章	▲ 5.4 程控多媒体计算机	▲ 5.5.5 程序的执行和退出
▲ 5.1 基本的数据结构	5.4.1 键盘是前提	5.5.5.1 初步解决地址问题
5.1.1 数组	5.4.2 搜索并显示	5.5.5.3 更好的人机交互方式
5.1.2 数据类型与ASCII码	5.4.3 实现简易计算器	5.5.5.4 程序的退出
5.1.3 结构体	5.4.4 录入和保存	5.5.5.5 使用外部设备和内存
5.1.4 数据怎么摆很重要	5.4.5 简易文件系统	▲ 5.5.6 多程序并发执行
▲ 5.2 高级语言	5.4.6 计时/定时	5.5.6.1 利用时钟中断来切换线程
5.2.1 简单的声明和赋值	5.4.7 发声控制	5.5.6.2 更广泛的使用中断
5.2.2 编译和编译器	5.4.8 图像显示	5.5.6.3 虚拟地址空间与分页
5.2.3 向编译器描述数据的编排方式	5.4.9 网络聊天	5.5.6.4 虚拟与现实的边界——系统调用
5.2.4 高级语言编程小试牛刀	▲ 5.5 程序社会	5.5.7 呼唤操作系统
5.2.5 人脑编译忆苦思甜	5.5.1 函数和调用	▲ 5.6 计算机操作系统
▲ 5.3 浮点数及浮点运算	5.5.2 设备驱动程序	▲ 5.6.1 内存管理
5.3.1 数值范围和精度	5.5.3 函数之间的联络站	5.6.1.1 分区式内存管理
5.3.2 浮点数的用处和表示方法	▲ 5.5.4 库和链接	5.6.1.2 分段+分区式内存管理
▲ 5.3.3 浮点数的二进制表示	5.5.4.1 静态库和静态链接	5.6.1.3 Intel CPU的分段支持机制
5.3.3.1 二进制浮点数转十进制小数	5.5.4.2 头文件	5.6.1.4 分页+分段式内存管理
5.3.3.2 十进制小数转二进制浮点数	5.5.4.3 API和SDK	5.6.1.5 DOS和Linux操作系统的内存管理
5.3.3.3 负指数和0的表示	5.5.4.4 动态库和动态链接	5.6.2 中断响应及处理
5.3.3.4 无穷与非规格化数的表示	5.5.4.5 库文件/可执行文件的格式	5.6.3 驱动程序与设备管理
5.3.4 浮点数运算挺费劲		5.6.7 文件系统和I/O协议栈
5.3.5 浮点数的C语言声明		5.6.5 系统调用
5.3.6 十六进制表示法		5.6.4 线程调度与管理
		5.6.6 进程间通信
		5.6.8 人机交互界面

写作思路：是时候引入程序是如何执行的了。从基本数据结构，介绍到高级语言，然后是编译原理，并手动编译一个简易程序示例。都是带着思考，先给出命题，为什么会这样，解决思路，具体设计，实际设计。最终一步步引出操作系统的概念。



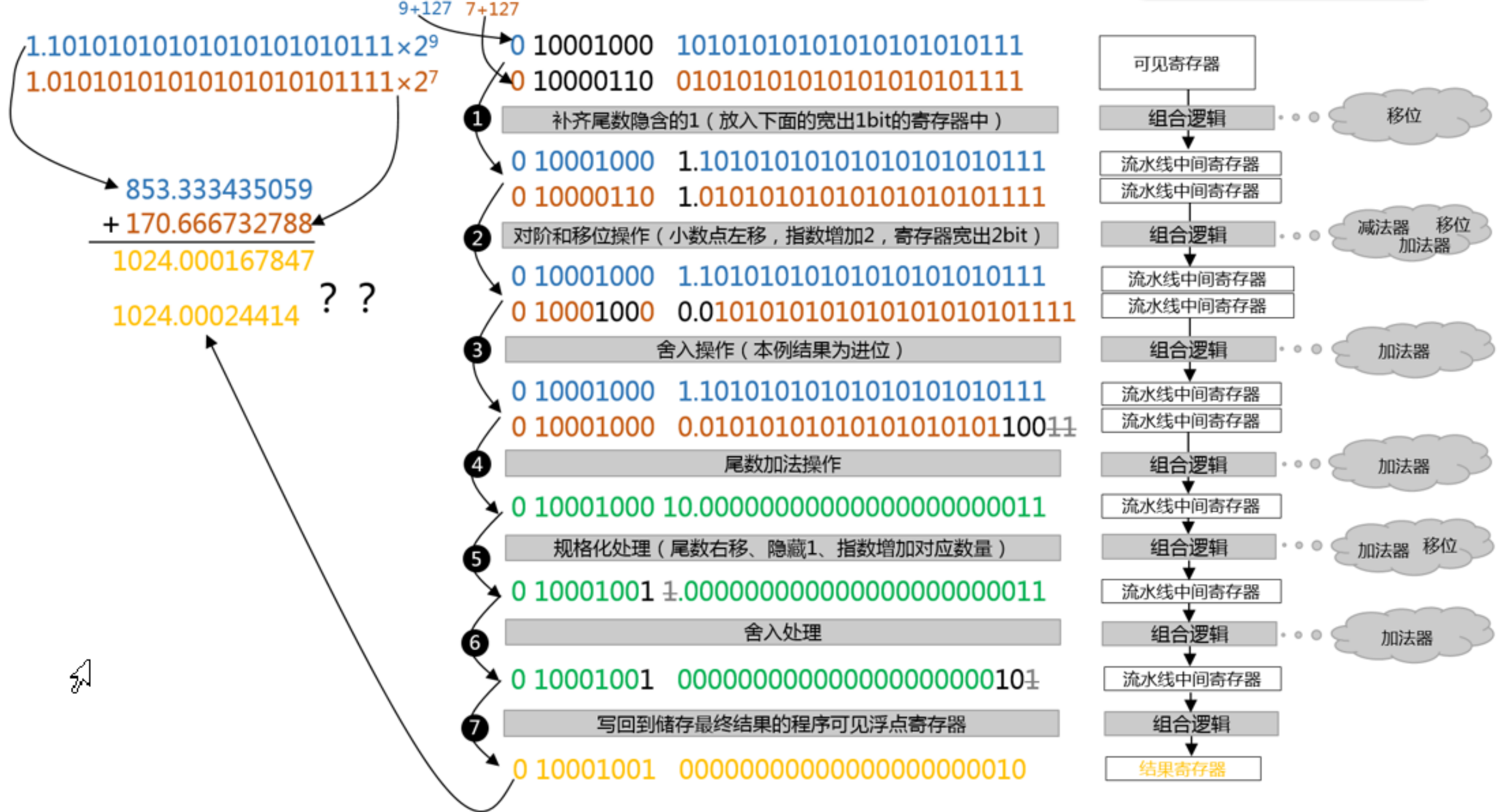
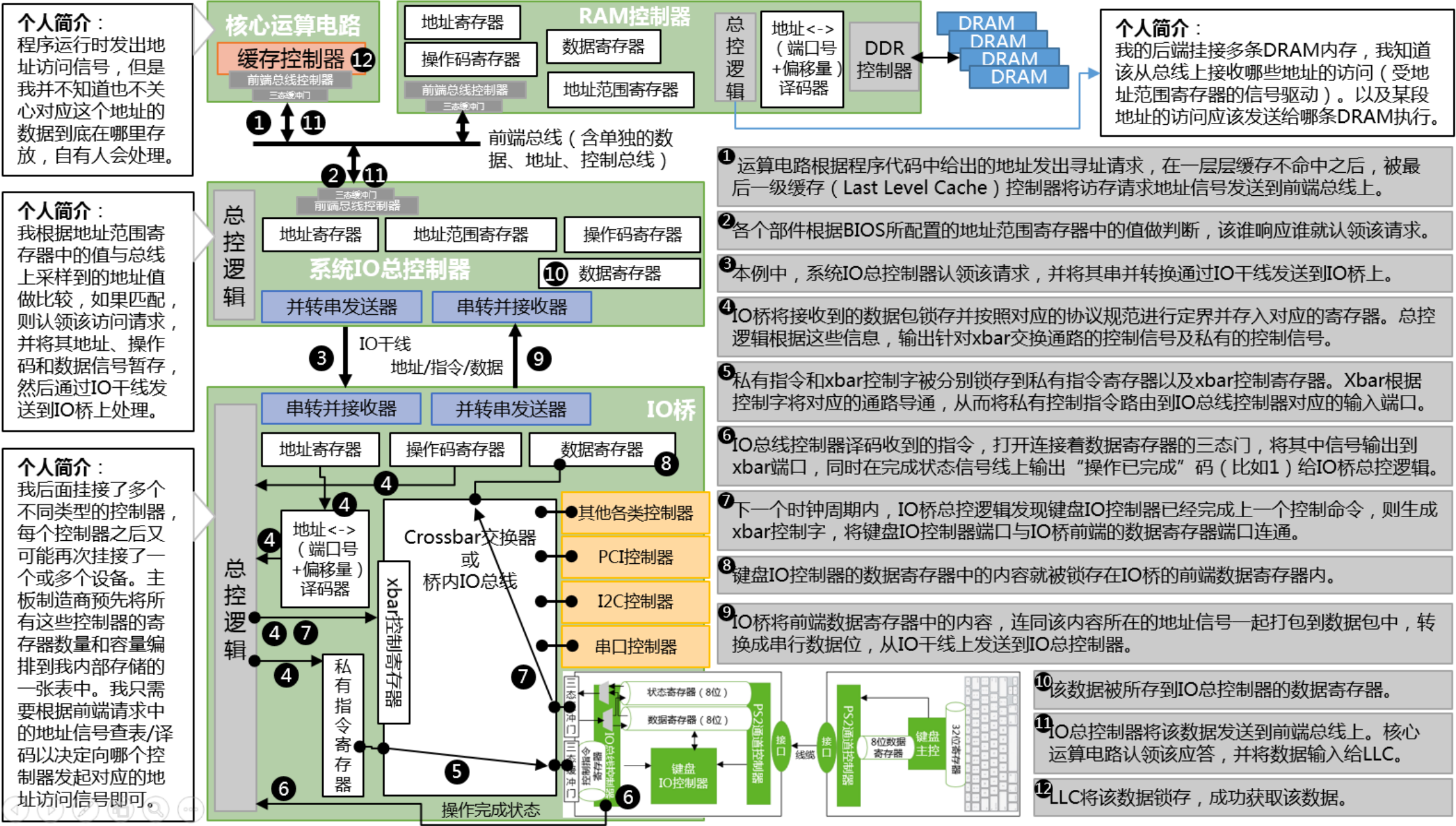
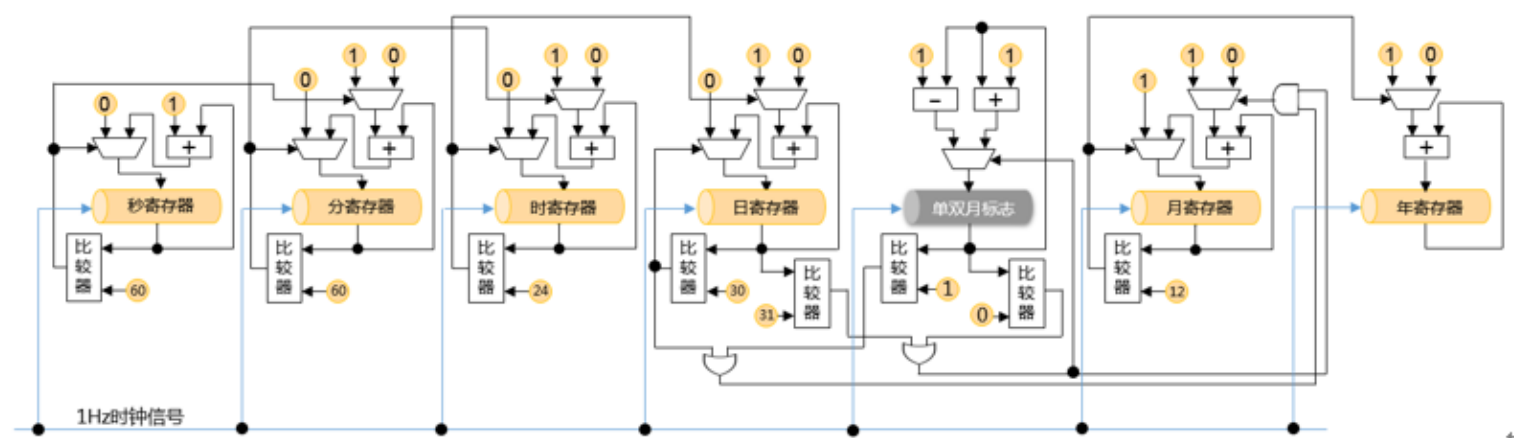


图 5-18 两个相同符号的浮点数相加的执行流程



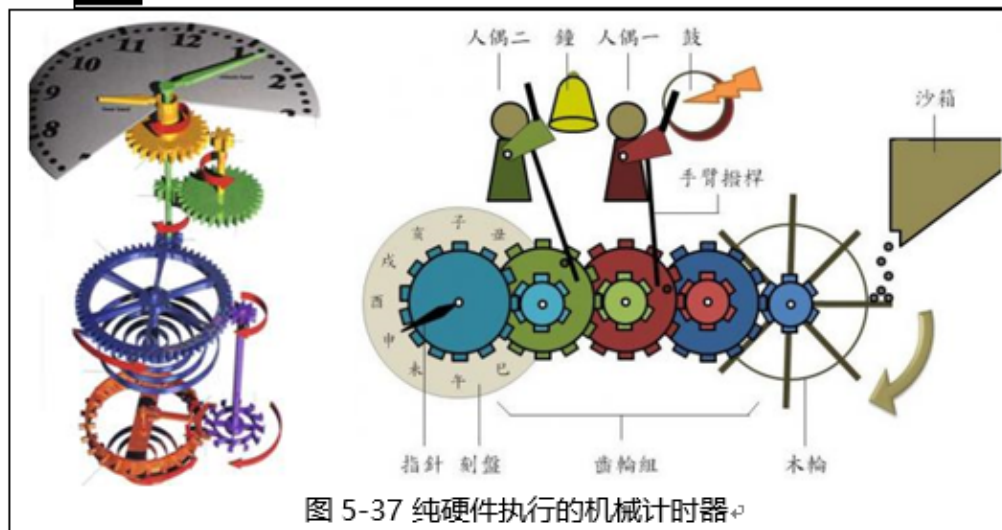




可以把这个纯数字逻辑计时器，辅之以一个 IO 控制器，然后与主 CPU 一侧的 IO 桥对接起来，也完全可以完成计时和时间读取的功能。↵

软和硬

怎么样，你现在应该彻底理解了什么叫“软件定义”、什么叫“硬件加速”了。所谓软件定义，就是利用通用 CPU，依靠一条条的机器指令来完成一个任务；而硬件加速，就是直接将这个任务中的每一步做成纯数字逻辑，不需要通用 CPU+代码，所有的逻辑、步骤都被固化到数字逻辑中，这样执行速度最快，效率最高，但是不具备灵活性，也就是“可编程”性，比如，一旦某个需求变化了，那么就需要重新设计数字逻辑，而并不能像软件定义那样直接改一改代码重新编译一下即可。冬瓜哥将在第 9 章中更详细的介绍硬件加速。



将图 5-36 与 5-37 中的机械计时器相比，是不是感觉到事物底层的相似性？最终，你又一次深刻领会到了 CPU+代码实现某个功能，与用纯数字逻辑实现之间的联系和区别。看到这里，你应当在脑海中自然回想起第 1 章的简易计算器，以及我们为了使之变得更灵活而设计的那个简易 CPU，并结合现在在计时器设计上我们反演了这个过程，如果能有一种奇妙的顿悟，那

read\_keyboard()函数会扫描该表中所有类型为“键盘”的设备从而让用户选择使用哪个键盘来输入，或者read\_keyboard()自行选择一个默认设备作为输入设备；刚才图 5-49 中的菜单，也是由程序扫描设备信息表中所有发声控制设备，从而让用户选择使用哪个设备输出声音。

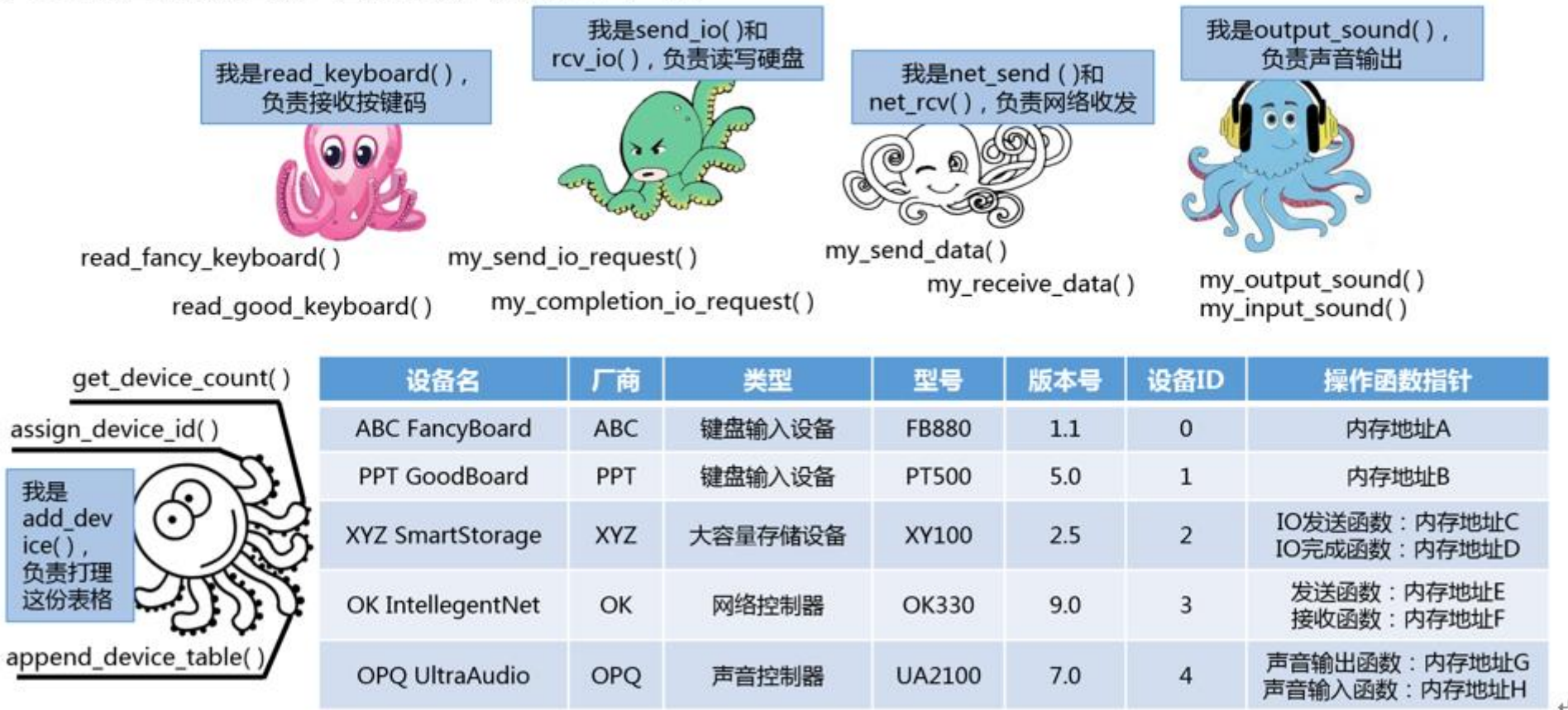
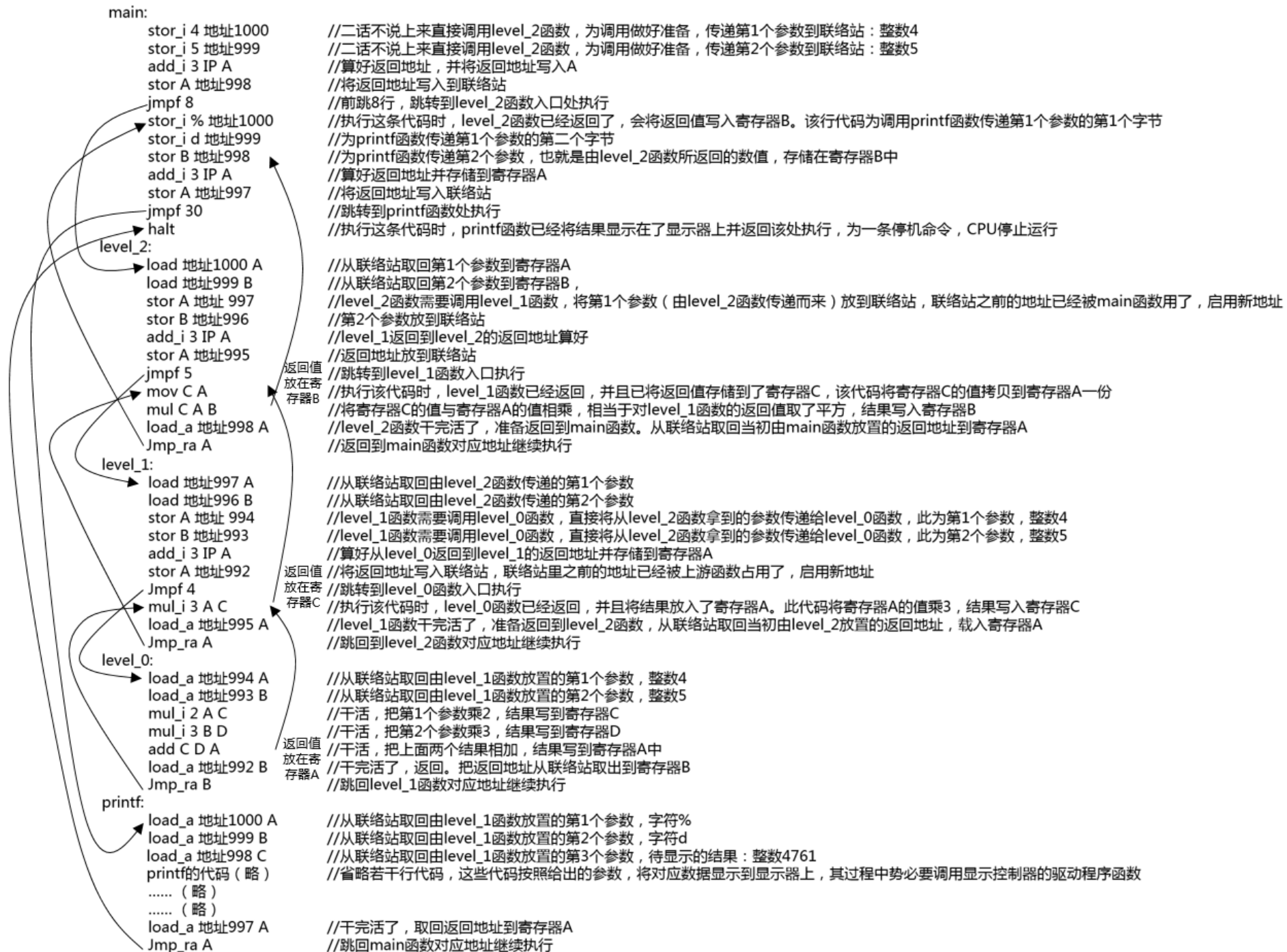


图 5-50 设备信息描述表





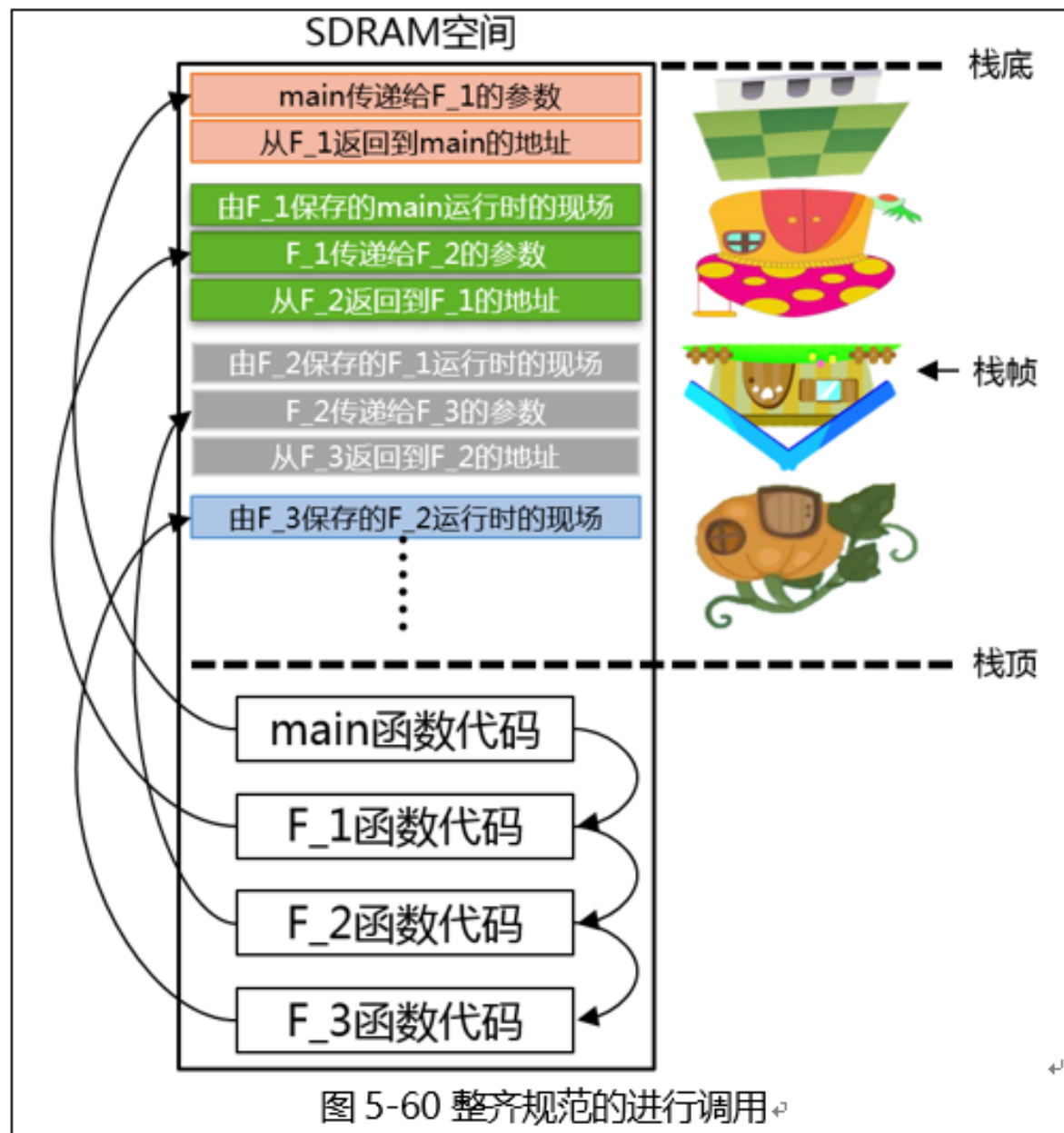


图 5-60 整齐规范的进行调用

基于这个思路，人们把整个联络站称为“栈 (Stack)”，就像客栈一样，给函数们提供一个歇脚和相互沟通的地方。而把栈中的每一间客房称为“栈帧 (Stack Frame)”，函数把自己要传达给下游的参数以及返回地址一行一行的放到栈帧里。

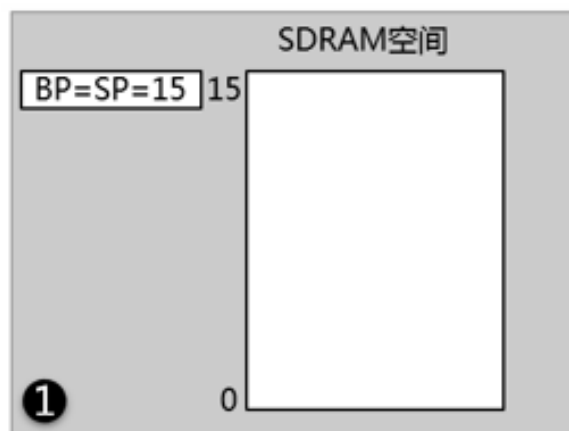
函数们来住店，首领 main 先住了进来，然后就开始派活给其他函数，这老哥们下达完指令就开始睡大觉了，他委派的人可能继续委托其他人，每个人委托其他人之后都开始睡大觉，然后等待下游主动执行跳转指令跳回来，然后醒来继续干活。可以看到，栈的空间是不断扩大的，自顶向下。程序员需要注意，如果级连调用的函数过多，或传递的参数过多，而导致栈空间大到触碰了 SDRAM 低地址区域的代码、数据，那么会导致出错。

另外，栈空间扩大时一定是一条条追加上去的。比如某函数被调用执行后，第一件事就是保存上一个函数的现场，将数个寄存器值追加

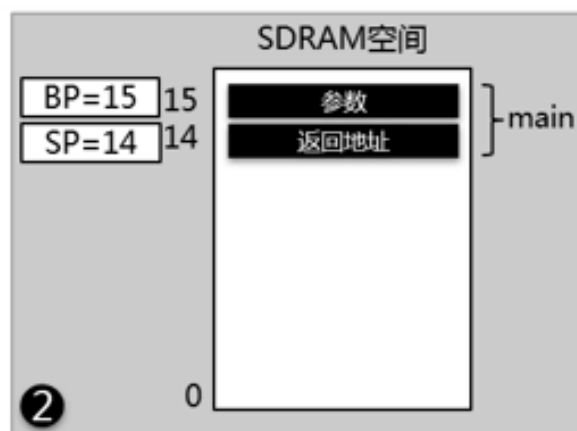
到栈顶部，当其需要调用其他函数时，再将数个参数追加写入栈顶部、返回地址写入栈顶部；如果其调用的函数继续向



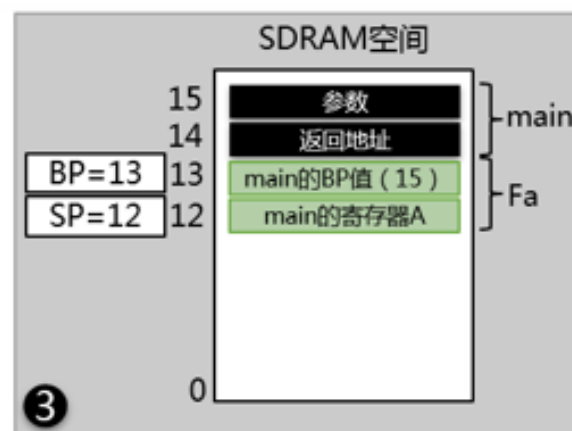
个过程示意图如图 5-61 所示。Main 调用 Fa，Fa 又调用 Fb，Fb 不再调用其他函数。



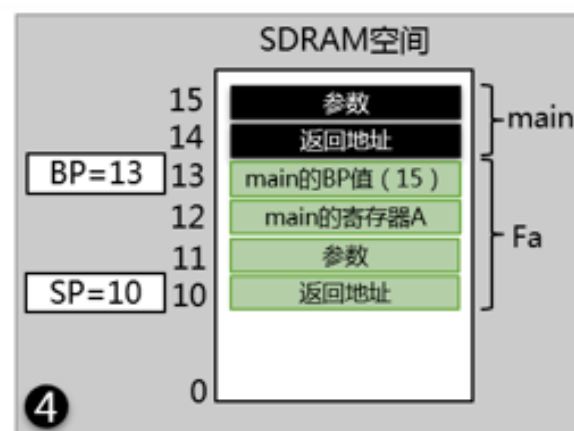
初始状态



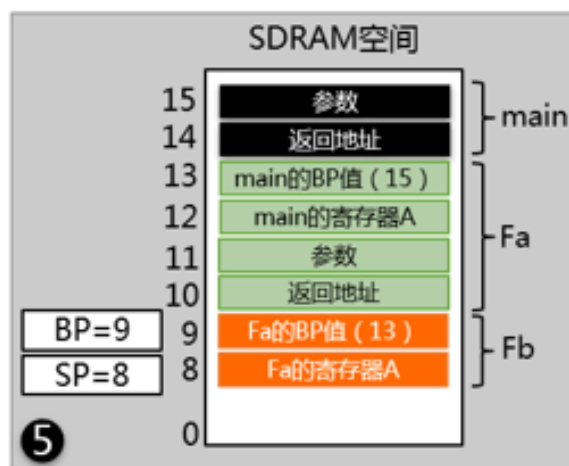
main准备调用下游函数



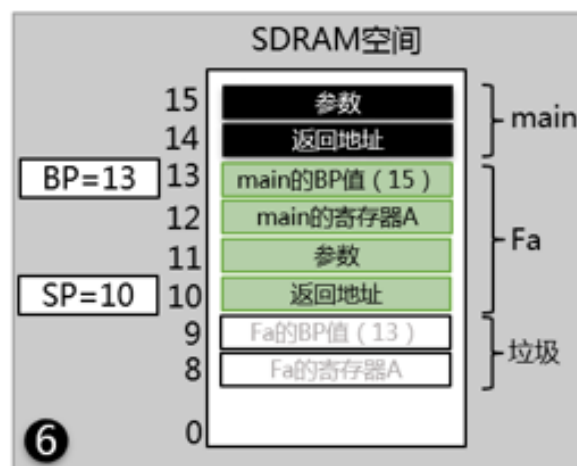
函数Fa开始执行，保护main的现场



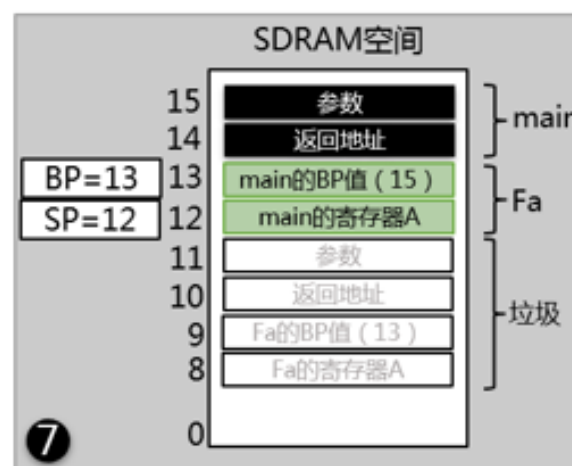
函数Fa准备调用Fb



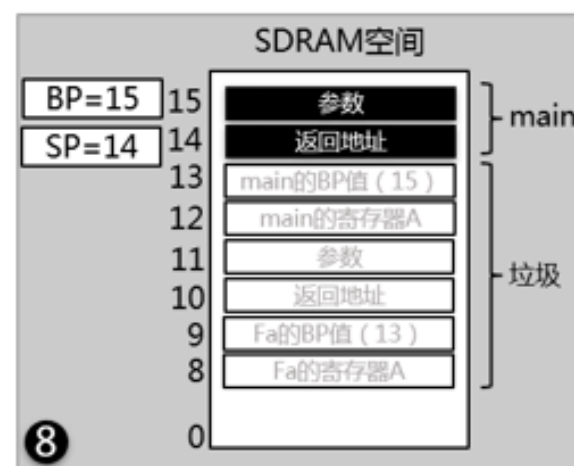
Fa调用Fb，Fb开始执行后



Fb返回到Fa跳转之前的瞬间



Fa回收栈空间中没用的空间



Fa返回到main跳转之前的瞬间

序。如图 5-64 所示，我们可以做个拼图游戏，看看编译器需要从 melonbro\_func.c 文件中将哪个零件拿出并贴上。

Melonbro\_func.c

int handle ( 参数 )  
{具体实现代码};

int otherfunc1 (参数)  
{具体实现代码};

void otherfunc2 (参数)  
{具体实现代码};

char generator( 参数)  
{具体实现代码};

void otherfunc3 (参数)  
{具体实现代码};

void datatransfer(参数)  
{具体实现代码};



让冬瓜哥如数家珍的函数代码文件



将prog.c调用到的位于melonbro\_func.c中的函数代码融合到prog.c之后所生成的程序代码文件prog.cfull



图 5-64 设想中的程序代码融合过程



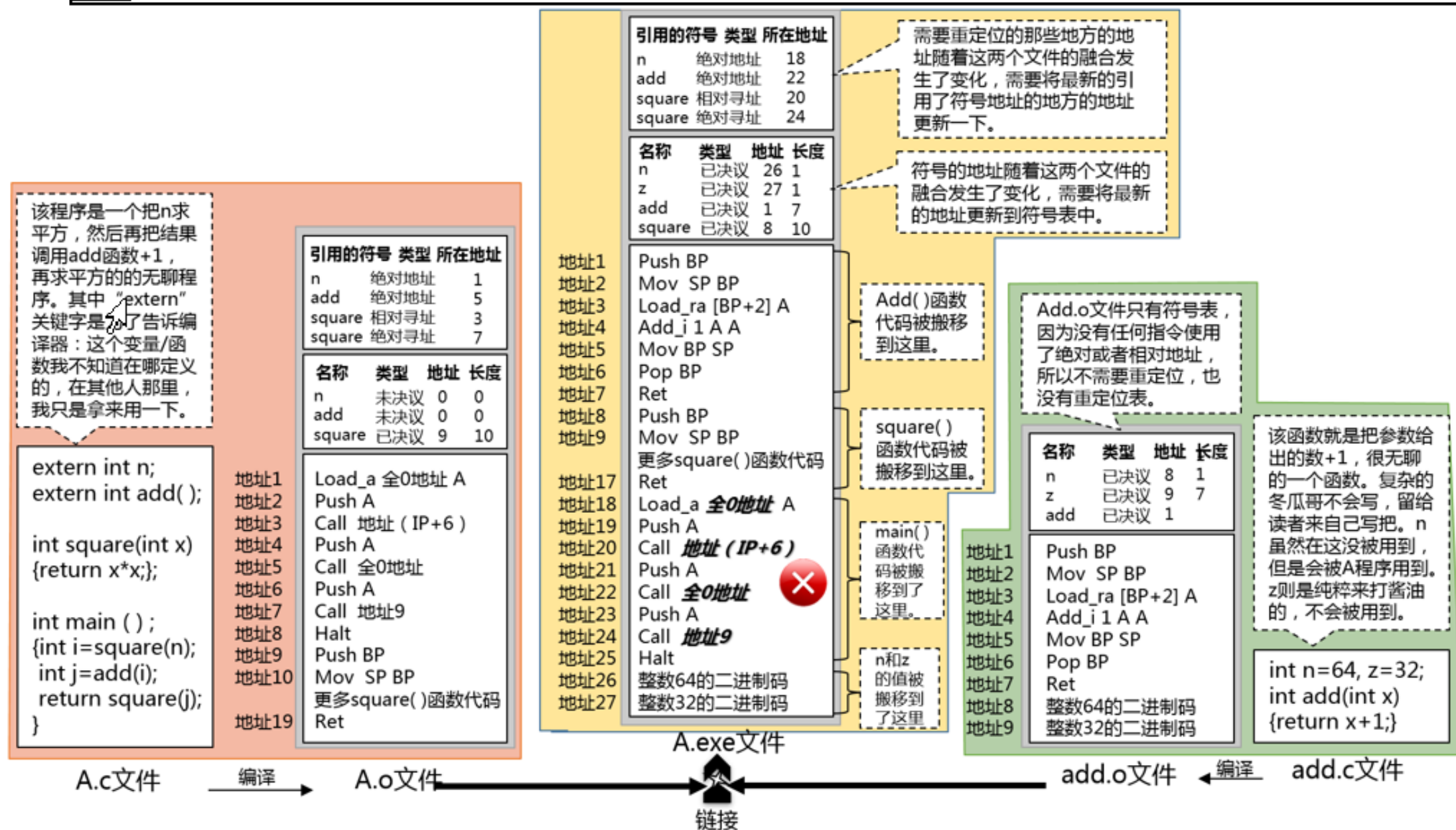
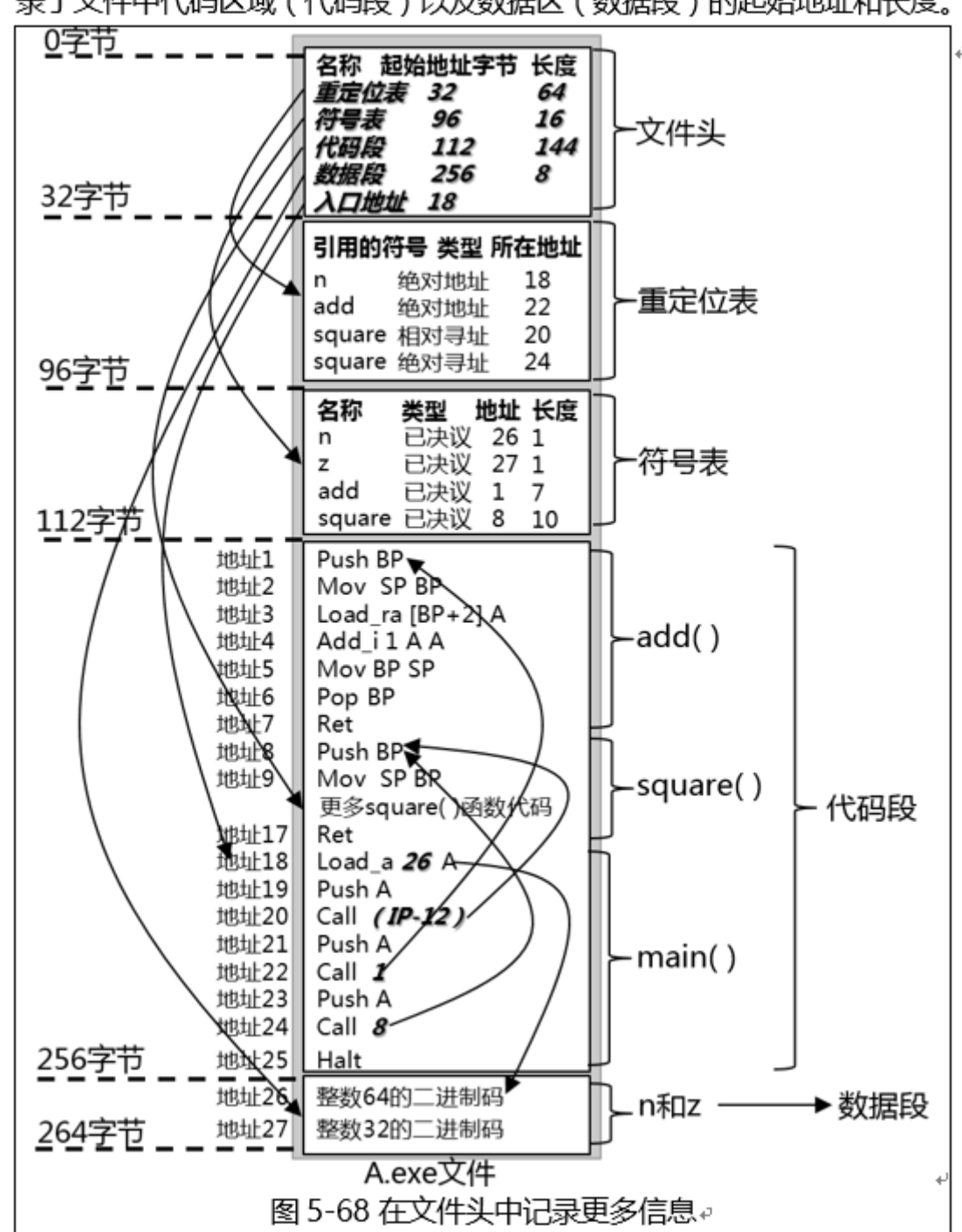


图 5-66 融合后按照符号定义和引用的最新位置更新符号表和重定位表中的地址列











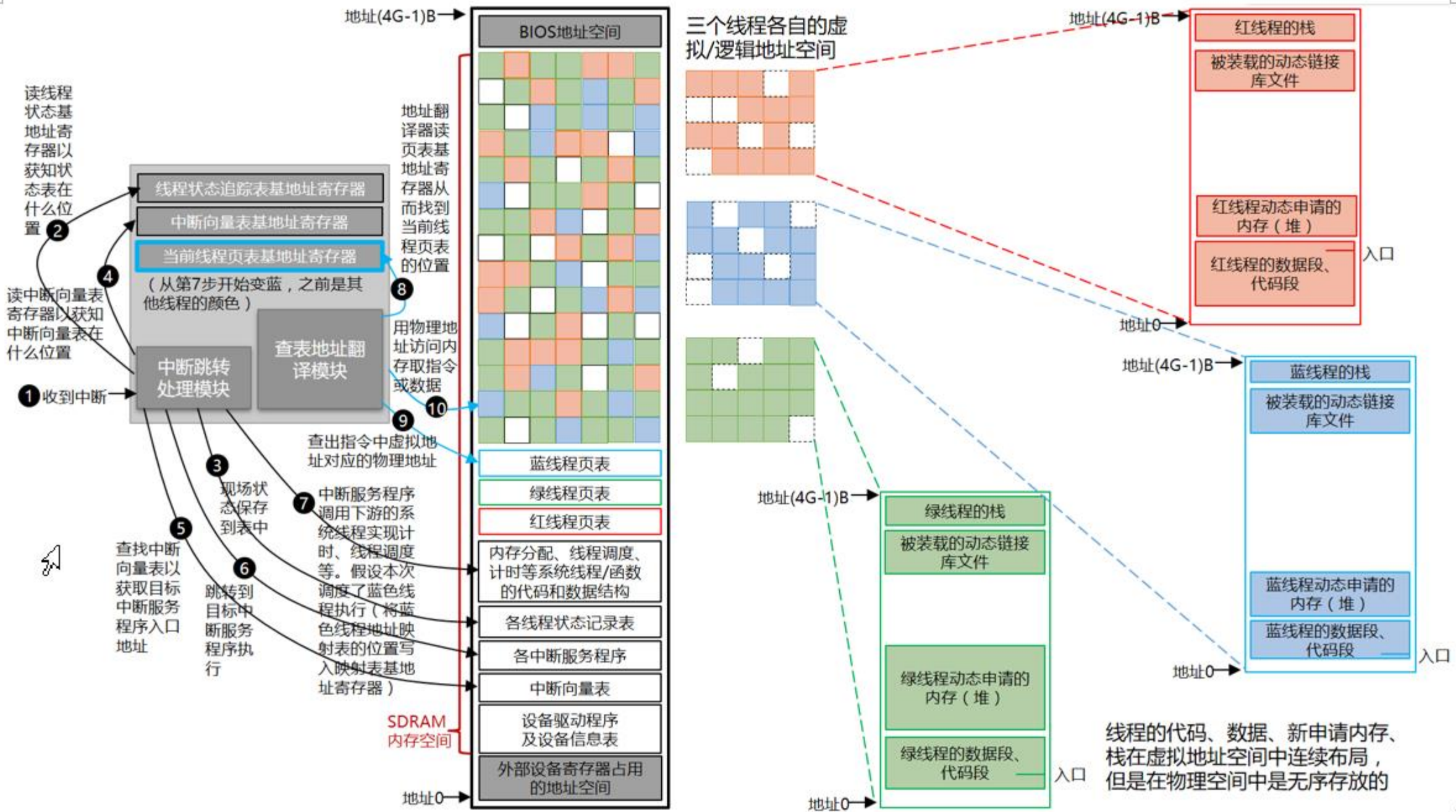


图 5-80 分页内存管理思想示意图

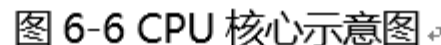
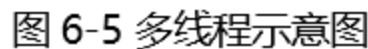
这里禁不住要问了，所谓“DOS 没有限制程序的访存范围”，这意思难道暗指 DOS 是可以去限制的？怎么限制？正如我们前文中所述，当操作系统的 Loader 程序让 CPU 跳转到用户程序运行之后，整个 CPU 就是在运行用户程序了，CPU 此时完全受到用户程序代码的控制，让它走东绝不往西，此时操作系统代码只是静静地呆在内存里起不到任何作用，此时只有靠 CPU 来检查和防止越界。要想实现这个功能，必须将当前执行的代码可访问的地址范围限制在某个区域中，比如从地址 1024 开始的长度 2048 字节的这 2KB 区域中，为了支持这个功能，CPU 必须提供至少两个寄存器，一个用于存放该区域的基地址，也就是上述的地址 1024，另一个用于存放长度，也就是上述的 2048。在操作系统 Loader 程序跳转到用户程序执行之前，必须使用对应的机器指令来更新这两个寄存器，告诉 CPU：“兄弟，后续任何代码只能在这个区域内执行，一旦越界你就报异常，反过来执行我提供的异常服务程序，昂~！”，然后再跳转到用户程序执行，此时用户程序就被框住了。然而，这一招只能防君子，却防不了小人。程序（或者说黑客们）是不会善罢甘休的，程序是不是也可以用对应的机器指令来更新这两个寄存器呢？比如将基地址更新为 0，长度更新为 1GB，从而逃脱限制？这就相当于马路上有个栏杆，守规矩的司机碰到栏杆就会避让，但是不守规矩的可以下车把栏杆往边上移动一下然后说：“你看，我没违规啊”。

设计者早就考虑到这个问题了，所以，当操作系统 Loader 程序让 CPU 跳转到用户程序执行之前，还要告诉 CPU 一件事：“兄弟，接下来任何想更改你的一些特殊控制寄存器的指令你都不执行，一旦发现则直接报异常转回来执行我提供的异常服务程序，昂~！”。这一招下去，相当于接下来执行的程序再也无法逃脱出这个框框了。正如前文中所述，CPU 的指令集中有一些属于特权指令，只有操作系统的代码可以执行特权指令。任何尝试执行特权指令的用户程序，CPU 此时会直接中断程序的运行，跳转到异常服务程序执行，后者直接把该程序终止掉，并弹出窗口或者文字告诉你“刚才这个程序不老实想搞事情，被终止了”，当然，说的好听一点是“非法操作：尝试访问了 xxx 地址”。



第6章	6.2.15 Load/Stor Queue与Stream Buffer	6.5 QPI片间互连网络简介	6.8.6 解决提前执行测错乱问题
6.1 从超线程到多核心	6.2.16 非阻塞Cache与MSHR	6.5.1 QPI物理层与同步异步通信原理	6.8.7 解决乱序执行错乱问题
6.1.1 超线程并行	6.2.17 缓存行替换策略	6.5.2 QPI链路层网络层和消息层	6.8.8 小结
6.1.2 多核心/多CPU并行	6.2.18 I_Cache/D_Cache/TLB_Cache	6.5.3 QPI的初始化与系统启动	6.9 解决多核心访存空间一致性问题
6.1.3 idle线程	6.2.19 对齐和伪共享	6.5.3.1 链路初始化和拓扑发现	6.9.1 基于总线监听的缓存一致性实现
6.1.4 乱序执行还是SMT ?	6.3 连起来，为了一致性	6.5.3.2 系统启动	6.9.1.1 Snarfin/Write Sync方式
6.1.5 逆超线程 ?	6.3.1 Crossbar交换矩阵	6.5.4 QPI的扩展性	6.9.1.2 Write Invalidate方式
6.1.6 线程与进程	6.3.1.1 点阵式Crossbar	6.6 基于QPI互联的高端服务器架构一览	6.9.2 推导MESIF状态机
6.1.7 多核心访存基本拓扑	6.3.1.2 基于复用器的Crossbar	6.6.1 某32路CPU高端主机	6.9.3 MOESI状态机
6.2 缓存十九式	6.3.1.3 Crossbar的级联	6.6.2 IBM x3850/3950 X5/X6主机	6.9.4 结合MESIF协议进一步理解锁和屏障
6.2.1 缓存是分级的	6.3.2 Ring环网	6.6.3 HP Superdome2主机	6.9.5 结合MESIF深刻理解时序一致性模型
6.2.2 缓存是透明的	6.3.3 NoC片上网络	6.6.4 Fujitsu PQ2K主机	6.9.5.1 终极一致性 ( UC )
6.2.3 缓存的容量、频率和延迟	6.3.4 众核Many-Core	6.7 理解多核心访存时空一致性问题	6.9.5.2 严格一致性 ( SC )
6.2.4 私有缓存和共享缓存	6.3.5 多核心程序执行过程回顾	6.7.1 访存空间一致性问题	6.9.5.3 顺序一致性 ( SEC )
6.2.5 Inclusive和Exclusive模式	6.3.6 在众核心上执行程序	6.7.2 访存时间一致性问题	6.9.5.4 处理器一致性 ( PC )
6.2.6 Dirty和Valid标记位	6.4 存储器在网络中的分布	6.7.2.1 延迟到达导致的错乱	6.9.5.5 弱一致性 ( WC )
6.2.7 缓存行Cache Line	6.4.1 CPU片内访存网络与存储器分布	6.7.2.2 访问冲突导致的错乱	6.9.6 缓存行并发写优化
6.2.8 全关联/直接关联/组关联	6.4.2 CPU片外访存网络	6.7.2.3 提前执行导致的错乱	6.9.7 Cache Agent的位置
6.2.9 Virtual Cache	6.4.2.1 全总线拓扑及南桥与北桥	6.7.2.4 乱序执行导致的错乱	6.9.8 基于共享总线的嗅探过滤机制
6.2.10 缓存Homonym问题	6.4.2.2 AMD Athlon北桥	6.8 解决多核心访存时间一致性问题	6.9.8.1 Bitmap粗略过滤
6.2.11 缓存Alias问题	6.4.2.3 常用网络拓扑及UMA/NUMA	6.8.1 有你没我，有我没你	6.9.8.2 Vector Bitmap精确过滤
6.2.12 Page Coloring页面着色	6.4.2.4 AMD Opteron北桥	6.8.2 让子弹飞，等响声来	6.9.8.3 Bloom filter布隆过滤器与哈希采样
6.2.13 小结及商用CPU的缓存模式	6.4.3 参悟全局共享内存架构	6.8.3 硬件原生保证的基本时序	6.9.8.4 JETTY filter
6.2.14 缓存对写入操作的处理	6.4.4 访存网络的硬分区	6.8.4 解决延迟到达错乱问题	6.9.8.5 Stream Register Filter
		6.8.5 解决访问冲突错乱问题	6.9.8.6 Counter Stream Register Filter

#### 6.9.11.8 小结



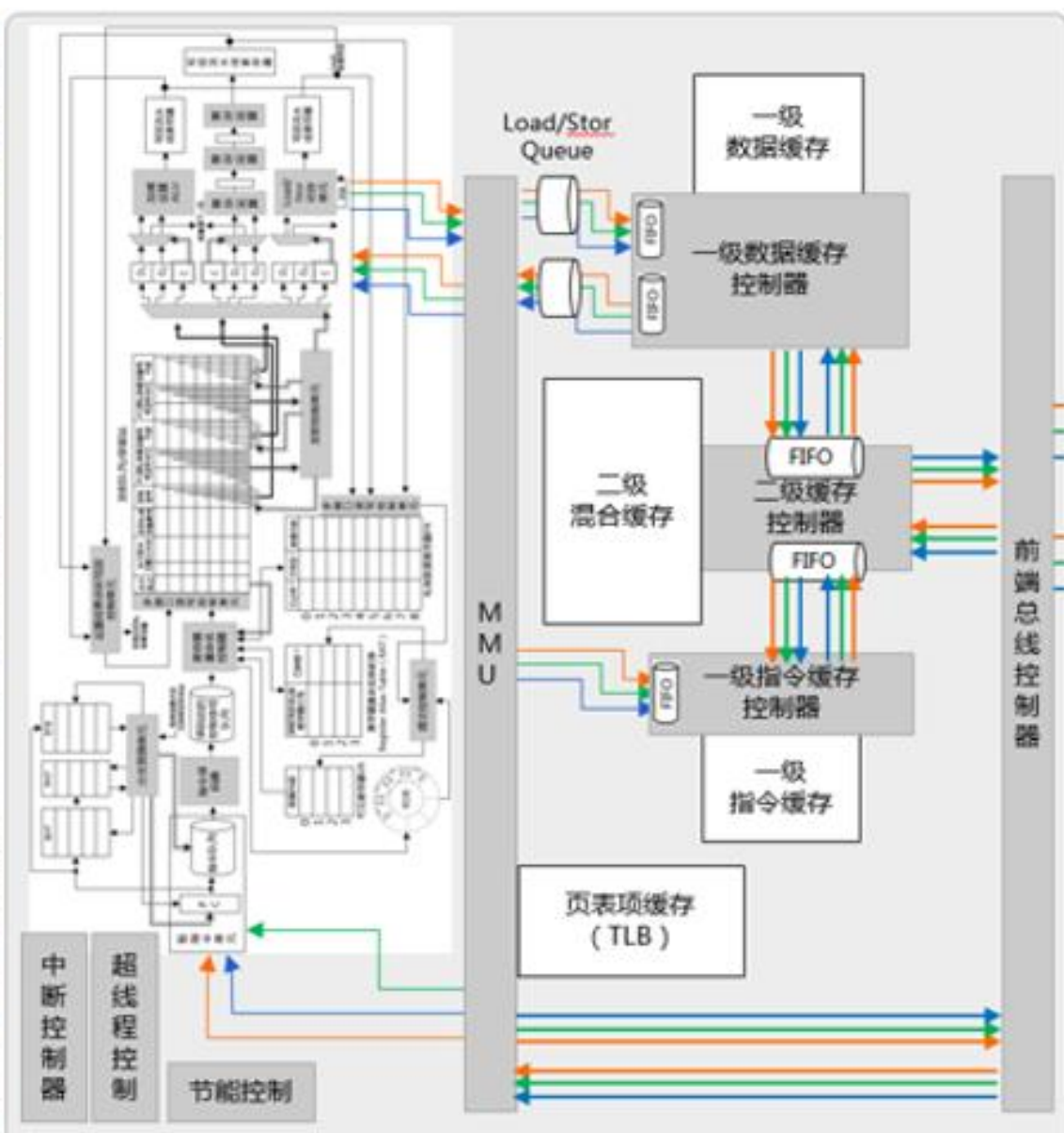
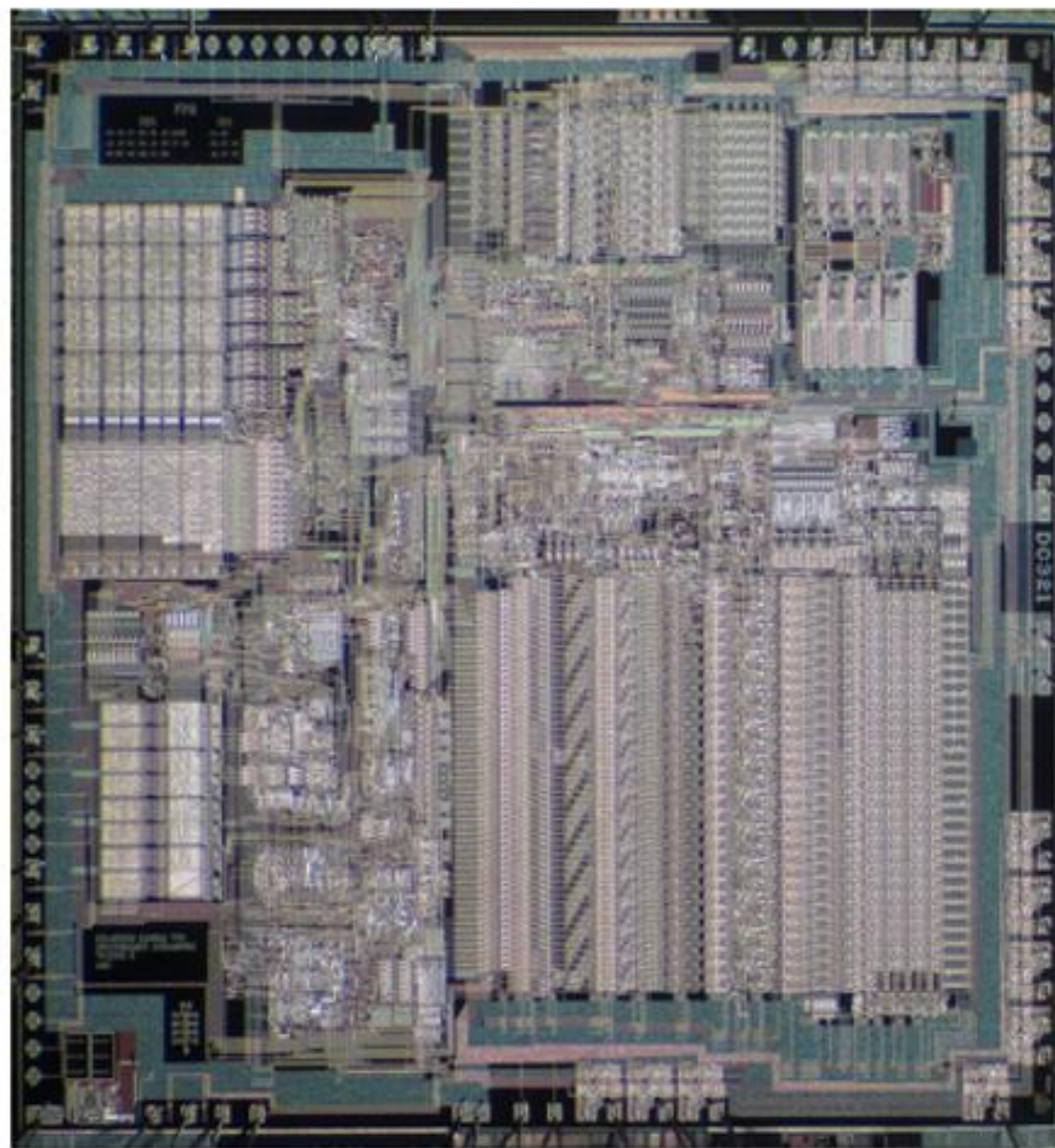
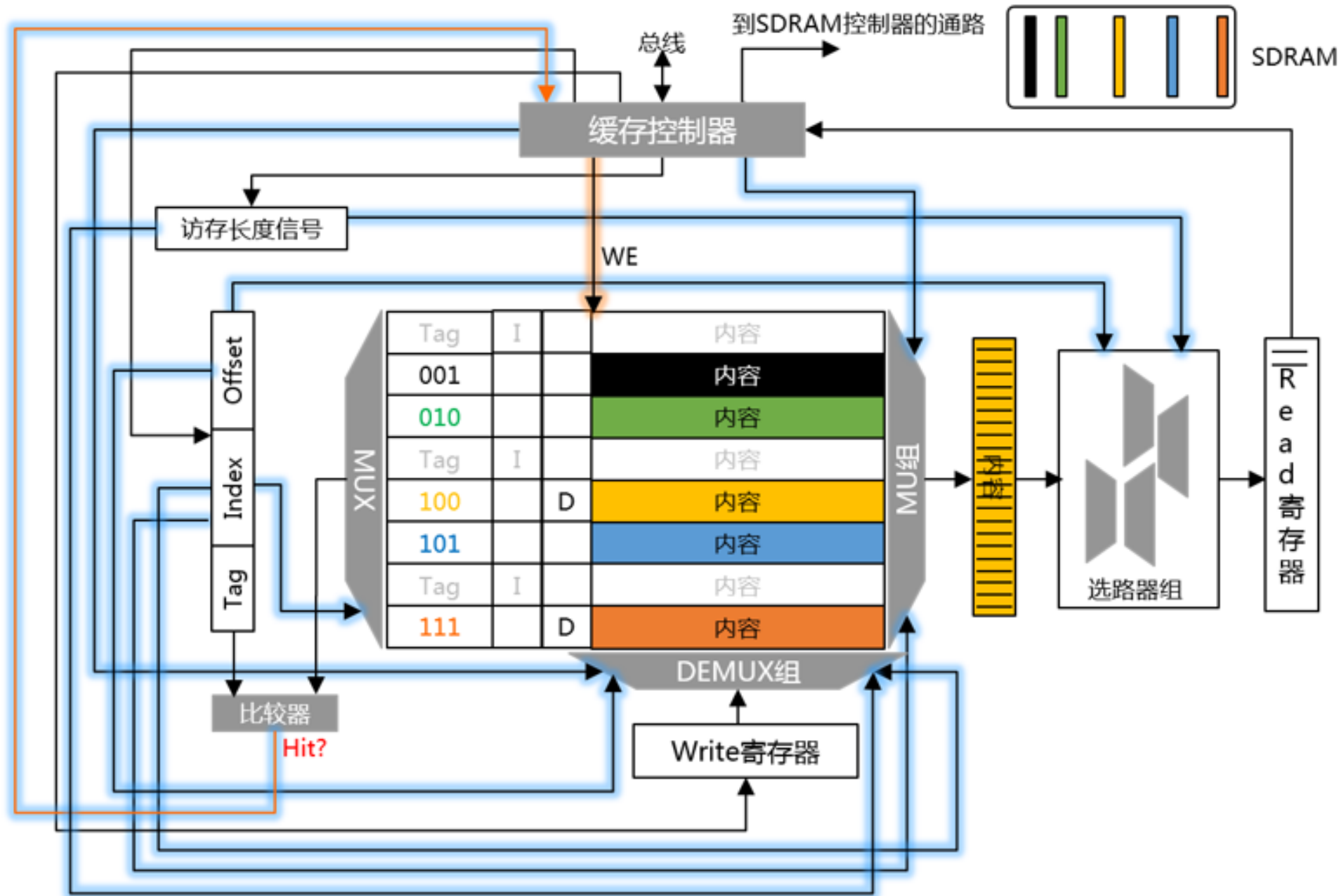


图 6-7 实际的 CPU 核心芯片显微图与架构示意图比较





不需  
地址  
直接  
出对  
然后  
段比  
出数  
SDR

路更  
的运  
较高  
种方  
存行  
置，  
行号

图 6-16 组关联缓存查找过程示意图

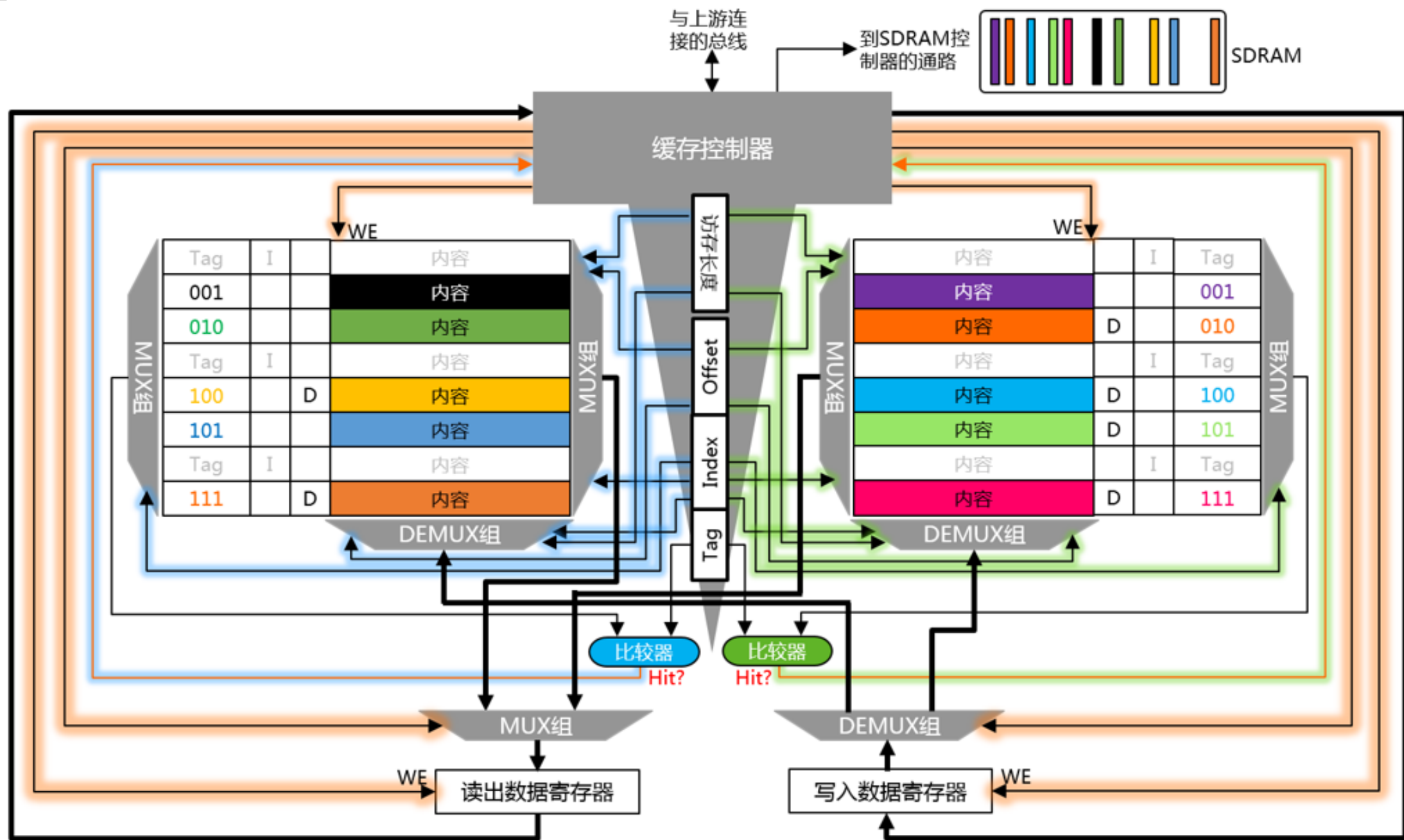
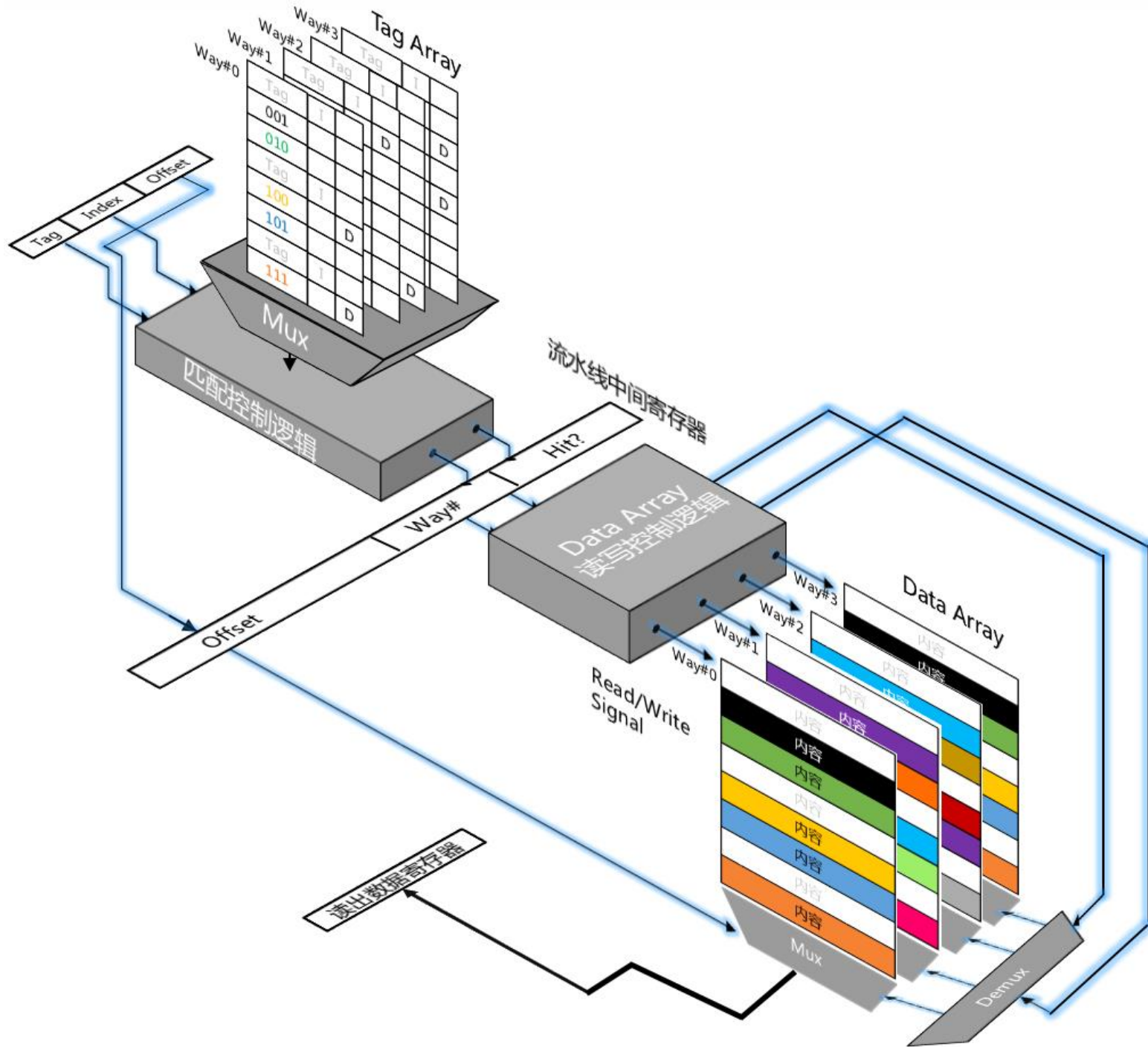


图 6-19 2 路组关联缓存的硬件架构示意图





## 6.2.16 非阻塞 Cache 与 MSHR

对于 CPU 核心，每一条机器指令对它来讲就是一个任务；对于缓存来讲，核心的 L/S 单元发出的访存请求也是一个任务。核心把计算任务交给 ALU 来执行，把访存任务交给 L/S 单元来执行，L/S 单元再把任务交给缓存控制器来执行。在 LSQ 中，可能有多个 Load/Stor 类指令正在等待结果返回，这些请求会按照顺序一条一条的被发送给 L1 缓存控制器去执行。那么假设，排在 LSQ 中队首的某个 Load 指令被发往了 L1 缓存控制器，结果控制器的查找结果是本次访存请求不命中，也就是 Miss 了，那么 L1 缓存控制器将会等待十个左右的时钟周期才能从 L2 缓存中拿到数据，那么在这 10 个左右时钟周期内，难道 L1 缓存个控制器就无事可做了么？显然不划算，我们多么希望 LSQ 中后续的访存请求能够利用这个间隙被缓存控制器执行！

显然，我想所有人都在暗想：如果 L1 缓存控制器能够记录一个状态，比如哪个请求未命中，该请求访问的是哪个地址等等；然后继续执行后续的 Load/Stor 请求，如果又 Miss 了，则继续记录一条追踪记录。当数据返回时，比较这些追踪记录看看是哪个访存请求终于返回了数据，然后缓存控制器就在于 L/S 单元连接的总线上向 L/S 单元返回对应的数据，并对所有的返回数据加以区分，以告诉 L/S 单元这次返回的是哪一条访存请求的数据，比如通过 ID 标记机制，每个请求都有各自的 ID，也可以直接在总线上通告该数据要被载入到的目标寄存器号。L/S 单元拿到返回的数据之后，根据 LSQ 中记录的访存请求的目标数据寄存器号，直接将数据载入对应的核心的数据寄存器从而完成访存操作。当然，这里其实是返回到了内部的私有寄存器中，会被流水线结果侦听单元捕获到，然后更新到保留站以及私有寄存器中的记录上，从而依赖这些数据的、依然等待在保留站中的其他指令就可以在下一个时钟周期内被发射控制单元判断为符合执行条件，从而被调度执行，这个过程建议回顾流水线那一章里的流程图。

这种支持前序访存请求即便 Miss 了也不影响后续访存指令继续执行的缓存被称为**非阻塞缓存，None Blocking Cache**。上面提到过的用于追踪各个 Miss 的访存请求的执行状态的硬件模块被称为 **MSHR ( Miss Status Holding/Handling Register )**。



图 6-27 MSHR 寄存器组

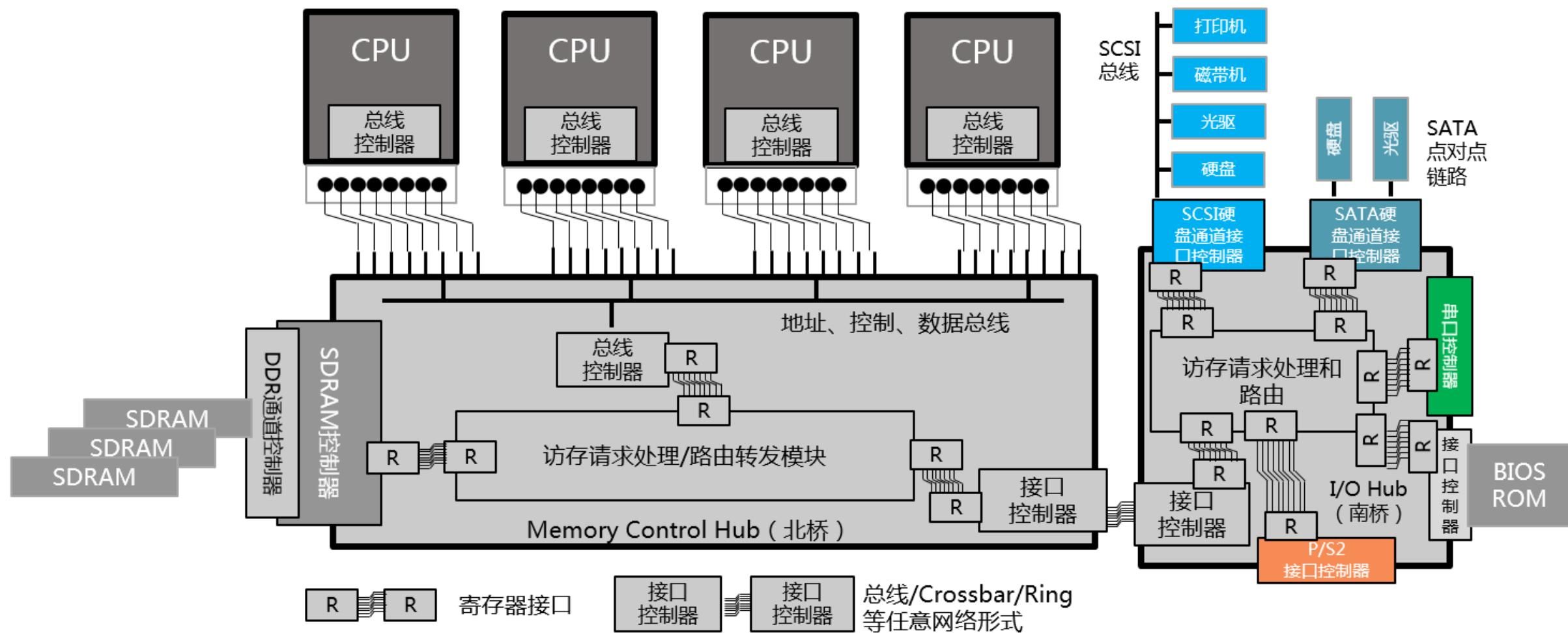




图 6-91 输入开关及跳线

在早期，很多参数都是通过这种办法输入到电路中去的，包括有些 I/O 总线上的节点地址 ID，都要用跳线帽来短接对应的引脚从而让电路区分不同的设备。如图 6-91 所示为早期主板上的用于配置这些参数的开关。这些开关对应的电路会直接输入到目

标模块中，目标硬件模块每次加电之后，直接从这些输入开关获取其所表示的 0/1 信号，将其直接锁存到内部的寄存器中，这样，与将这些寄存器映射到系统全局地址空间，访存网络中配好路由信息，然后用机器指令的方式操纵 CPU 来写入 ID 到这些寄存器的效果是一样的。

这种操作非常笨拙，人们总是不断追求更方便的方式。于是人们想出了另外一种方式，用另一台计算机运行某个程序，该计算机一方面负责通过 Web 页面等方式与用户交互（当然，用户也需要在用一台计算机打开浏览器来访问这台计算机的 Web 页面），另一方面则与需要被配置的硬件电路对接，对接方式可以通过 i2c 等串行接口方式，从而可以将用户发送过来的配置请求，转换为在 i2c 总线上传递的信息，发送给目标电路，目标电路从 i2c 总线上接收这些配置信息，然后配置到自己的寄存器中，完成分区或者各种其他操作。由于这台计算机需要用 i2c 与待配置的电路模块相连，所以它的实际形态是一片 SoC（System on Chip），集成了 CPU、SDRAM、i2c 接口控制器、以太网控制到一个单一芯片中，该芯片本身就是一台小计算机，与它要配置的那台大计算机放在同一张主板上，其内部运行着程序，包括底层的支撑程序（操作系统、各种控制器驱动程序）、Web 网页服务端程序、后端向 i2c 总线发送配置信息的程序等。Web 服务程序通过以太网接收用户通过网页访问方式发来的配置请求。经过这样的设计，用户再也不用去打开机箱，按照说明手册搬动开关去向电路输入参数了。



这个专门用于对计算机做更加底层的配置或者监控的片上小计算机，被称为 **BMC**，**Baseboard Management Controller**，**基板控制器**。如图 6-92 所示为当前比较流行的 BMC 芯片，可以看到其下方的一片 1GB 容量的 SDRAM，是供 BMC 运行使用的内存。

当前的高端服务器计算机全都使用 BMC 来做基本配置和监控，其不仅仅是配置这些初始参数，其还可以直接对计算机机箱内的其他模块做控制和监控，比如风扇、电源等。只要系统电源有达标的电压







核心1

a = 1

Barrier()

b = 1

核心2

while(1) {

if (b == 0) { continue;}

if (a == 1) { do();} }

核心1

a = 1

Barrier()

b = 1

核心2

while(1) {

if (b == 0) { continue;}

**Barrier()**

if (a == 1) { do();} }

b 0 M a 0 S

核心1

Stor\_i 1 a

WrtIvld

a 1 M

Sfence

b 1 M

Stor\_i 1 b

RdRspDC

b=1所在缓存行

b 1 S

b 1 a 0 S

核心2

先回应后干活

a 0 S

Load b 寄存器A

b 0 F

b 1 S

Cmp 寄存器A 0

Jmpz\_b 2

Load a 寄存器A

a 0 S

Cmp 寄存器A 1

a 1

Jmpz do()

Jmp\_b 6

到这都还没作废

终于作废但是为时已晚

访存共享总线

b 0 M a 0 S

核心1

Stor\_i 1 a

WrtIvld

a 1 M

Sfence

b 1 M

Stor\_i 1 b

RdRspDC

b=1所在缓存行

b 1 S

b 1 a 0 S

核心2

先回应后干活

a 0 S

Load b 寄存器A

b 0 F

b 1 S

Cmp 寄存器A 0

Jmpz\_b 2

**Lfence**

改了再往下走

Load a 寄存器A

a 1 S

Cmp 寄存器A 1

Jmpz do()

Jmp\_b 6

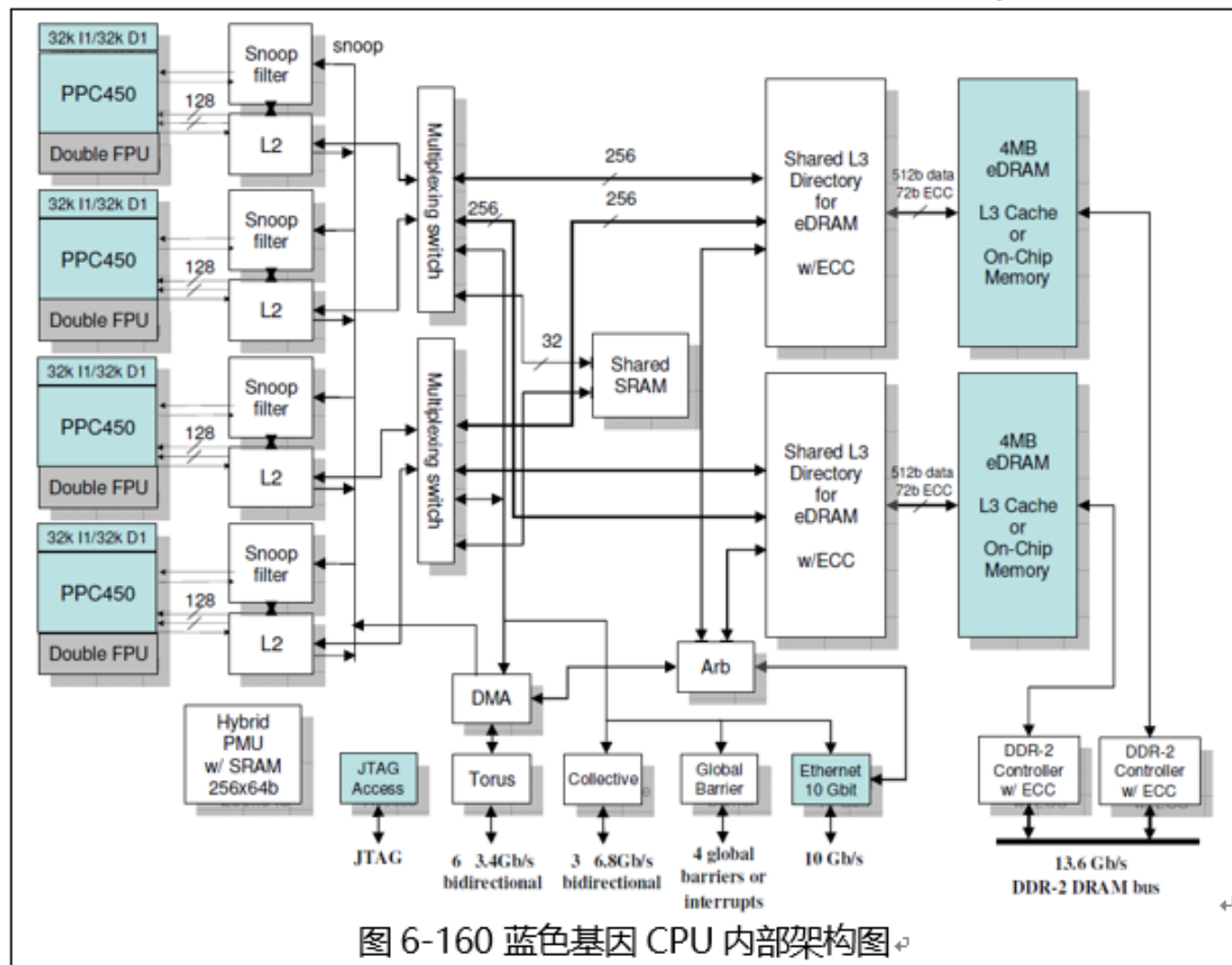
符合条件跳到do函数执行

访存共享总线



### 6.9.8.7 蓝色基因/P 中的 Snoop Filter

IBM 在其蓝色基因/P 超级计算机的 CPU 中使用了三种 Snoop Filter 来过滤广播。



如图 6-160 所示为蓝色基因 CPU 内部架构图。其 CPU 的 L1 缓存采用 Write Through 模式透写到 L2 缓存及共享的 L3 缓存，也就是从 L1 直接穿透到 L3，但是 L2 中也会留存一份。所以不需要在 L1 而只需要在最后一层私有缓存处也就是 L2 缓存控制器前端需要过滤装置即可。

如图 6-161 左侧所示，每个 PowerPC 450 CPU 核心前方都放置一个 L2 缓存控制器和一个 Snoop Filter。所有的 L2 缓存控制器与 4 个 Snoop Filter 采用点对点方式广播 CC 请求。

可以看到缓存控制器之间并没有直接通路，所以核心之间无法通过 L2 缓存直接相互转发数据，但是由于 L1 是写透到 L3 的，所以 L3 里始终都是最新的数据，当某个 CPU 收到其他 CPU 的 L2 缓存控制器发送的 Invalidate 消息从而作废了其自己 L2 缓存内某缓存行之后，后续针对该行的读取可直接到 L3 读取，不需要直通转发。

## 没有远程目录时无法源过滤

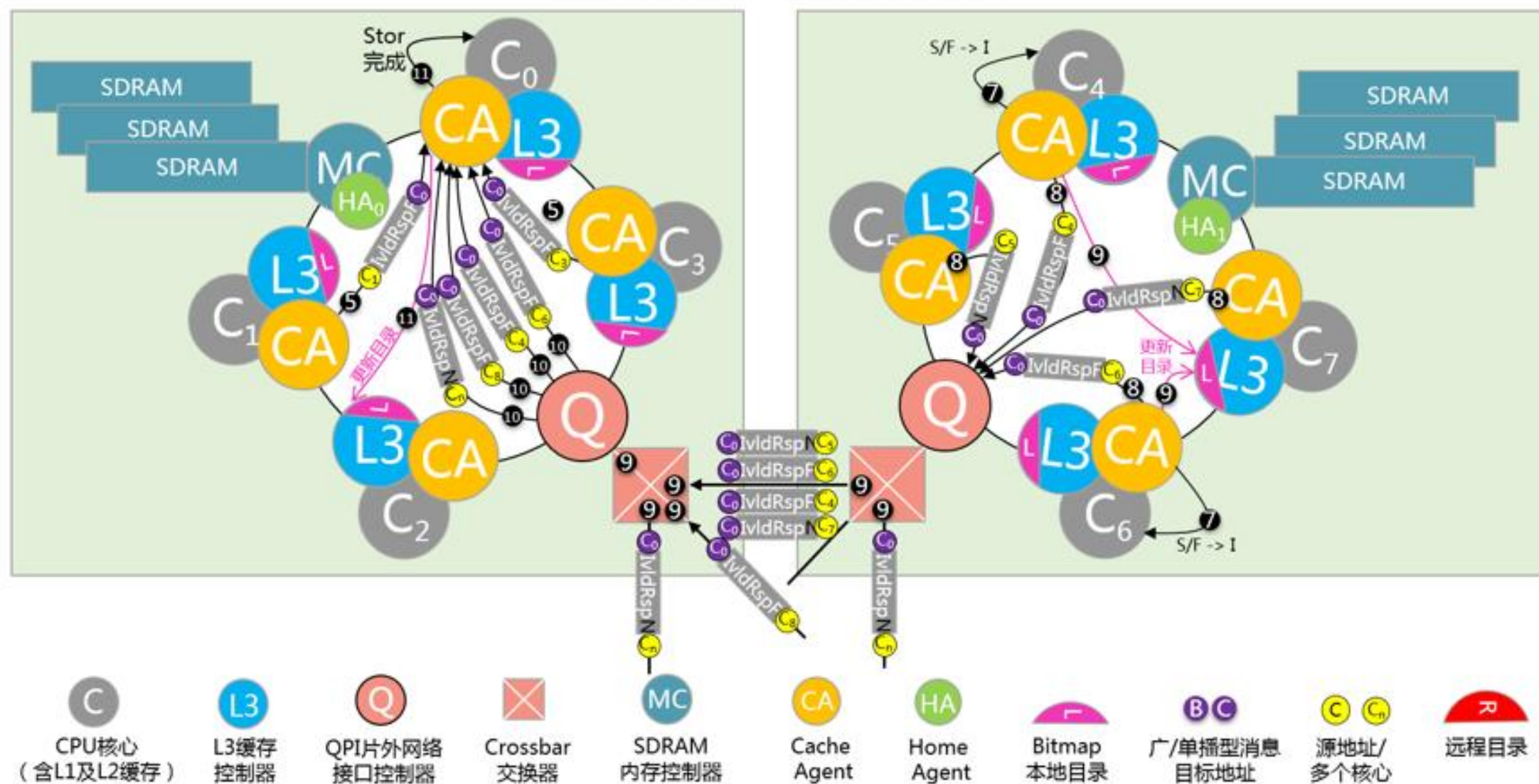


图 6-168 没有远程过滤目录时的 CC 示意图 (2)

## ▲ 第7章

### ▲ 7.1 计算机I/O的基本套路

7.1.1 Programmed IO+Polling模式

7.1.2 DMA+中断模式

7.1.3 DMA与缓存一致性

7.1.4 Scatter/Gather List ( SGL )

7.1.5 使用队列提升I/O性能

### ▲ 7.1.6 固件/Firmware

7.1.6.1 固件与OS的区别与联系

7.1.6.2 固件的层次

7.1.6.3 固件的格式

7.1.6.4 固件存在哪

7.1.6.5 固件如何加载运行

7.1.7 网络I/O基本套路

7.1.8 接入更多外部设备

7.1.9 一台完整计算机的全貌

7.2 中断处理

### ▲ 7.3 网络通信系统

#### ▲ 7.3.1 OSI七层标准模型

7.3.1.1 应用层

7.3.1.2 表示层

7.3.1.3 会话层

7.3.1.4 传输层

7.3.1.5 网络层

7.3.1.6 链路层

7.3.1.7 物理层

7.3.1.8 传送层

7.3.1.9 小结

#### ▲ 7.3.2 底层信号处理系统

7.3.2.1 AC耦合电容及N/Mb编码

7.3.2.2 加扰的作用

7.3.2.3 各种线路编码

7.3.2.4 各种模拟调制技术

7.3.2.5 频谱宽度与比特率

7.3.2.6 数字信号处理与数字滤波

#### ▲ 7.3.3 以太网——高速通用非访存式后端外部网络

7.3.3.1 以太网的网络层

7.3.3.2 以太网的链路层和物理层

7.3.3.3 以太网I/O控制器

### ▲ 7.4 典型I/O网络简介

#### ▲ 7.4.1 PCIE——高速通用访存式前端I/O网络

7.4.1.1 PCI网络拓扑及数据收发过程

7.4.1.2 PCI设备的配置空间

7.4.1.3 PCI设备的枚举和配置

7.4.1.4 PCI设备寄存器的物理地址分配和路由

7.4.1.5 中期小结

7.4.1.6 PCIE网络拓扑及数据收发过程

7.4.1.7 PCIE网络的层次模型

7.4.1.8 NTB非透明桥

7.4.1.9 PCIE Switch内部

7.4.1.10 在PCIE网络中传递消息

7.4.1.11 在PCI网络中传递中断信号

7.4.1.12 在PCIE网络中传递中断信号

7.4.1.13 MSI/MIS-X中断方式

7.4.1.14 PCIE体系中的驱动程序层次

7.4.1.15 小结

#### ▲ 7.4.2 USB——中速通用非访存式后端I/O网络

7.4.2.1 USB网络的基本拓扑

7.4.2.2 USB设备的枚举和配置

7.4.2.3 USB网络协议栈

7.4.2.4 USB网络上的数据包传送

7.4.2.5 USB网络的层次模型

7.4.2.6 小结

#### ▲ 7.4.3 SAS——高速专用非访存式后端I/O网络

7.4.3.1 SAS网络拓扑及设备编号规则

7.4.3.2 SAS网络中的Order Set一览

7.4.3.3 SAS的链路初始化和速率协商

7.4.3.4 SAS网络的初始化与设备枚举

7.4.3.5 SAS和SCSI的Host端协议栈

7.4.3.6 形形色色的登记表

7.4.3.7 SAS网络的数据传输方式

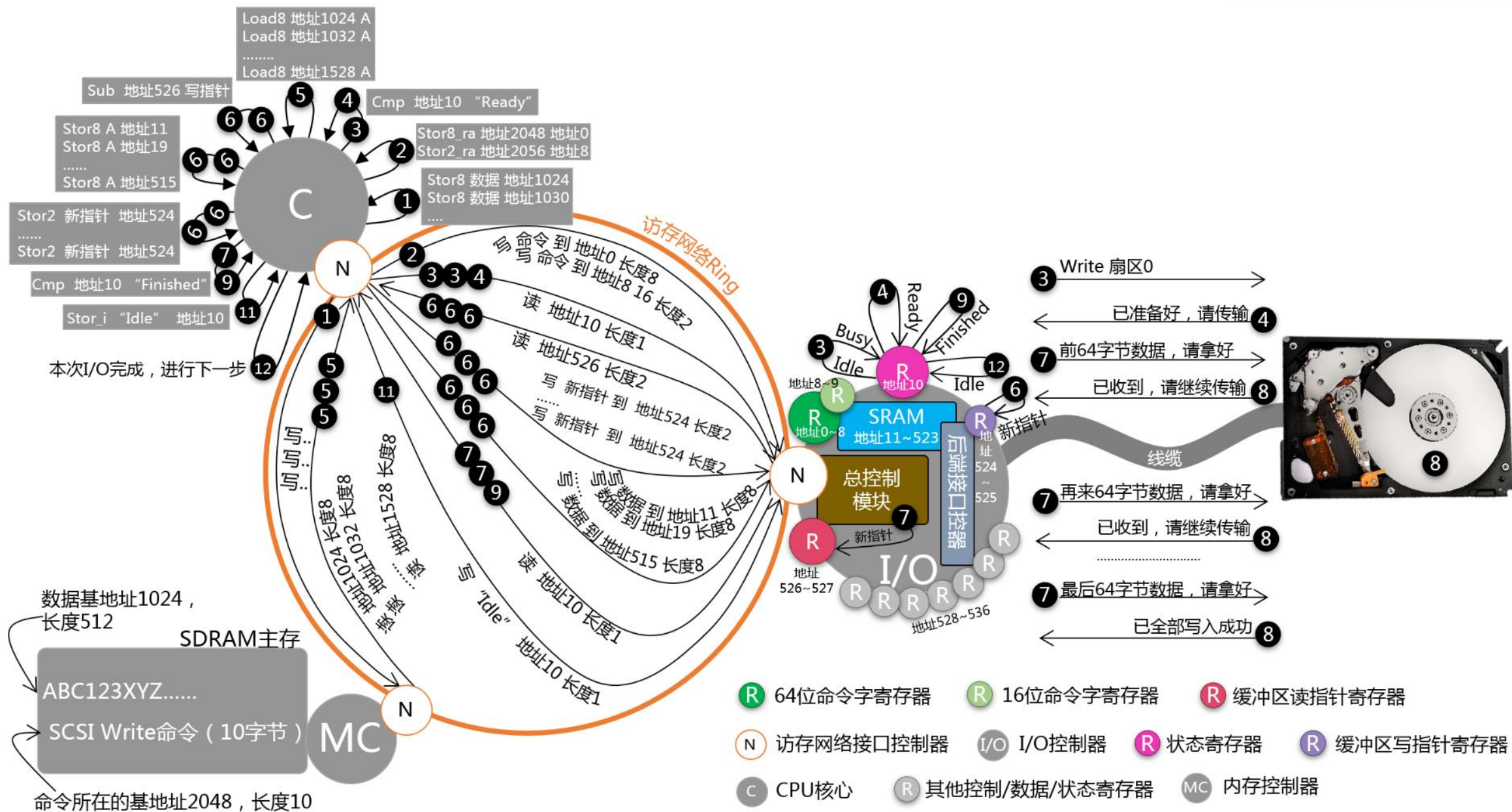
7.4.3.8 SAS网络的层次模型

7.4.3.9 SAS控制器内部架构

7.4.3.10 SAS Expander内部架构

7.4.3.11 基于SAS体系的存储系统形态





上角部分中的紫色立方体)，然后将该数据包的元数据写回到接收队列中的对应条目的对应位置（图中绿色区块）。然后，网络 I/O 控制器发出中断信号，Host 端程序对该数据包进行后续处理，并更新队尾指针。至于“后续处理”都处理的什么东西，怎么个流程，我们会在后文中的网络协议栈相关章节详细介绍。

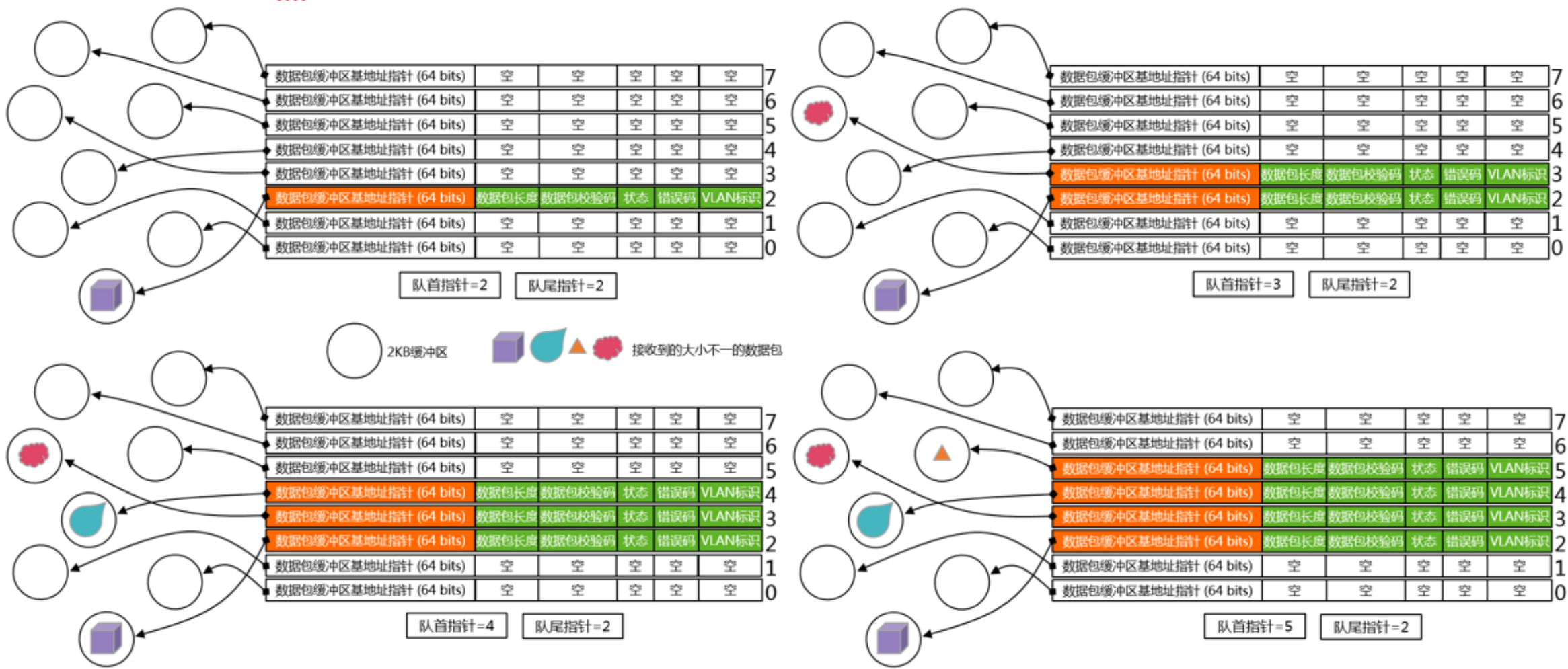


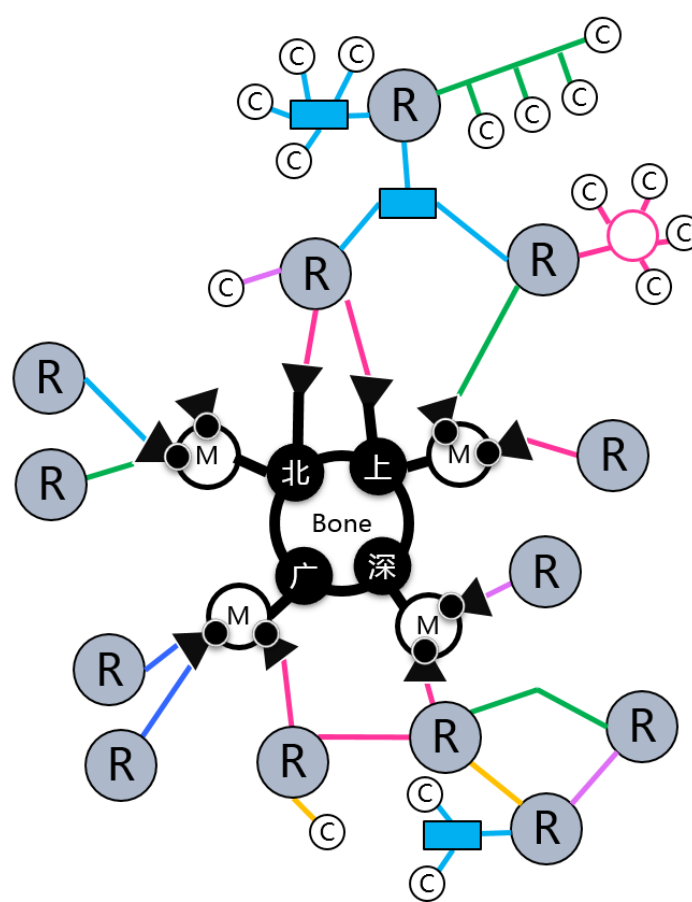
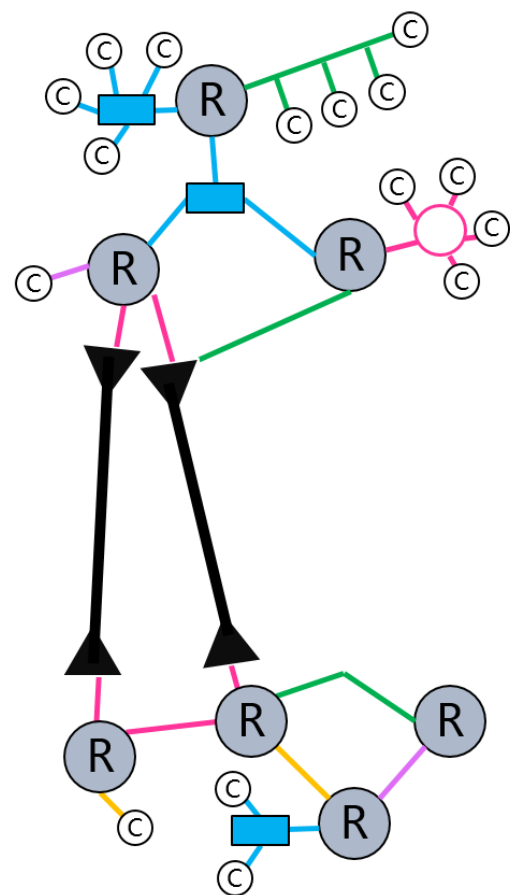
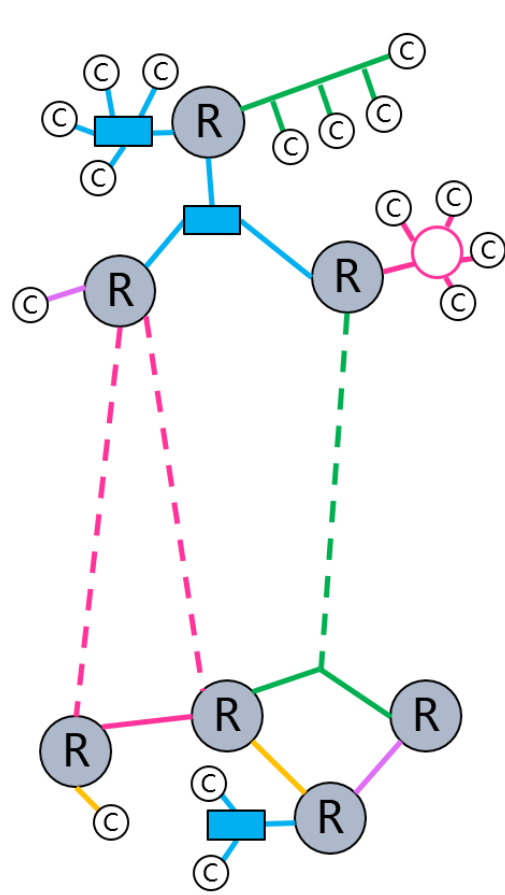
图 7-12 网络 I/O 控制器接受数据包并写入主存中缓冲区的过程

我们假设 Host 端程序还没有来得及处理接收到的数据包，结果网络控制器又接收到了多个数据包。如上图中右上角、左下角、右下角所示。此时接收队列中一共积压了 4 个数据包。









▼ 接入网传输设备

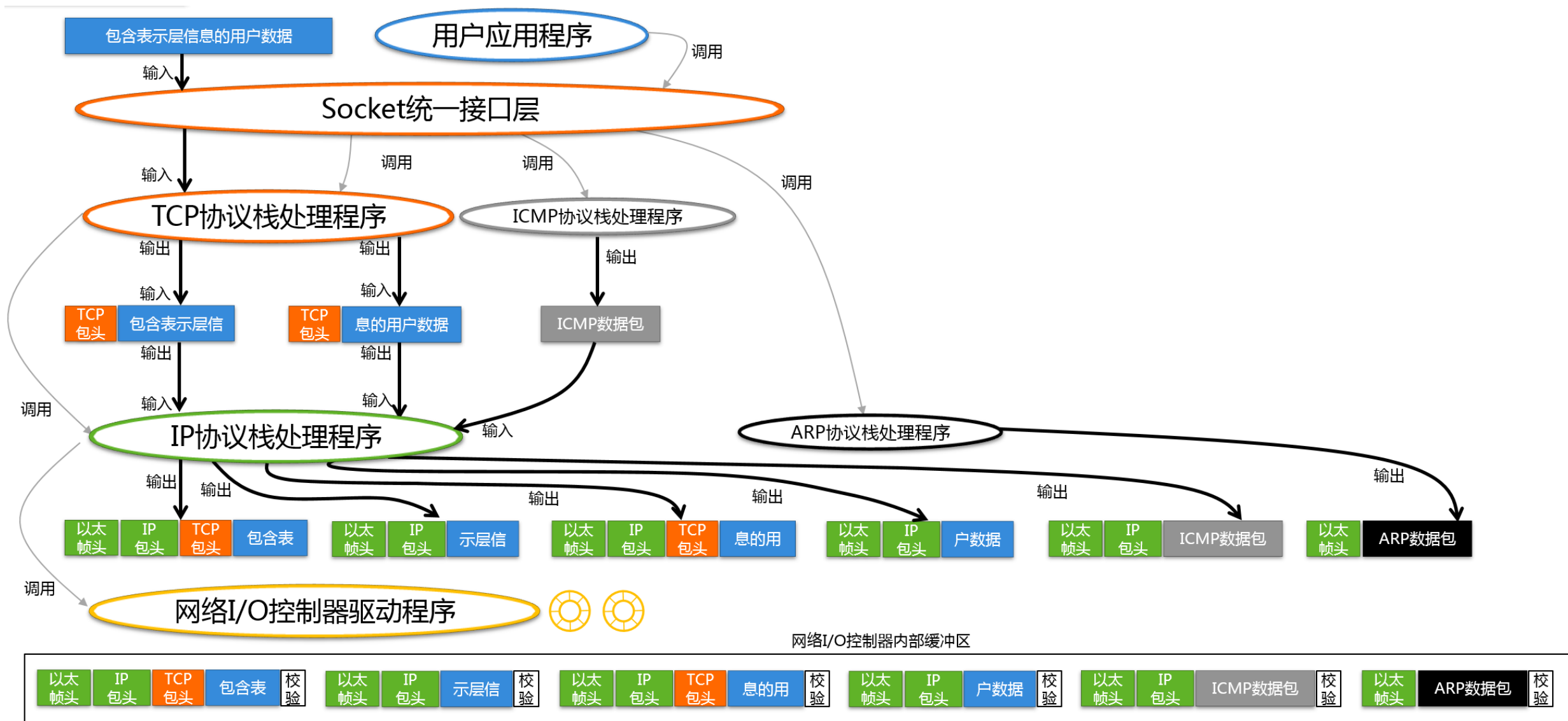
⊙(M) 城域网

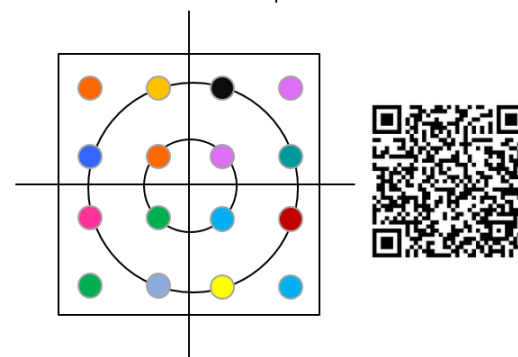
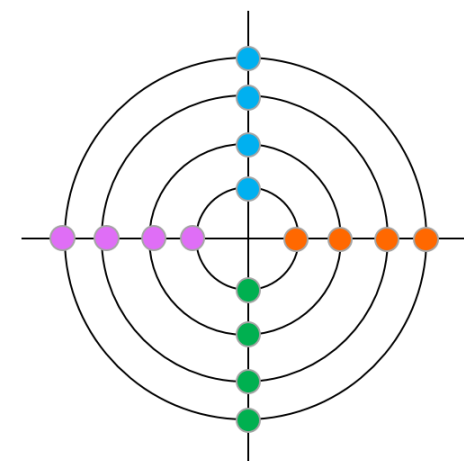
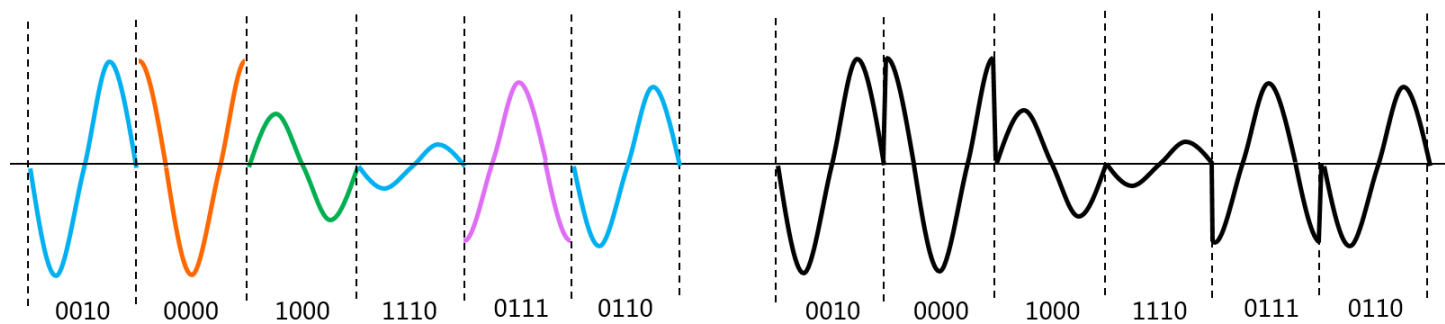
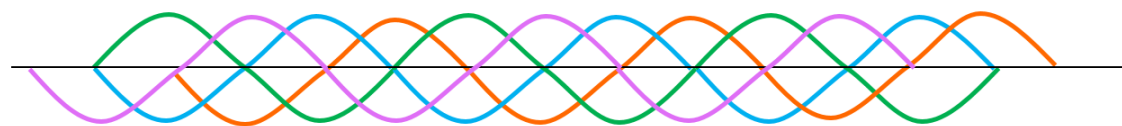
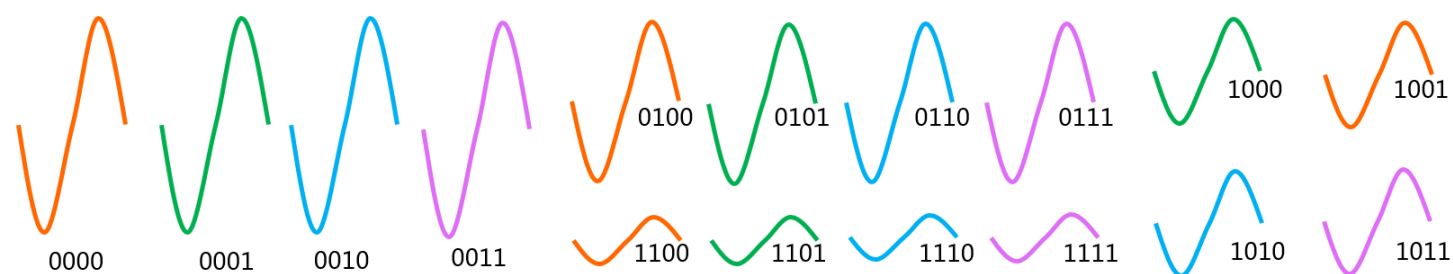
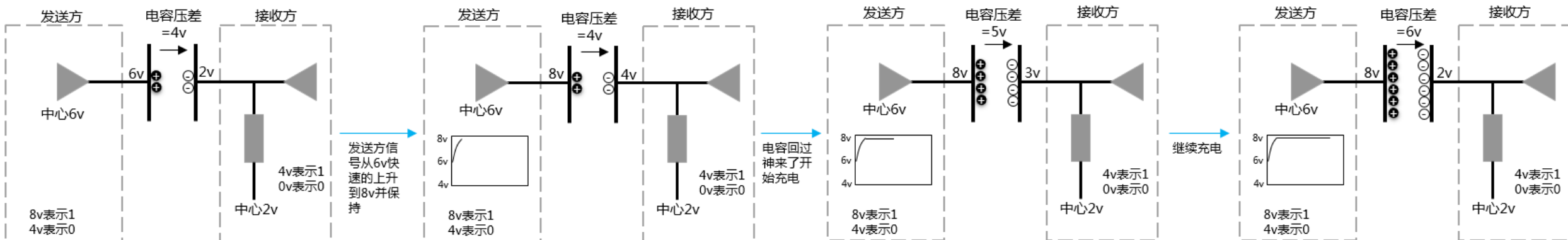
● 城域网传输设备

● 骨干网传输设备

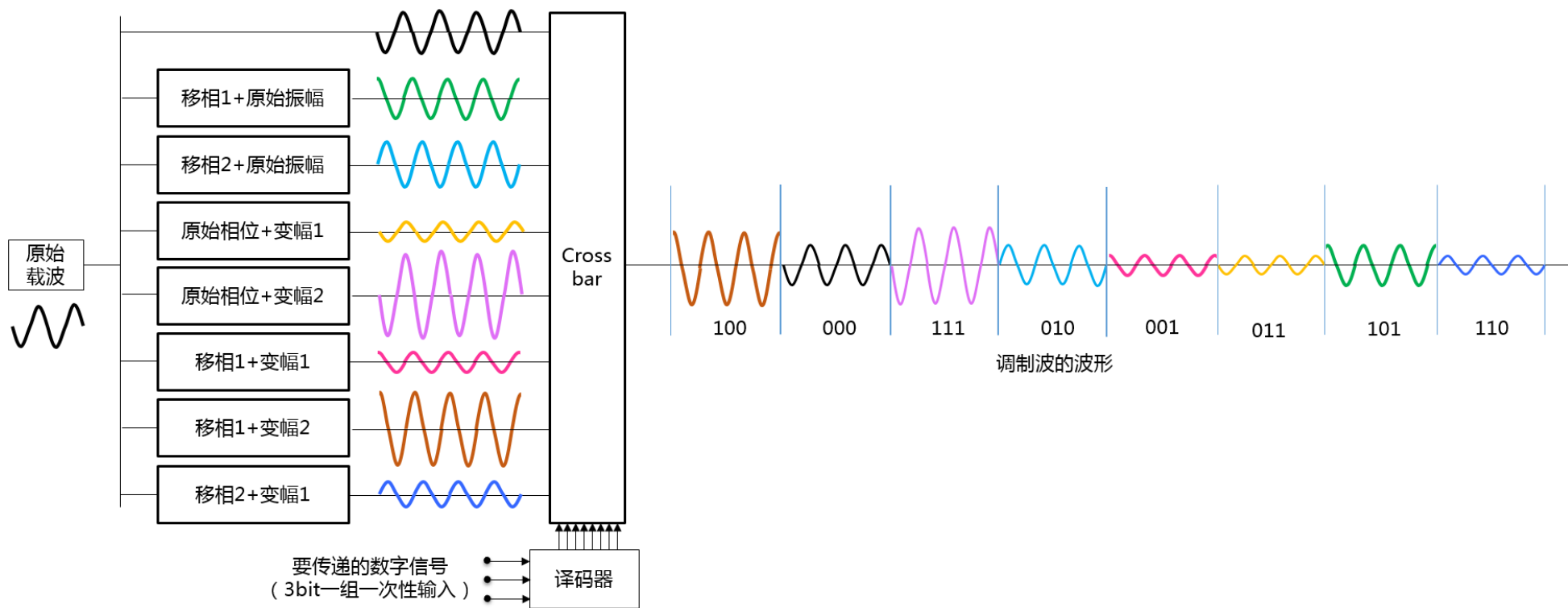
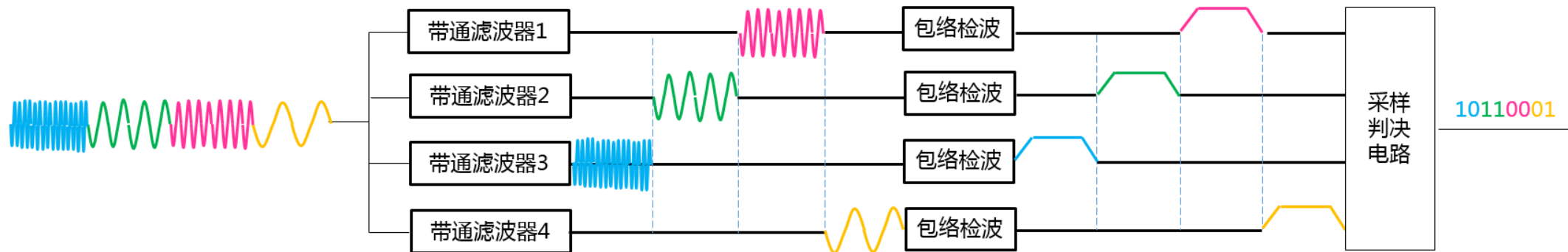
○ 骨干网

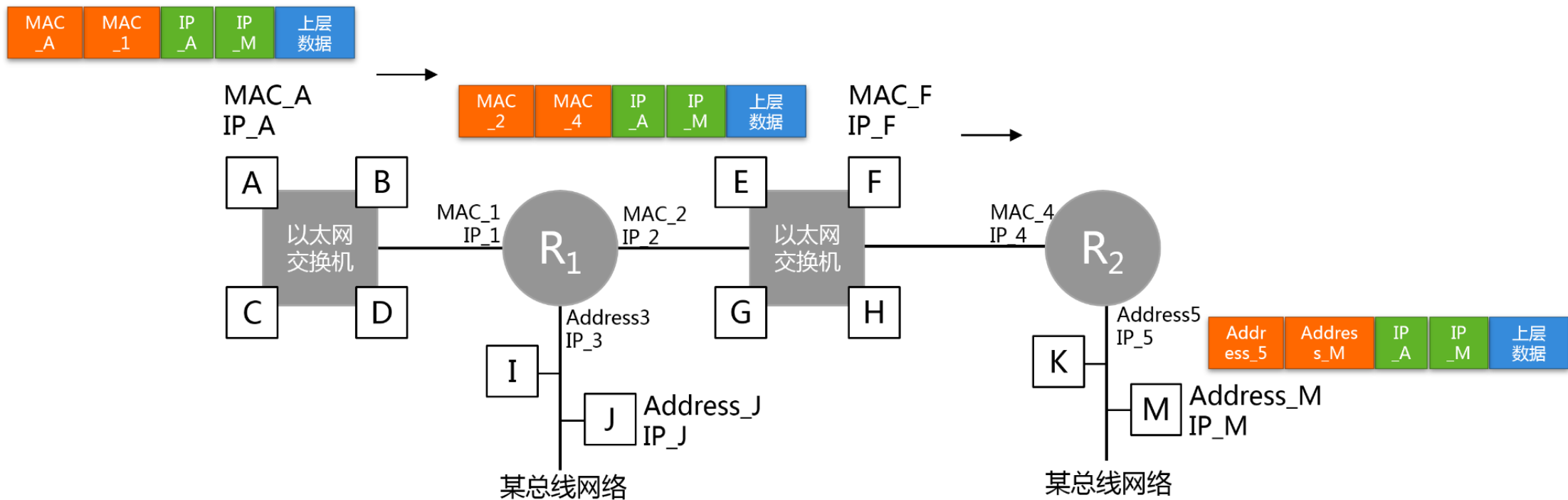
圆形部分整体称为传输网











程序用Stor指令向CONFIG\_ADDRESS寄存器写入配置读指令，包含设备编号、配置寄存器号等信息

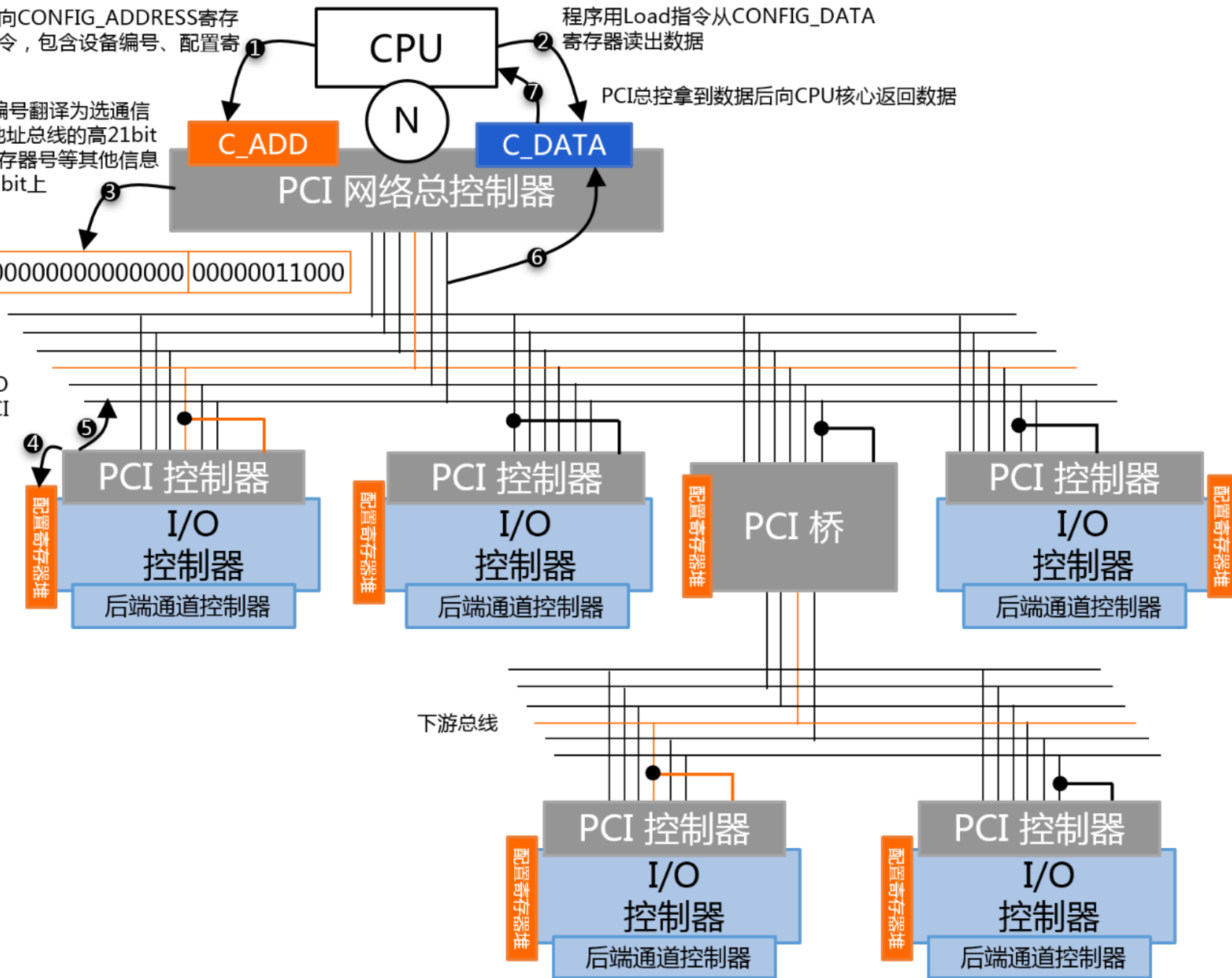
程序用Load指令从CONFIG\_DATA寄存器读出数据

PCI总控将设备编号翻译为选通信号放置在数据/地址总线的高21bit上，并将配置寄存器号等其他信息放到总线的低11bit上

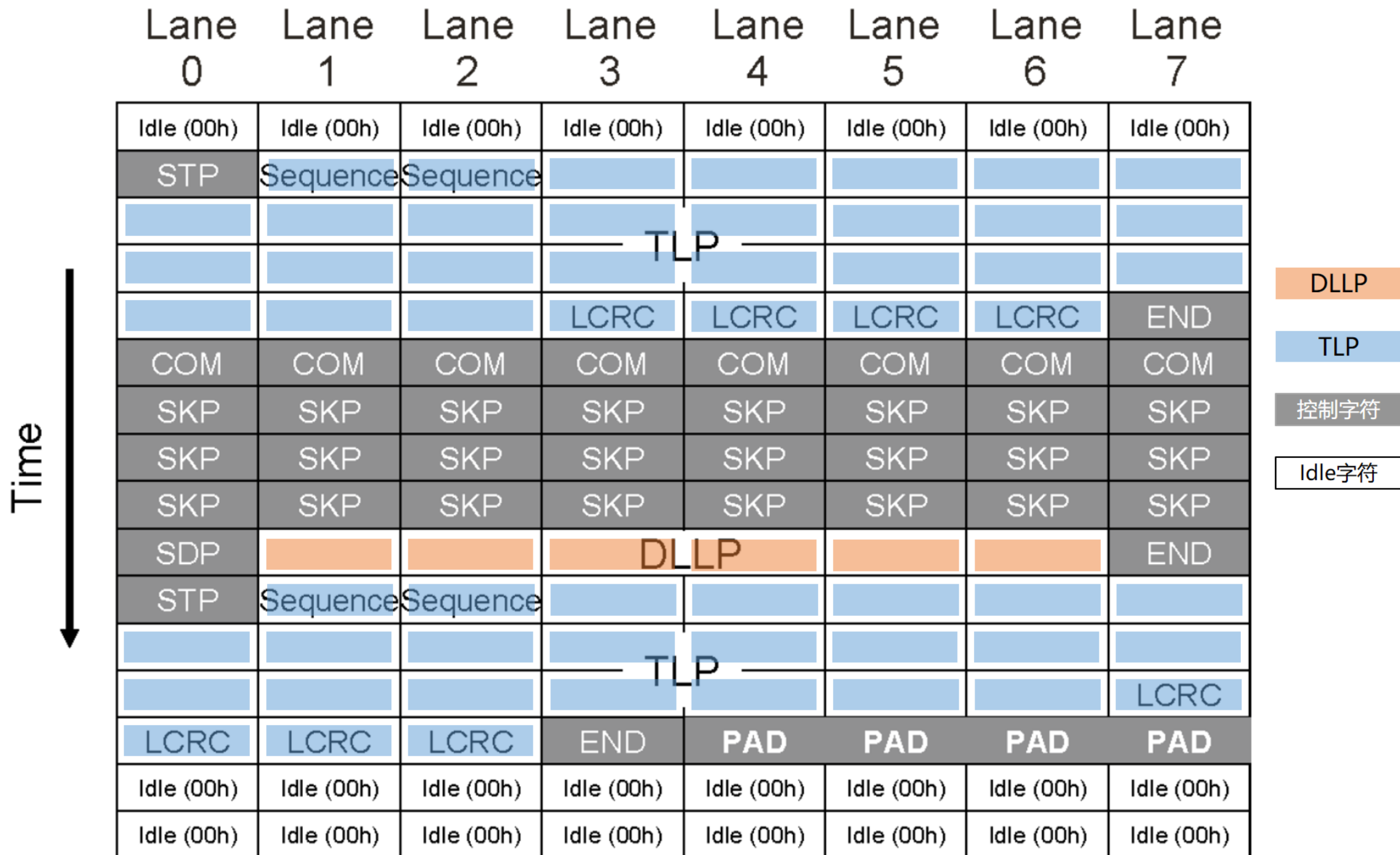
PCI总控拿到数据后向CPU核心返回数据

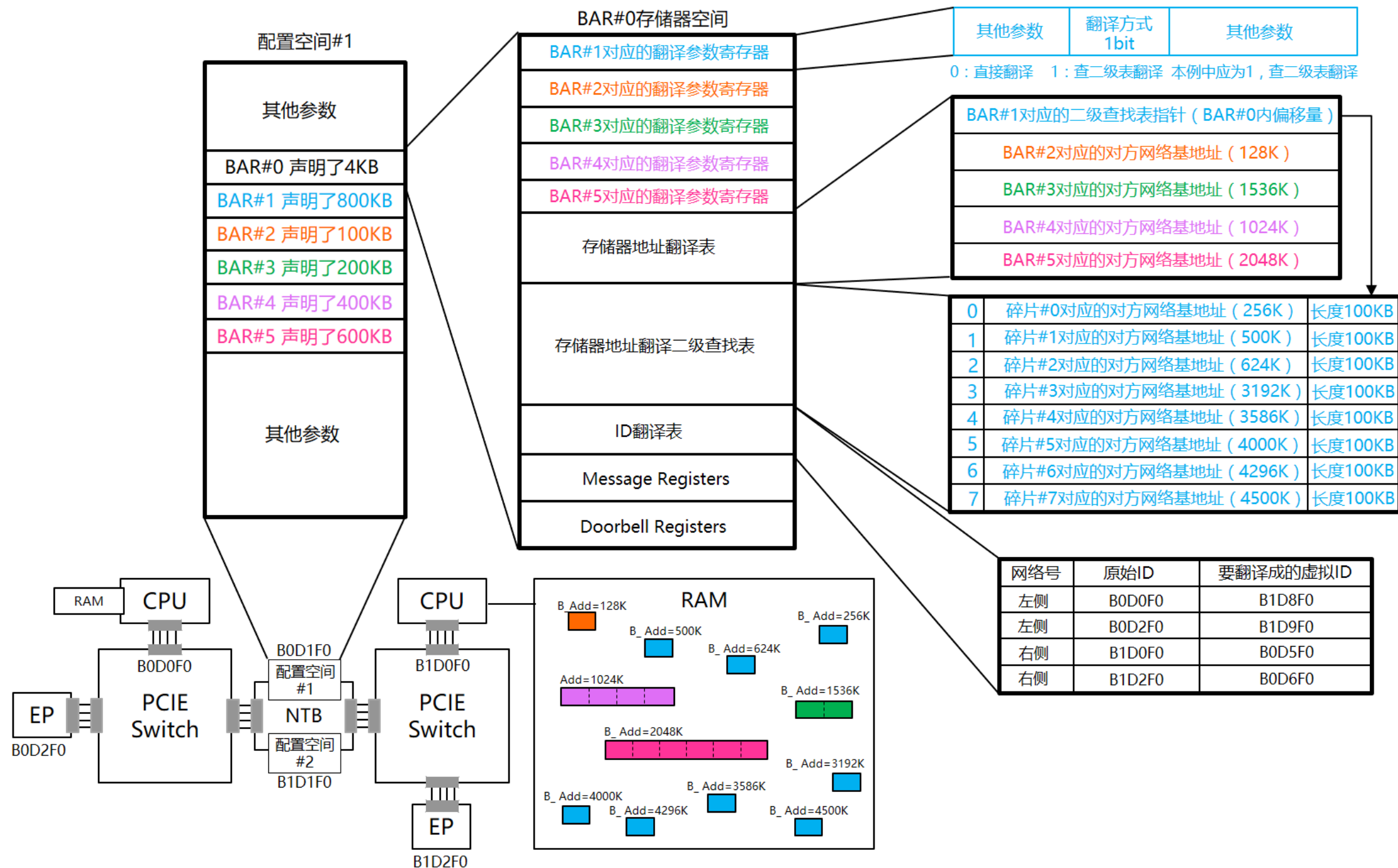
被选通的那个I/O控制器一侧的PCI控制器从总线上接收命令，根据寄存器号从配置寄存器堆中读出对应值并返回到总线

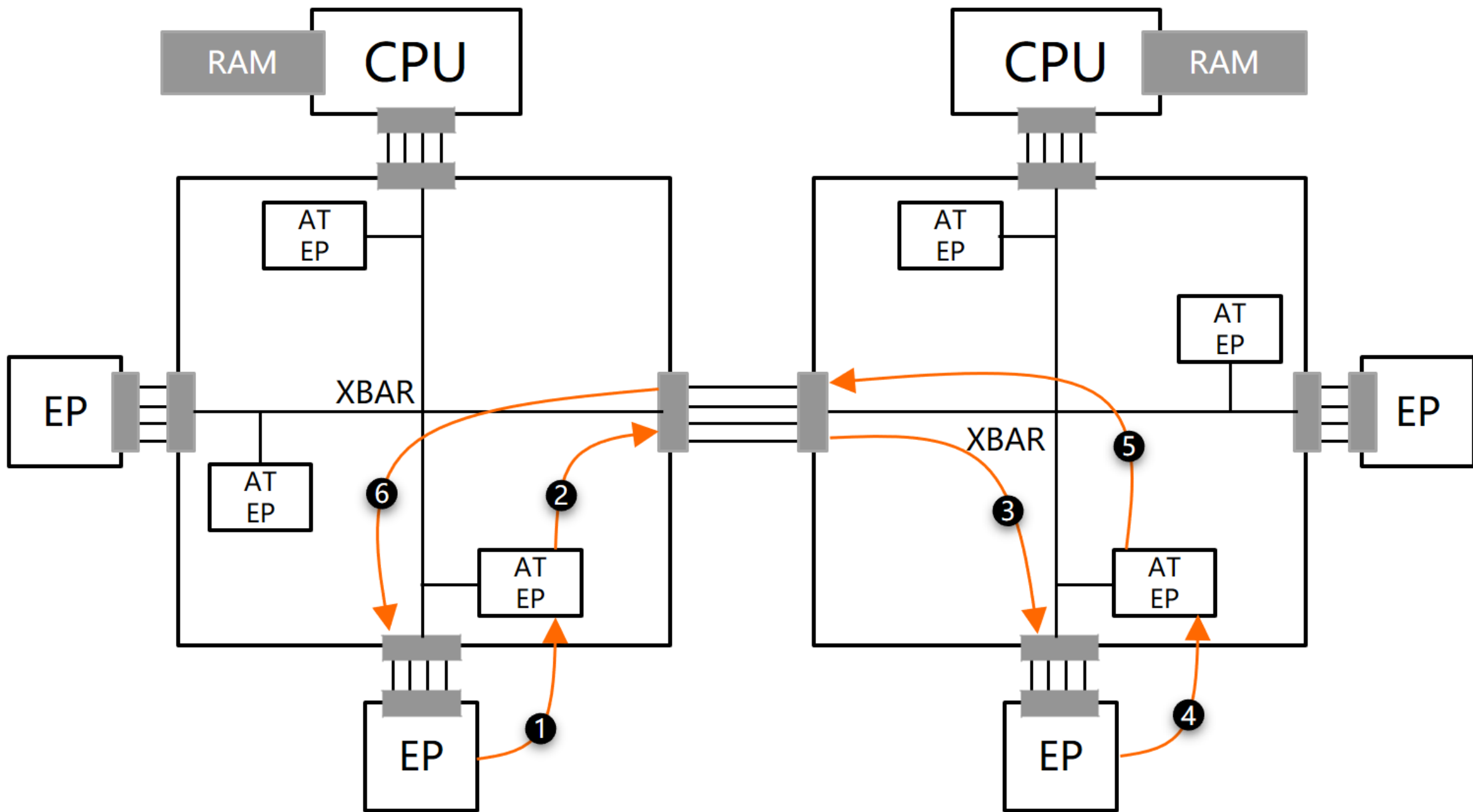
00010000000000000000 00000011000

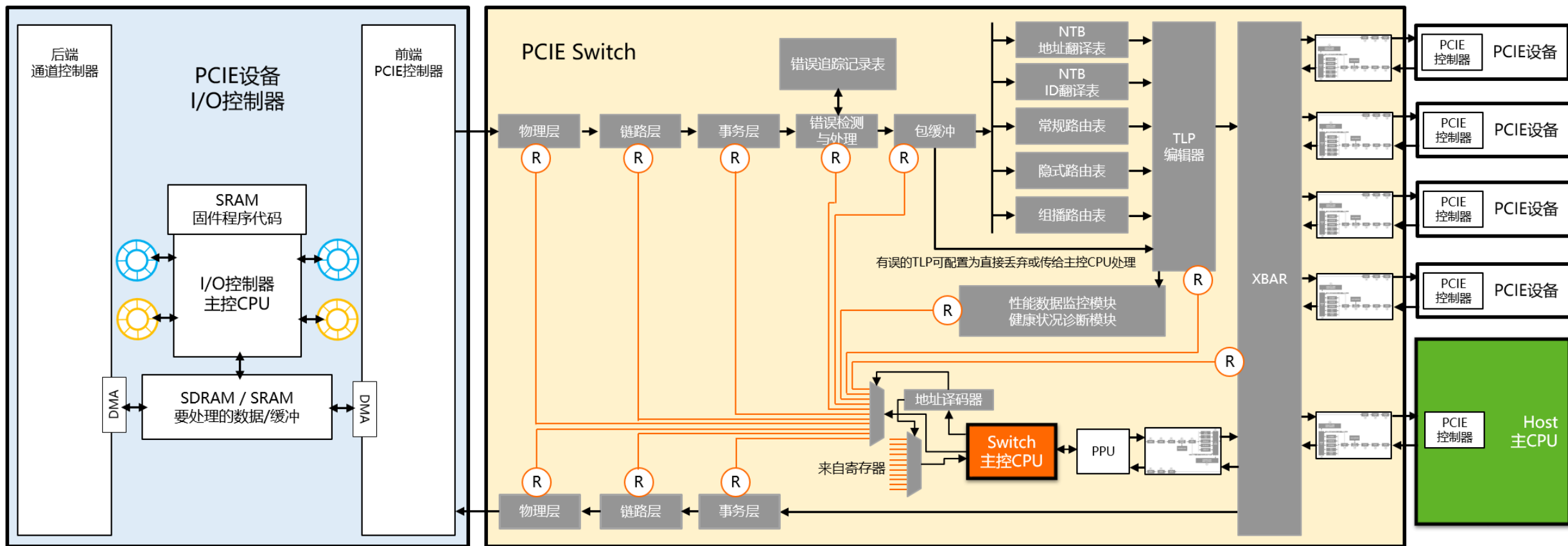














Byte			
3	2	1	0
Device ID		Vendor ID	
Status Register		Command Register	
Class Code			Revision ID
BIST	Header Type	Latency Timer	Cache Line Size
Base Address 0			
Base Address 1			
Base Address 2			
Base Address 3			
Base Address 4			
Base Address 5			
CardBus CIS Pointer			
Subsystem ID		Subsystem Vendor ID	
Expansion ROM Base Address			
Reserved			Capabilities Pointer
Reserved			
Max_Lat	Min_Gnt	Interrupt Pin	Interrupt Line

## 64字节的标准配置空间



## 最大4K字节的扩展配置空间



## 第8章

### 8.1 声音处理系统

#### 8.1.1 让蜂鸣器说话

#### 8.1.2 音乐是可以被勾兑出来的

##### 8.1.2.1 可编程音符生成器PSG

##### 8.1.2.2 音乐合成器

##### 8.1.2.3 FM合成及波表合成

#### 8.1.3 声卡发展史及架构简析

#### 8.1.4 与发声控制相关的Host端角色

#### 8.1.5 让计算机成为演奏家

#### 8.1.6 独立声卡的没落

### 8.2 图形处理系统

#### 8.2.1 用声音来画图

#### 8.2.2 文字模式

##### 8.2.2.1 向量文本模式显示

##### 8.2.2.2 用ROM存放字形库

##### 8.2.2.3 点阵文字显示模式

##### 8.2.2.4 Monochrome Display Adapter ( ...

##### 8.2.2.5 点阵作图与ASCII Art

#### 8.2.3 图形模式

##### 8.2.3.1 Color Graphics Adapter ( CGA )

##### 8.2.3.2 Enhanced Graphics Adapter ( EG...

##### 8.2.3.3 Video BIOS ROM的引入

##### 8.2.3.4 Video Graphics Array ( VGA )

##### 8.2.3.5 VGA的后续

##### 8.2.3.6 当代显卡的图形和文字模式

### 8.2.4 2D图形及其渲染流程

#### 8.2.4.1 2D图形加速卡PGC

#### 8.2.4.2 2D图形模型的准备

#### 8.2.4.3 对模型进行渲染

#### 8.2.4.4 矢量图和bitmap

#### 8.2.4.5 顶点、索引和图元

#### 8.2.4.6 2D图形动画

#### 8.2.4.7 坐标变换及矩阵运算

#### 8.2.4.8 2D图形渲染流程小结

#### 8.2.4.9 2D绘图库以及渲染加速

### 8.2.5 3D图形模型和表示方法

#### 8.2.5.1 3D模型的表示

#### 8.2.5.2 顶点的4个基本属性

### 8.2.6 3D图形渲染流程

#### 8.2.6.1 顶点坐标变换阶段/Vertex Transform

#### 8.2.6.2 顶点光照计算阶段/Vertex Lighting

#### 8.2.6.3 栅格化阶段/Rasterization

#### 8.2.6.4 像素着色阶段/Pixel Shading

#### 8.2.6.5 遮挡判断阶段/Testing

#### 8.2.6.6 混合及后处理阶段/Blending

#### 8.2.6.7 3D渲染流程小结

### 8.2.7 典型的3D渲染特效简介

#### 8.2.7.1 法线贴图 ( Normal Map )

#### 8.2.7.2 曲面细分与置换贴图 ( Tessellation

#### 8.2.7.3 视差/位移贴图 ( Parallax Map )

#### 8.2.7.4 物体投影 ( Shadow )

#### 8.2.7.5 抗锯齿 ( Anti-Aliasing )

#### 8.2.7.6 光照控制纹理 ( Light Mapping )

#### 8.2.7.7 纹理动画 ( Texture Animation )

### 8.2.8 当代3D游戏制作过程

### 8.2.9 3D图形加速渲染

#### 8.2.9.1 3D图形渲染管线回顾

#### 8.2.9.2 固定渲染管线3D图形加速

#### 8.2.9.3 可编程渲染管线3D图形加速

#### 8.2.9.4 Unified可编程3D图形加速

#### 8.2.9.5 深入AMD R600 GPU内部执行流程

### 8.2.10 3D绘图API 及软件栈

#### 8.2.10.1 GPU内核态驱动及命令的下发

#### 8.2.10.2 GPU用户态驱动及命令的翻译

#### 8.2.10.3 久违了OpenGL与Direct3D

#### 8.2.10.4 Windows下图形软件栈

### 8.2.11 3D图形加速卡的辉煌时代

#### 8.2.11.1 街机/家用机/手机上的GPU

#### 8.2.11.2 SGI Onyx超级图形加速工作站

#### 8.2.11.3 S3 ViRGE时代

#### 8.2.11.4 3dfx Voodoo时代

#### 8.2.11.5 NVidia和ATI时代

## 8.3 结语和期盼



### 8.1.1 让蜂鸣器说话

后来，个人电脑上逐渐使用了频响范围较宽的蜂鸣器，这使得其发声时可以带有更多谐振频率，让音质变得相对更柔润了一些。这样就可以驱动蜂鸣器发出任意单一频率或者任意频率叠加之后的杂波了，意味着，其可以发出人声！



图 8-1 蜂鸣器、宽频响喇叭以及第一个利用其发出中文语音的游戏

第一个用 PC 蜂鸣器发出中文语音的游戏程序，是 1991 年智冠科技开发的《三国演义》，张飞：“匹夫，出来与我决一死战！”，吕布：“汝何人也，不配与吾交手！”。在这之前的一两年内，国外的一些游戏已经开始可以用 PC 喇叭发出语音了。要知道在 25 年前，多数 PC 电脑用户是根本没体验过电脑发出人声的，那时候能够看到电脑屏幕上出现图形，而不是一串串的字符，就已经兴奋满足不已了，何况能听到声音？不过，在当时已经出现了独立声卡，只不过非常昂贵，一般只有发烧友级别的消费者才会去购买，所以能用任何 PC 都有的蜂鸣器发出声音便是当时让人难忘的记忆。用 PC 蜂鸣器发出人声，是前人们使用非常聪明的做法完成的。

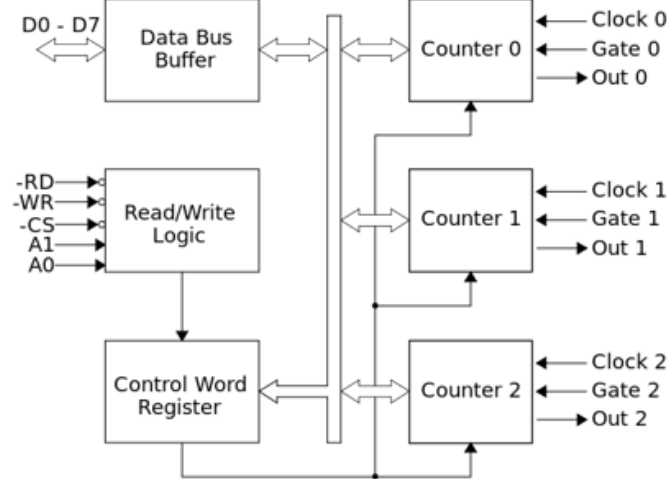


图 8-2 Intel 8235 时钟发生器架构图

当时 IBM 的 PC 兼容机广为流行，其中使用了 Intel 生产的时钟发生器，比如 Intel 8253/8254 型时钟发生器。该发生器的架构如图 8-2 所示。其中内含三个计数器，程序可以将一个数值写入到计数器中，然后计数器就开始自动倒计时，当数值降低到 1 之后，便会将对应的 Out# 信号置为高电平。利用这个时钟发生器，可以定时产生中断，从而让系统记录时间以及实现其他功能。其可以被配置为多种模式，程序可以向 Control Word Register 中写入对应的控制字来切换其实现各种模式。这些寄存器中的信号会控制下游电路中的逻辑从而切换到对应模式。比如其中一种模式就是，程序向 Counter2 中写入一个初始值，倒计时结束后电路会向 Out2 发出一个高电

平，然后回到低电平，然后继续载入（电路自动保存和载入）使用刚才的初始值继续倒计时，继续向 Out2 发出高电平，这样便会在 Out2 上产生一个恒定周期的方波脉冲，周期的大小可以通过 Counter2 的初始值来给定。

那么，这个时钟发生器与 PC 蜂鸣器又有何关系呢？原来，当时 IBM PC 普遍被设计为使用时钟发生器的 Counter2 的 Out2 输出端来驱动蜂鸣器发声，比如遇到了内存没有插，需要间隔 0.2ms 发出 3 声嘀声，每一声持续发声 0.5ms。这个过程怎么做到？根据上述原理，可以先向 Counter1 中写入一个初始值，只要按照该时钟发生器的运行频率，精确计算出经过 0.2ms 时间该计数器会变化多少次，那就写入对应的数值进去，到了时间自然会触发一次高电平，该高电平会触发一次中断，中断服务程序运行之后，可以向 Counter1 再写入一个能够持续计数 0.7ms 的值，让 Counter1 开始倒计时，然后立即向 Counter2 中写入一个能够持续极短时间的值，比如持续 0.001ms 的值，这样，Counter2 每隔 0.001ms 便会发出一个高电平，让蜂鸣器振动一下，在 0.5ms 内，蜂鸣器会振动 500 次，那么其振动频率也就是 1KHz，这就能够让人耳感受到声响了，所以人耳听到的就是一声持续 0.5ms 的长鸣。然后 Counter1 会发出中断，程序可以继续这个过程。

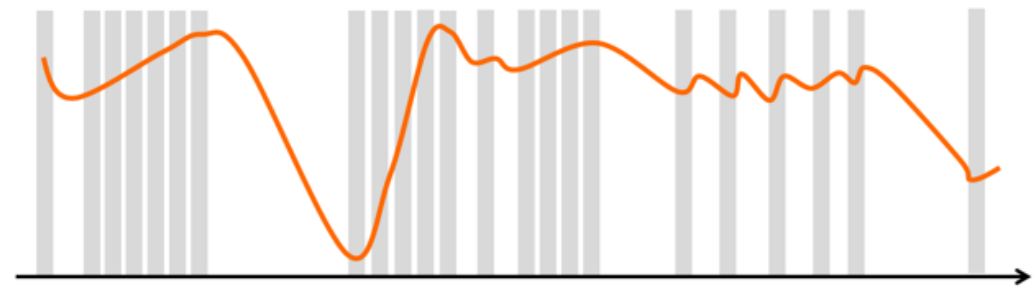
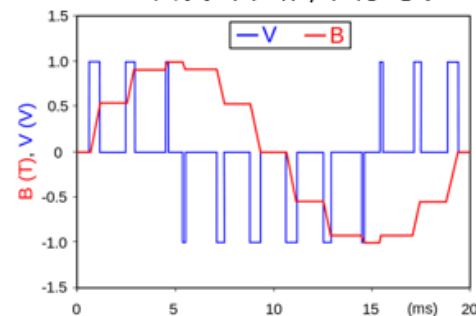
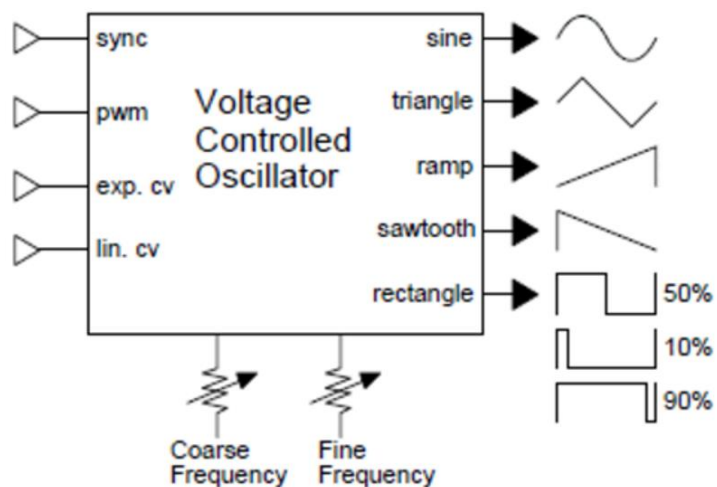


图 8-3 PWM 脉冲宽度调制示意图





Sine正弦波形



Ramp陡降波形



Triangle三角波形



Square方波形

- √ 声音、视频和游戏控制器
  - ➔ Realtek High Definition Audio
  - > 鼠标和其他指针设备
  - > 通用串行总线控制器
  - > 图像设备
  - > 网络适配器
  - > 系统设备
  - > 显示适配器
- √ 音频输入和输出
  - Microphone Array (Realtek High Definition Audio)
  - ➔ Speaker/HP (Realtek High Definition Audio)

设备管理器

文件(F) 操作(A) 查看(V) 帮助(H)

系统设备

- ACPI 盖子
- ACPI 固定功能按钮
- ACPI 热区域
- ACPI 睡眠按钮
- ➔ High Definition Audio 控制器
- Intel Collaborative Processor I
- Intel(R) Management Engine I
- Lenovo PM Device
- Microsoft ACPI 兼容的嵌入式控
- Microsoft ACPI-Compliant Sys

High Definition Audio 控制器 属性

常规 驱动程序 详细信息 事件 资源

High Definition Audio 控制器

设备类型: 系统设备

制造商: Microsoft

位置: PCI 总线 0、设备 27、功能 0

设备状态

这个设备运转正常。

吉他内部的合成器，合成器加入一些基本的效果比如回啊、六吗守，来模拟个吉他的六吗腔效果，然后输出最终的模拟信号给功放放大后输出到音响发声。当然，一般吉他内部不会带有较强功能的合成器，所以输出的音色比较单调，如果想要更多的音色和后期处理，需要另外制备独立合成器/效果器。有些高级产品利用两组拾音器产生方向相反、大小相等的两路电流，从而防止各种电磁干扰，这就与第 6 章中介绍的高速通信 PHY 底层的差分信号的本质是类似的。



图 8-17 电吉他上的拾音器

后来，人们采用了野路子，直接对各种乐器的声波采样并保存采样值，只要精度够高，就可以还原出原始波形。采样时只需要针对乐器发出的每个音调的一个大周期采样即可。所有乐器的所有音调的采样值被保存在一张表中，这个表被称为**波表**。波表保存在合成电路旁边一片 ROM 存储器中，当然也可以保存在计算机硬盘中，前者被称为**硬波表**，后者则为**软波表**。早期，由于存储器的成本非常高，人们只能降低采样精度来节省波表占用空间，带来的损失则是音质不高，后来一些合成器逐渐采用大容量存储器，随之而来的音质也就非常高了。有些合成器可以在运行时实时的通过 PCI 接口直接访问被载入 Host 端主存中的波表信息，从而节省卡上的 ROM。



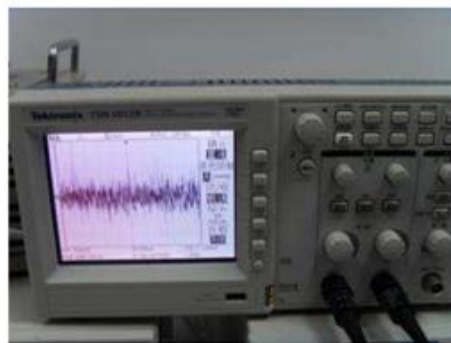
定要组装一台 4 路顶级显卡交火 (SLI/CrossFire) 的机器, 以最高特效, 4K 分辨率, 流畅运行所有游戏。



图 8-62 游戏《战争机器 4》(左)及《神秘海域 4》(右)截图

然而, 在欣赏这些画面的同时, 冬瓜哥心底一直藏着的那个大大的问号, 也终于忍不住爆发了: 计算机到底是怎么

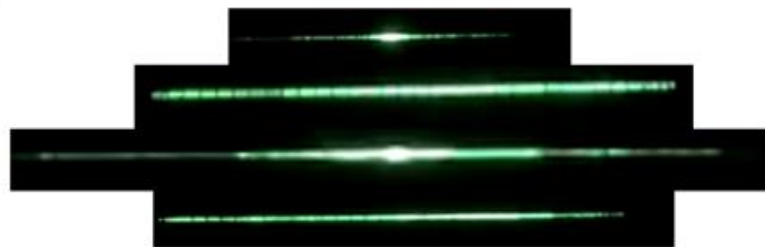
用你的音频播放器播放任意音乐，就可以了。



吐 8-70 用示波器展示任意波形

现在我们把玩一下，让 X 轴上不再安分分的加锯齿波，而是加上正弦波或者方波等，让屏幕内的这个波形世界的时间不再匀速流逝，而是扭曲着流逝，甚至可以回退！会发生什么？你可以先冥想一下。如图 8-70 右侧所示的波形。嗯？这镜头好像也不算火爆啊，冬瓜哥这个导演一般啊！

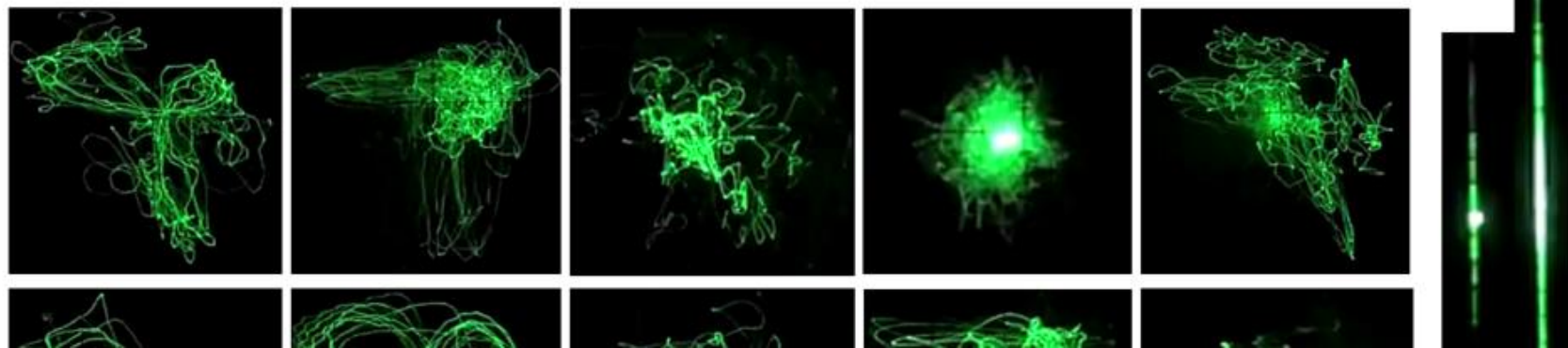
嗯，但是下面的剧情不会让你失望的。如图 8-71 所示为一名高手打算用示波器+声波信号画出蘑菇的图案。建议先不要扫码观看视频，先看完文字。



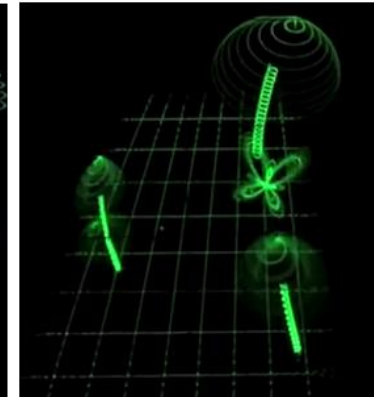
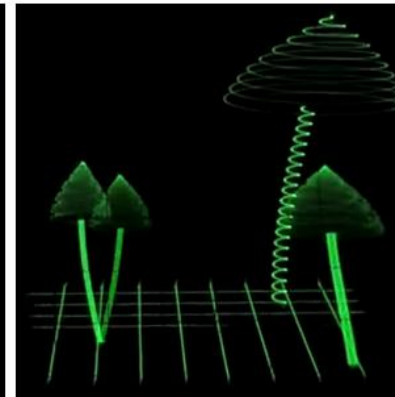
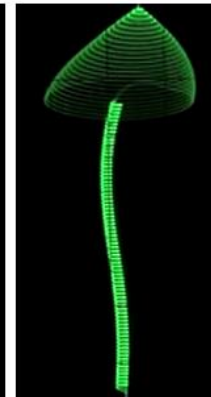
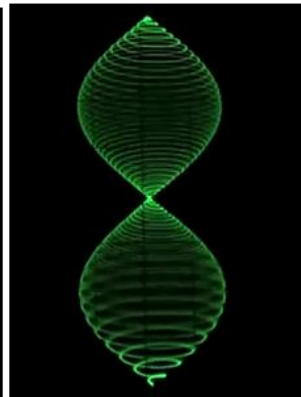
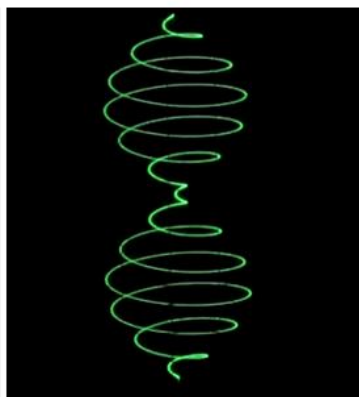
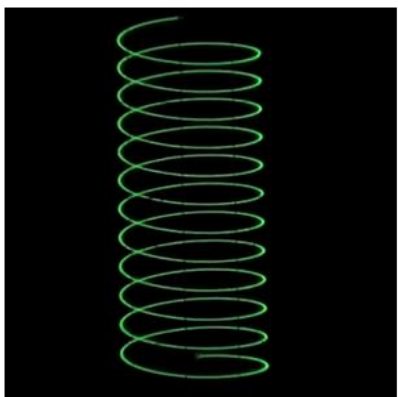
视频地址



声音文件下载地址







### 8.2.2.1 向量文本模式显示

人们是这样做的，如图 8-76 所示，比如对于字符“A”，人们先将其抽象的放入一个栅格中，然后将该字形跨越的栅格节点的坐标位置记录下来，取其中可以组成直线或者近似直线的若干个点，将这些点在栅格中的相对坐标值记录下来，并保存到某个地方，比如某个文件中，所有的可能的、常用的字符，都这样处理一番，最后形成的所有字符的形状样点坐标值所保存到的文件，被称为**字库**，或者**字符集**。字库/字符集可以有不同的设计风格，就像每个人写出来的字体都不同一样，有些人写的看上去就那么的端庄，而冬瓜哥写出来就无法让人直视。用户可以替换字库文件以实现装载不同风格字体的字库。字形库，就犹如声音波表一样，它们俩的本质目的是相同的。如图 8-76 所示为需要为向量显示器保存的字形库中的 A 字符的字形信息。

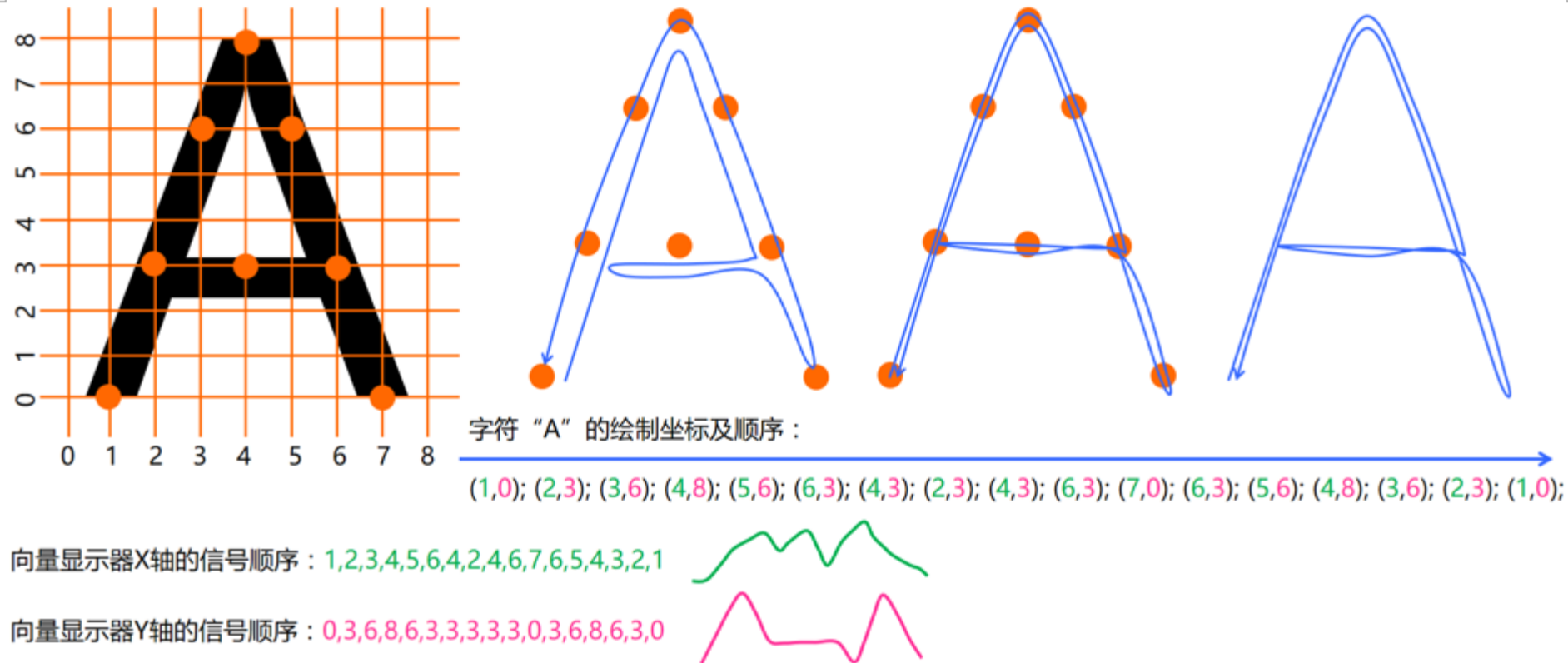


图 8-76 向量模式的字符“A”所需要保存的字形信息

(扫描线)。每个字符所占的方形位置又被称为一个 **Cell**，一横行 Cell 组成一个 **Character Line**。

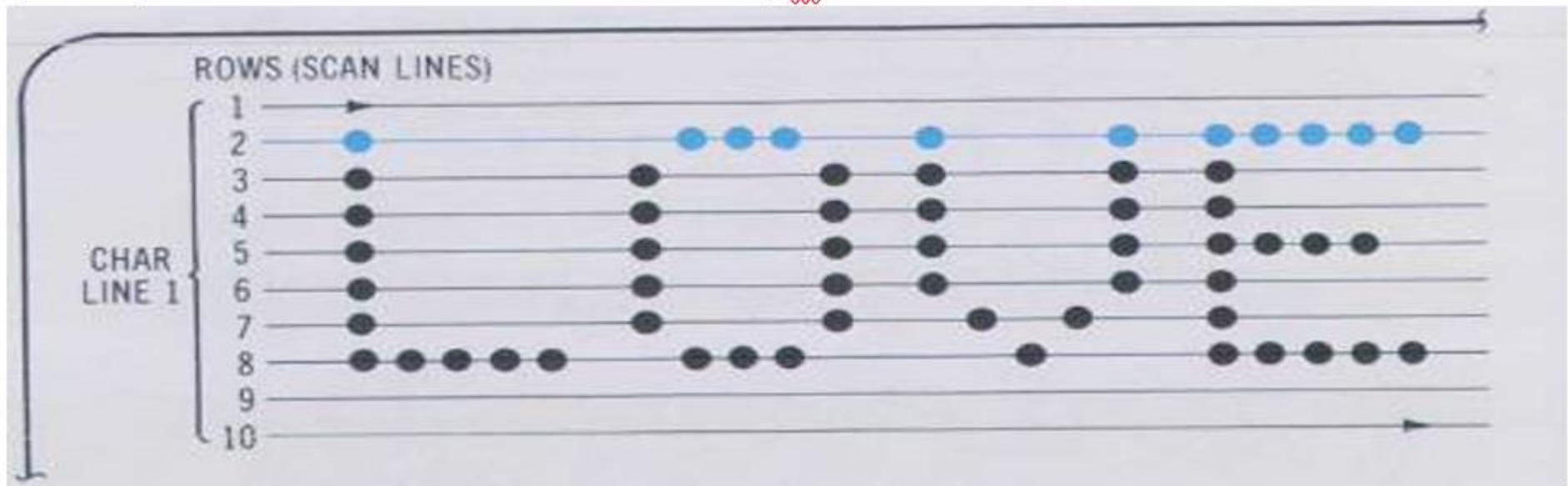
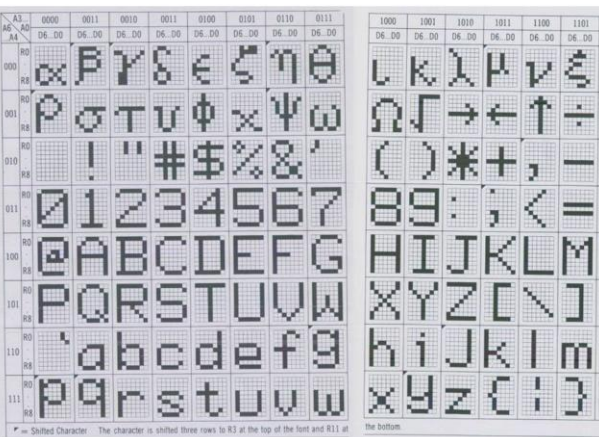
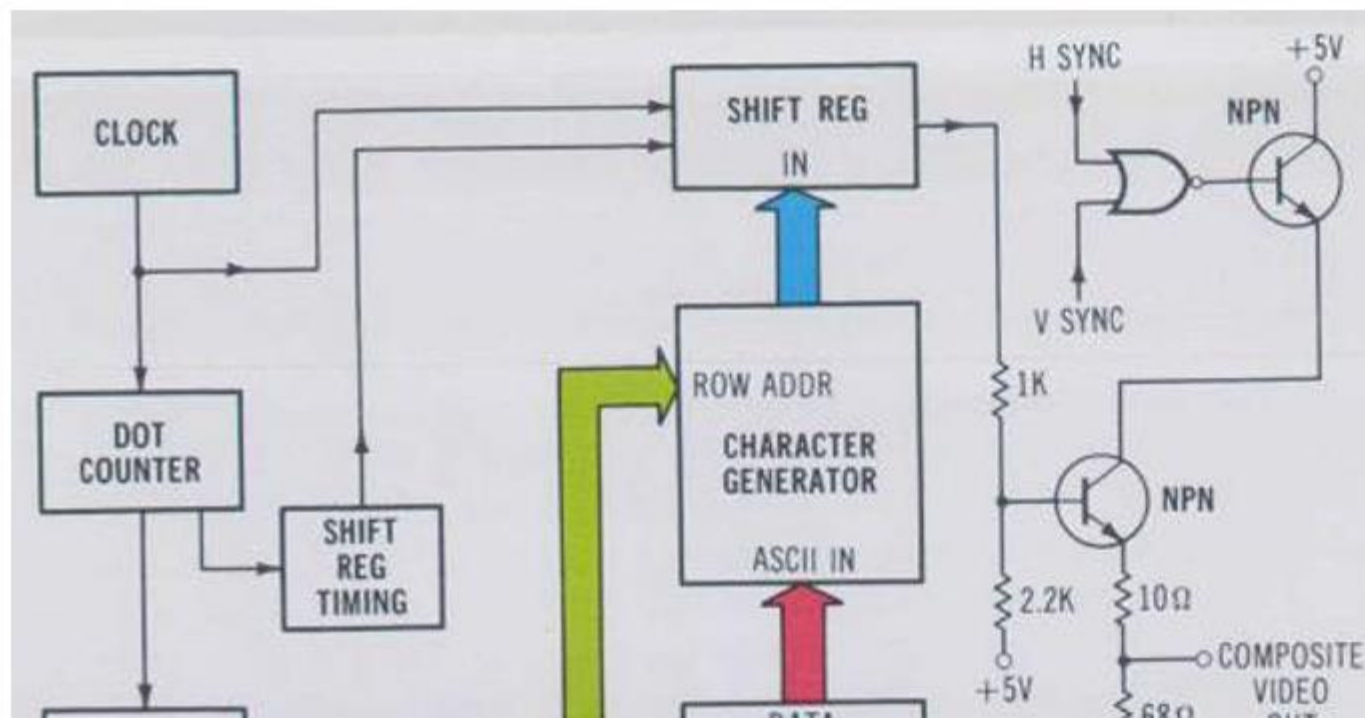


图 8-81 扫描一行会跨过多个字形



如图 8-82 所示为一个可以向像素显示器播放字形的完整的显卡架构实现。且看它的字形生成器，有两个主要输入，分别为从显存（图中 Display RAM）传送来的用于控制字形选择的 ASCII 码输入，以及从 Cell Line Counter 发来的用于控制让 ROM 输出对应字形的哪个 Scan line

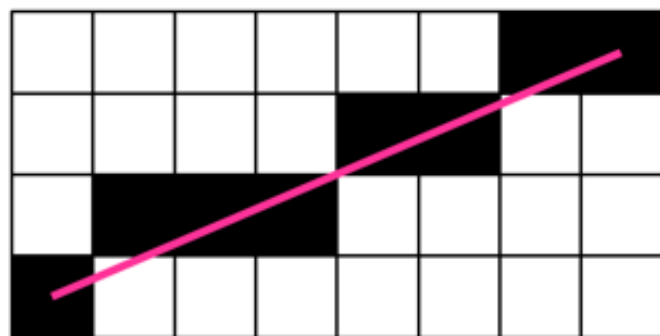


### 8.2.4.3 对模型进行渲染

现在想象你就是 `draw_line`(起始坐标, 结束坐标)函数的编写者, 你怎么根据这两个参数来确定屏幕上哪个像素组成这条线, 然后写入对应像素颜色值到 Frame Buffer 对应该坐标的字节? 还记得初中学过的, 给出两点坐标, 求该直线的表达式么? 冬瓜哥反正已经忘了。如果早知道今天要研究计算机作图, 当年一定不会忘。不过, 自己推导一下也能知道, 那就是先算出这个直线的斜率, 也就是  $m = \tan\theta = \text{两点纵坐标差} / \text{两点横坐标差}$ , 然后  $y = mx + a$ 。这样, 以起始坐标点开始, `for (x=起始坐标点, x<结束坐标点+1, x++) { 求  $y = mx + b$  }`, 即可求出该直线的所有纵坐标像素, 将其写入到对应的数组中, 也就确定了该直线的所有像素点, 这算是真正的图形“渲染”过程中最重要的一步, 相当于将图形的骨架勾勒出来, 然后就是对每个像素点上色和各种后期特效的加入。

给出顶点, 根据顶点位置计算连接顶点的线段所跨越的像素点的过程, 称为**插值 (Interpolation)**。

提示 实际上, 按照上述公式来计算的话, 底层函数需要做乘法, 而对数字电路来讲, 乘法的代价要高于加法。为此, 人们使用加法来计算 Y 坐标点。X 每次+1, y 每次增加的其实是 m, 那只需要将上一次计算所得的 y 值+m 即可, 而不需要用当前的 x 去乘以 m, 这就大大节省了计算量。这些函数底层做了大量的这类优化操作。当然, 对上层而言, 这些底层机制可能可能很少有人去关心了。



且慢! 由于我们使用的并非向量显示器而是像素/栅格显示器, 所以其每个像素是占有一定面积的, 所以直线跨越的像素点并不平滑, 如图 8-114 所示, 当一条直线跨越了两个像素点时, 到底该选哪个像素点呢? 这方面就对应了一些不同的算法了, 比如根据某种规则来对 Y 值取整数, 看与哪个 Y 坐标接近就选哪个。

确认对应图形所占用所有像素坐标的过程, 被称为**栅格化**

( **Rasterization** ) 有人称之为光栅化。因为目前的显示器每个栅格的



### 8.2.5 3D 图形模型和表示方法

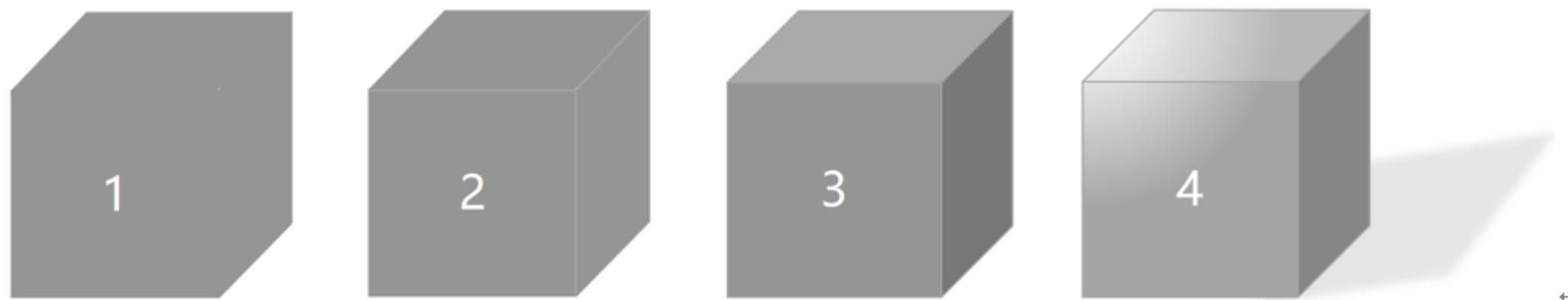
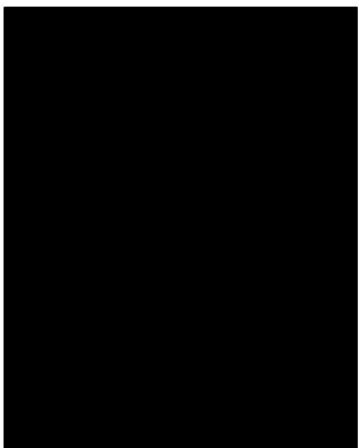


图 8-128 人脑对三维的理解

试想一下，人脑是如何感知三维世界的？如图 8-128 所示。如果挡住 2/3/4 号图形，只观察 1 号，你能判定它是个三维正方体么？不能。对于 2 号图形，多了三条线，产生了三个面，根据经验，已经可以让你判断出大概了。对于 3 号图形，不同的面的亮度不同，已经基本可以确定其意义了。4 号就更不用说了，左上角的光照和右下角的阴影说明了一切。其实，这四个图形在显示器平面上占用的像素面积、形状，是完全一样的。

如果仔细思考一下的话就能体会到，人脑接受的视神经的输入，永远都是 2D 的图像，也就是 3D 世界在人脑中的 2D 投影。人脑之所以理解了 3D，是因为观察一个物体的视角是可以变化的，你可以绕到该物体的侧面重新观察该物体，发现其 2D 投影也会跟着变化，经过对不同视角的投影进行组合识别分析，最终形成一幅 3D 图形的立体感观，所以人自从出生起就不断的接受到周围事物的各种角度投影的训练，形成了经验，才产生了理解 3D 世界的能力。

而对于 3/4 号图形，即便你不用去绕到它侧面去观察，只通过观察其不同位置对光线的反射率/明暗程度，根据经验，也可以判断出明显有一个 z 轴方向的存在，光线就是大自然赐予的参照物和测量尺。但是，设想一下，如果你从来没有尝试过绕到物体侧面观察过，假设你生来就是静止的，动弹不得，连眼珠都转动不了，你有可能只会认为该图形左上角的亮度高一些，右下角有个暗色斑而已，并不会认为该图形存在一个 z 轴方向，你根本不理解阴影为何会存在，产生不了立体感观，所以 4 号图形在你眼中只是一个正方形和两个菱形，和一个一角被挡住的躺着的菱形的组合罢了。



初始阶段



栅格化+环境光



环境光+顶点漫反射光照计算



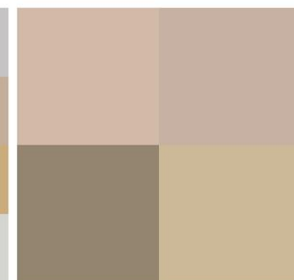
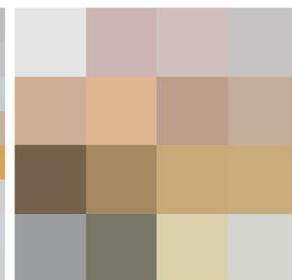
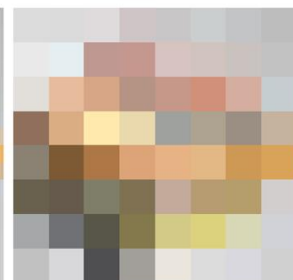
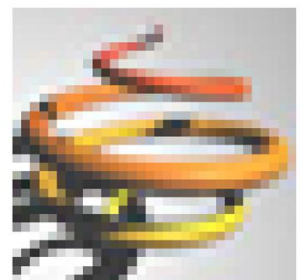
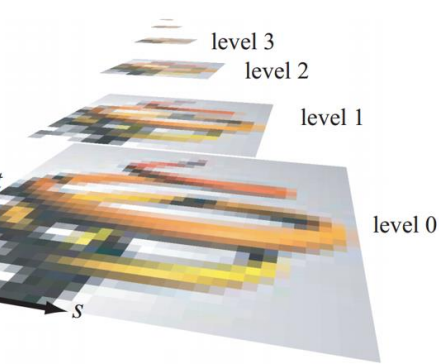
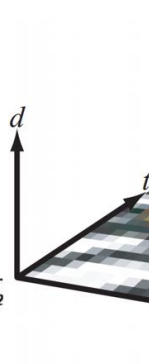
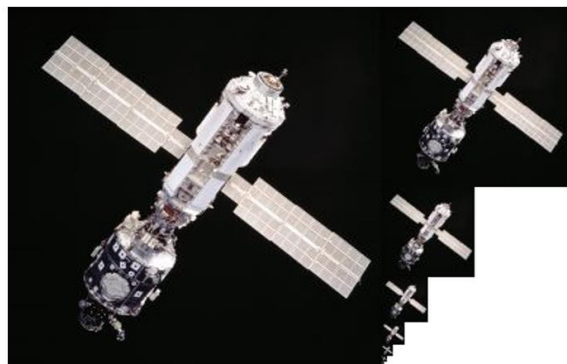
环境光+漫反射插值涂抹



光照+贴图后



灯光减弱



视差贴图的机制是，首先在法线贴图的基础上，像置换贴图一样引入高度图，但是并不用高度图去调制顶点坐标，然后利用高度信息，求出观察者在当前视角下应当观察到的是哪个纹素，然后将相应纹素值写到对应的屏幕像素中，从而形成对应位置上的纹素被抬高的假象。如图 8-215 所示。

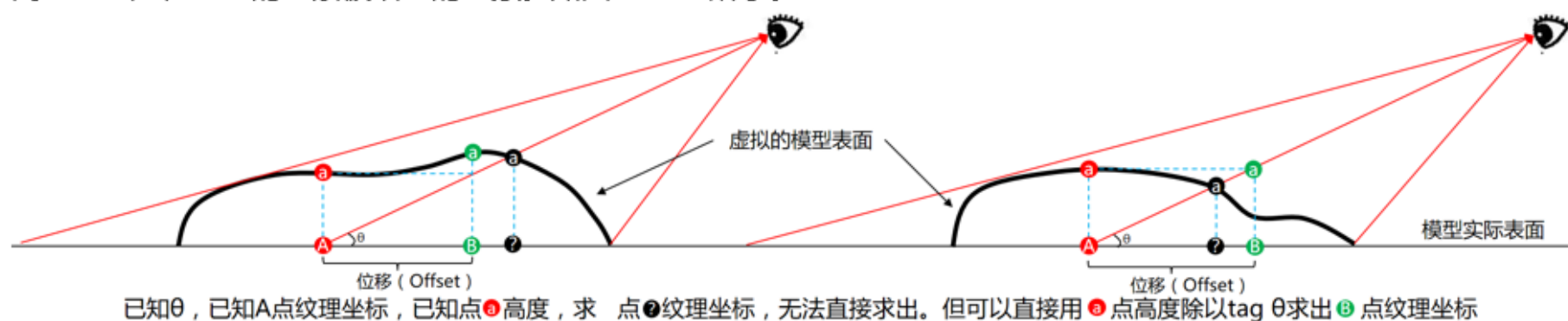


图 8-215 视差贴图基本原理示意图

先假设表面平整，没有高度信息，那么当观察者观察A点的时候，就需要对其呈现A点纹理坐标对应的纹素，这也是前文中所有贴图方式的做法。但是当要让观察者感觉到高度的时候，当观察者以同样的视角 $\theta$ 观察物体时，其应当观察到的是物体表面的a点，该点对应的纹素坐标在b点，显然，我们需要求出b点的位置，然后读出对应的纹素写

入到屏幕的A点，因为观察者最终还是看到屏幕上的A点，a是虚拟出来的，实际上不存在，但是只要将b的纹素贴到A点，观察者就会误认为a点被抬高了。这就是视差贴图的奇妙之处。这也是为何称之为视差/位移贴图的原因。

问题是，b点的坐标无法直接迅速的求出。给定任意一点A，其与观察者的连线与物体表面的交点位置，看似是固定的，一定能求出，但是，由于物体表面曲线是无法用方程来描述的，其完全没有规律，其被高度图中的值控制着，高度图是美工人员使用 2D 图像编辑软件，在屏幕上用鼠标刷出来的，就像往脸上抹粉一样，你不可能把你某次抹的粉的分布状况用一个公式描述出来。如图中右侧所示，如果物体表面高度变化了，那么b点位置也会变化，你根本不知道a点在哪里。找到光线与表面的交点位置的过程，被称为求交，对不规则表面求交的过程非常复杂。

但多数时候，纵使表面被抬高了，但是被抬高的地方有相当比例的位置依然是平坦的，那么这些位置就可以用公式来描述，那就是直线公式，当然，我们抛开晦涩的数学表示法不谈，转看图 8-216 所示。如果被抬高的表面平坦，那么求交操作的结果很容易计算出来，用a点高度（直接用A点坐标从高度图中读出）除以  $\tan \theta$ ，可得出b点与A点之间的 Offset 位移，然后再用A点坐标+Offset 就可得出b点的纹理坐标，然后将对应纹素读出贴入像素即可。



精度。

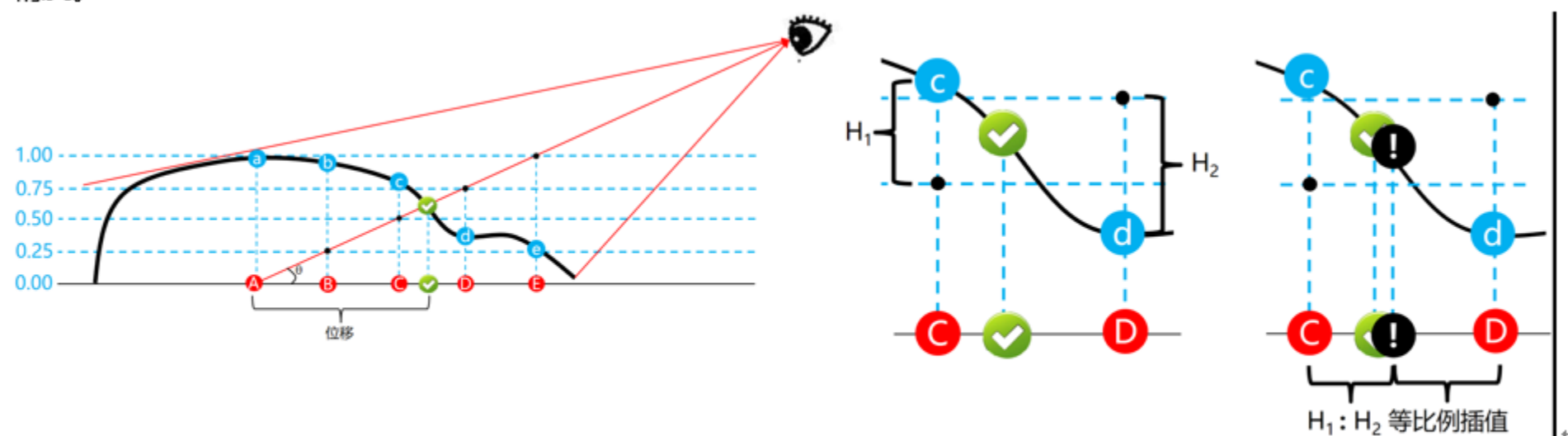


图 8-218 逐次逼近法求交示意图

虽然光线与假表面的交点无法直接求出，但是可以用笨办法来更进一步接近的求出。如图 8-218 所示，在假表面与真表面之间将整个高度切分为多个层级，然后将 A 点到观察者的光线分别与每一层的等高线进行求交，由于等高线为直线，所以很容易求得 A B C D E 各点的纹理坐标，然后从高度图中对应纹理坐标处读出这些纹素对应的假表面高度 a b c d e。显而易见的一个规律是，光线穿越表面的交点的两边的点，也就是本例中的 c 和 d，它们的表面高度值和层级高度值是相反的，也就是说，c 点的表面高度值大于其对应的层级高度值，而 d 点的表面高度值小于其对应的层级高度值。交点一定位于被这两个点夹住的表面区间内。所以，具体计算时并不需要将所有等高线对应的交点全求出，而是可以先求最高层级的点，考察其对应的表面高度值与层级高度值的关系，如果表面高度值小于层级高度值，则继续与下一个层级高度求交并做相同比较，直到找到发生反转（表面高度值大于层级高度值）的那个点，也就可以确定两个关键点了。或者采用二分法，先从中间层级高度开始求交计算，如果发现表面高度大于层级高度，则从再高一层的高度求交并计算、比较。

然后，利用这两个点各自的表面高度与层级高度之间的差值，用两者的比例去差值 C 和 D 两点间的距离，最终得





Packet Name	IT_OPCODE	Description
<b>Draw Packets</b>		
DRAW_INDEX_IMMD_BE	0x29	Draw packet used w 16 bit index data in
INDEX_TYPE	0x2A	Sends the current in
DRAW_INDEX	0x2B	Draw packet used w required
DRAW_INDEX_AUTO	0x2D	Draw packet used w
DRAW_INDEX_IMMD	0x2E	Draw packet used w packet
NUM_INSTANCES	0x2F	Sends the number o
MPEG_INDEX	0x3A	MPEG Packed Reg
<b>State Management Packets</b>		
SET_CONFIG_REG	0x68	Write Register Data
SET_CONTEXT_REG	0x69	Write Render State
SET_ALU_CONST	0x6A	Write ALU Const
SET_BOOL_CONST	0x6B	Write Boolean Con
SET_LOOP_CONST	0x6C	Write Loop Const
SET_RESOURCE	0x6D	Write Resource Co
SET_SAMPLER	0x6E	Write Sampler Con
SET_CTL_CONST	0x6F	Write Control Cons
SURFACE_BASE_UPDATE	0x73	Inform the CP whic around for CP on R
<b>Wait/Synchronization Packets</b>		
MEM_SEMAPHORE	0x39	Sends Signal & Wa
WAIT_REG_MEM	0x3C	Wait Until a Regist
MEM_WRITE	0x3D	Write DWORD to ?
CP_INTERRUPT	0x40	Generate Interrupt f
SURFACE_SYNC	0x43	Synchronize Surfac
COND_WRITE	0x45	Conditional Write t
EVENT_WRITE	0x46	Generate an Event v

**第 1 步。**上层绘图程序首先创建一个 batchbuffer，在 Host 端，程序的一切行为都要符合规范，申请任何内存空间都要像内存管理程序模块申请，可以将这个过程封装成函数，每次申请时调用即可。这个函数，GPU 厂商为您开发！实际上 Intel 开发了一整套 Graphic Execution Management ( GEM ) 函数库，就是为您提供方便的，谢谢！申请 batchbuffer 的函数名为 gem\_create( )，向该 buffer 中填入命令时，您可以调用 gem\_write( )。这里我们创建一个名为 gem\_exec 的 batchbuffer 并向其中填入大批 MI\_NOOP 命令然后结尾跟着一条 MI\_BATCH\_BUFFER\_END 命令。我们的最终目标是把这个 batchbuffer 的指针通告给 GPU 让它读取执行。↵

**第 2 步。**任何命令，不管是单条，还一批，都必须被放入一个 Descriptor 中，这个套路我们在第 7 章已经介绍的足够充分了。所以，在这一步中，绘图程序需要填个表，把这个 Descriptor 描述出来。GEM 函数库中，该描述表的格式遵循 execbuffer2 结构体（至于为什么会带一个 2 这个冬瓜哥真的不知道）。程序只需要复印这张表，并给其命名为 execbuf，并向其中填入需要执行的命令，如果是批量命令则要将上一步生成的 batchbuffer（名为 gem\_exec）填入，然后填入一些其他的控制信息，比如从 batchbuffer 中的哪个位置开始执行指令（batch\_start\_offset）等，最终这张样式为 execbuffer2 的、名称为 execbuf 的表，需要被传送给 GPU 的驱动程序。↵

**第 3 步。**上层绘图程序调用 drmioctl( )将刚才填好的&execbuf（还记得&符号的作用么？在 C 语言中他表示某个对象所在位置的地址指针）表传递给位于内核态的 GPU 驱动程序。Drmioctl( )函数等其实是执行了 ioctl 系统调用（回顾第 5 章的 5.5.6.4 一节，系统调用的更多细节我们将在第 11 章中介绍），从而将信息传入内核态。↵

**第 4 步。**位于操作系统内核态的 i915 显卡驱动（i915.ko 程序模块，Linux 操作系统下的内核态的设备驱动程序扩展名为.ko）对应的代码 i915\_gem\_do\_execbuffer( )被系统调用激活，从而将第 2 步生成的 execbuf 表格拿到手。↵

**第 5 步。**驱动程序将该命令描述表整理好之后直接 Dispatch 到 Ring Buffer 中，当然，派发过程中其实对应着很多细节步骤，比如修改 Ring Buffer 的控制指针，以及更新 GPU 一侧用于记录 Ring Buffer 首尾指针的寄存器等，这些都由驱动自动完成。整个命令下发过程宣告完成，后续就是 GPU 的渲染流程了。↵



- 第9章
  - 9.1 科学计算到底在算些什么
    - 9.1.1 蛋白质分子的故事
      - 9.1.1.1 氧气运输的故事
      - 9.1.1.2 更复杂的生化逻辑是如何完成的
    - 9.1.2 如何模拟蛋白质分子自折叠过程
    - 9.1.3 如何将模拟过程映射为多线程并行计算
    - 9.1.4 其他科学计算场景
  - 9.2 大规模系统共享内存之向往
    - 9.2.1 UMA/NUMA/MPP
    - 9.2.2 OpenMP并行编程
  - 9.3 基于消息传递的非共享内存系统
    - 9.3.1 采用消息传递方式同步数据
    - 9.3.2 MPI库基本函数简介
    - 9.3.3 MPI库聚合通信函数简介
  - 9.4 超级计算机
    - 9.4.1 IBM蓝色基因
      - 9.4.1.1 中央处理器CPU
      - 9.4.1.2 计算节点和I/O节点
      - 9.4.1.3 三个独立网络同时传递数据
      - 9.4.1.4 蓝色基因Q的网络控制和路由实现
      - 9.4.1.5 节点卡及整机架布局
      - 9.4.1.6 整体系统拓扑
      - 9.4.1.7 Service Node
      - 9.4.1.8 Service Card
      - 9.4.1.9 时钟同步
      - 9.4.1.10 系统启动
      - 9.4.1.11 软件安装与用户认证

- 9.4.1.12 状态监控
    - 9.4.1.13 计算任务的执行
    - 9.4.1.14 操作系统
    - 9.4.1.15 Login Node
    - 9.4.1.16 存储系统
  - 9.4.2 圣地亚哥Gordon
  - 9.4.3 Fujitsu PrimeHPC FX10
    - 9.4.3.1 SPARC64 IXfx CPU
    - 9.4.3.2 水冷主板
    - 9.4.3.3 Tofu六维网络互联拓扑
    - 9.4.3.4 ICC互联芯片
- 9.5 利用GPU加速计算
  - 9.5.1 Direct3D中的Compute Shader
  - 9.5.2 OpenCL和OpenACC
  - 9.5.3 NVidia的CUDA API
    - 9.5.3.1 CUDA基本架构
    - 9.5.3.2 一个极简的CUDA程序
    - 9.5.3.3 CUDA对线程的编排方式
    - 9.5.3.4 GPU对CUDA线程的调度方式
    - 9.5.3.5 CUDA程序的内存架构
    - 9.5.3.6 基于CUDA的PhysX库效果
- 9.6 利用PLD和ASIC加速计算
  - 9.6.1 PAL/PLA是如何工作的
  - 9.6.2 CPLD是如何工作的
  - 9.6.3 FPGA是如何工作的
  - 9.6.4 FPGA编程及应用形态
  - 9.6.4 Xilinx FPGA架构及相关概念
  - 9.6.6 ASIC与PLD的关系
- 9.7 小结：软归软 硬归硬



### 9.1.3 如何将模拟过程映射为多线程并行计算

那么上述过程具体应该如何映射到多个线程并行运算？比如，可以以原子为单位，每个线程负责计算每个原子在 10 飞秒后将移动到哪个位置，也就是空间坐标。当然，需要先初始化好对应的数据结构，比如每个原子都用一张表来追踪它们的各种属性，包括：相邻原子的表的指针、本原子的空间坐标值、初速度、当前速度、当前受力值等等。

每个线程的入口函数的输入值是周边原子的作用力，由于原子间力是短程作用力，所以距离较远的原子作用力就可以忽略，一般只考虑其化学键链条上的 3 个原子对其的作用力，超过 3 个以后的认为作用力为 0。输入值还包括该原子的初始三维坐标位置以及该原子的速度（初始速度为 0 或者某一既定速度）。输出值则是在牛顿力学公式的作用之下经过 10 飞秒加速之后该原子的新的空间坐标和速度矢量，并将这些输出值记录在各自的表中。

该线程中又会有大量函数相互作用，比如：有函数会专门根据当前原子的化合价以及与其化合的其他原子都是谁，然后计算静电力，有的函数则负责计算氢键力、范德华力等，这些力的计算公式有些非常复杂，需要较大的运算量。最后，求合力  $F$ ，算出初速度为 0，合力  $F$ ，质量  $m$ ， $t=10\text{fs}$  之后的该原子的位置和速度，最后这一步相信高中物理及格的朋友都可以算出来了。

那么，第一个 10 飞秒的模拟运算结束之后（注意，并不是说运算过程持续了 10 飞秒，而是说算出原子运动 10 飞秒后的速度和坐标，这个运算过程耗费的时间远大于 10 飞秒）。应该怎么办呢？当然是要把计算完的值更新到每个原子对应的结构体的对应项目中。然后呢？当然是每个线程都需要从本原子相邻的其他原子（不超过 3 级）的记录表中取出它们各自的空间坐标，然后根据各种力计算出本原子在新的位之下所受的新的合力矢量，然后继续开始运算再一个 10 飞秒之后本原子的坐标和速度值了。就这样，以 10 飞秒为步进一直向前推进。

这个过程中会产生一个潜在问题，在多核心计算机系统中，多线程在时间上是物理并发执行的，如果某个线程运算的较快，先结束了，而其相邻原子对应的线程运算较慢尚未结束，那么如果先结束的线程直接去读相邻原子的记录表取出坐标，那么取出的将会是旧值，从而导致运算错误。必须实现一种方式，让相互依赖的线程之间形成一种等待关系，只有所有这些线程都执行到某个步骤之后，才能继续下一步的执行，这种方式被称为屏障，Barrier。这个概念我们在第 6 章介绍访存的时空一致性时中初步涉及过。互斥锁、屏障，都是多线程同步的方式，各自应用场景也不同。

**MPI\_Bcast(void \*buffer,int count,MPI\_Datatype datatype,int root,MPI\_COMM\_WORLD)**。该函数用于线程将数据广播给通信域内的所有除了发送方线程之外的节点。参数 root 表示发送广播的源线程标识。发送方调用 MPI\_Bcast(), 接收方也必须在同一位置调用 MPI\_Bcast。嗯? 接收端为什么也要发起广播? 不是的。接收端调用该函数的作用是为了接收广播, 该函数会判断当前调用者线程的标识是否与参数 root 相同, 如果相同, 则表明当前线程是发出数据, 如果不同, 则表明当前线程之所以调用该函数是为了接收 root 线程发来的广播。这就是 MPI 编程方式的特点之一, 也就是发送方和接收方永远都要配合起来, 一唱一和。如果有人先唱到了这一段, 那就阻塞等待, 等对方也唱到这一段之后, 继续往下唱。下面的函数如无特殊说明, 均为发送方和接收方同时调用。MPI\_Bcast() 函数底层其实也调用了 MPI\_Send() 和 MPI\_Recv() 函数, 其本质就是将向所有节点的单播过程封装为广播。同理, 下面介绍的所有聚合通信类函数, 都是用基本的 MPI\_Send() 和 MPI\_Recv() 函数封装而成的, 你也可以自己封装一个聚合通信类函数出来。如图 9-16 所示为 MPI\_Bcast() 的执行过程示意图。

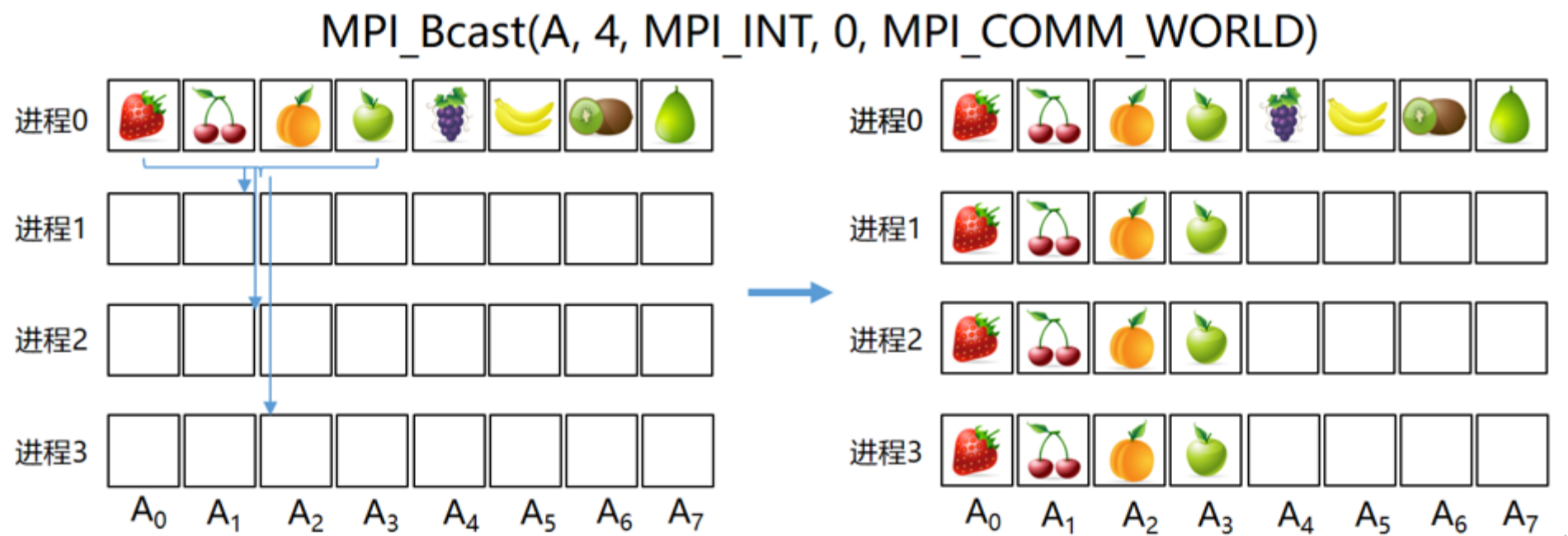


图 9-16 MPI\_Bcast() 的执行过程示意图

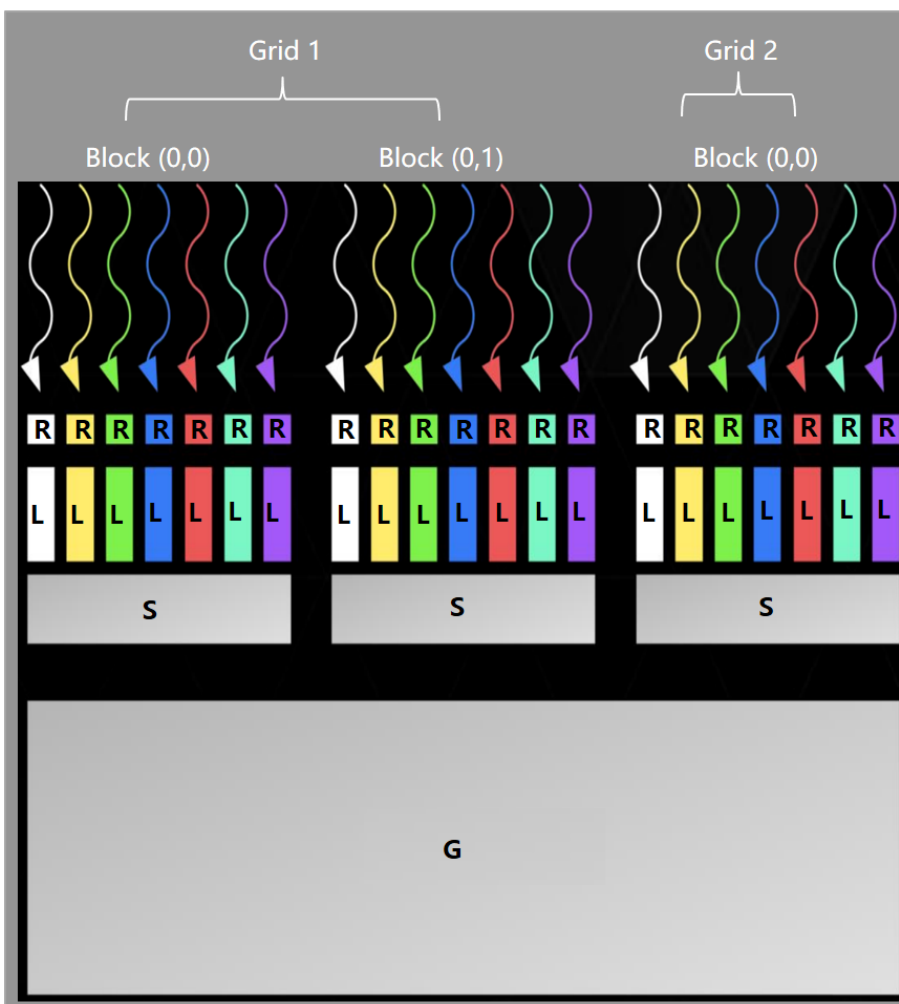
### ▪ 9.5.3.2 一个极简的 CUDA 程序

下面我们通过一个极简单的 CUDA 程序来向大家展示一下 CUDA 程序的基本编写方式。

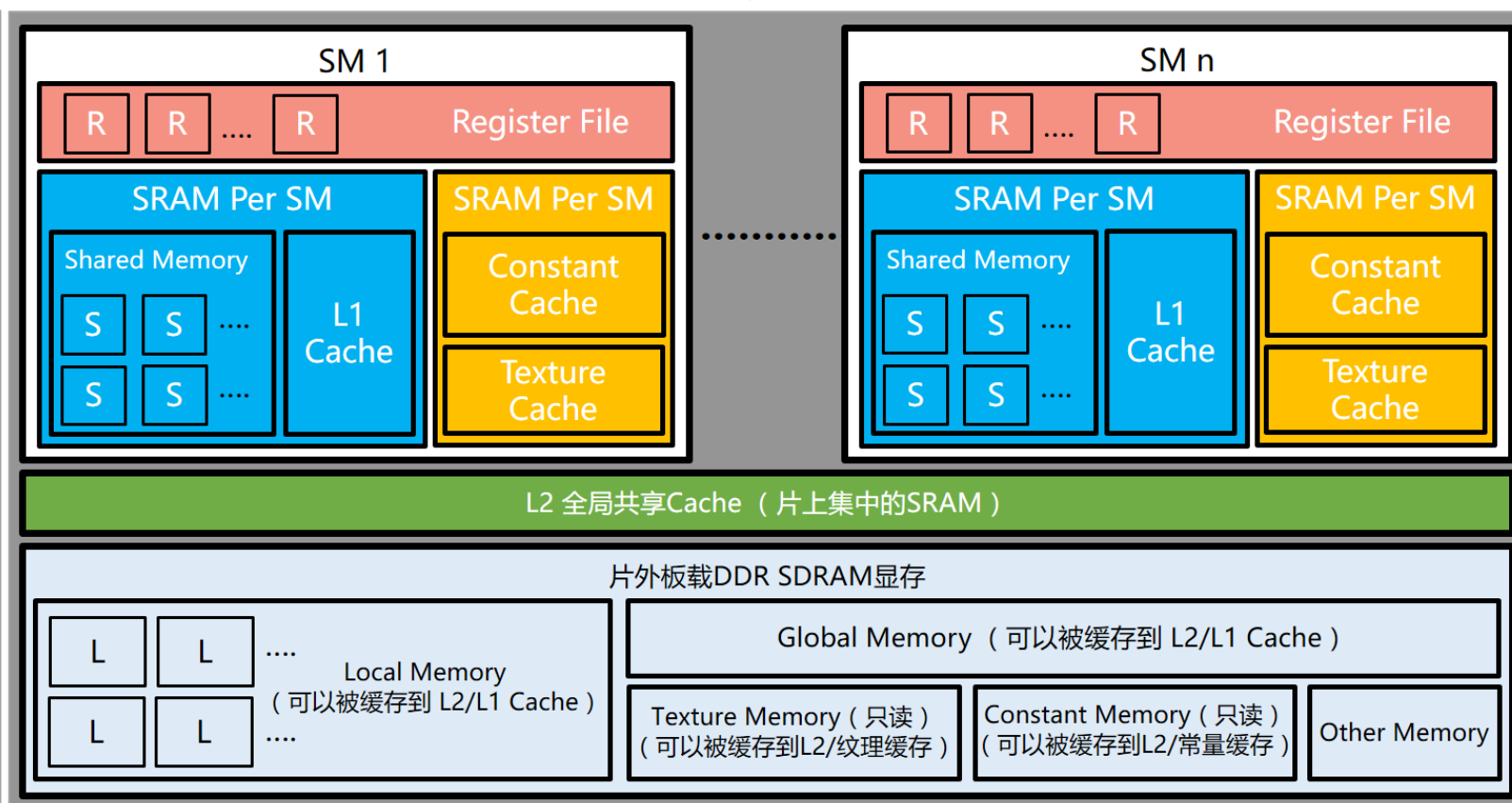
```
1.  #include<stdio.h>
2.  #define N 8 //代码中以 N 代替 8
3.  __global__ void add(int *a,int *b,int *c) //声明一个函数，该函数逻辑需要在 GPU 上执行
4.  {
5.      int tid=blockIdx.x;
6.      if(tid<N)
7.          c[tid]=a[tid]+b[tid];
8.  }
9.  int main() {
10.     int arr1[N],arr2[N]; int sum[N]; //定义三个数组
11.     for(int i=0;i<N;i++) {arr1[i]=i; arr2[i]=i+1;} //对数组进行数据填充
12.     int *a; int *b; int *res; //定义三个指针用于保存下面分配好的显存的指针
13.     cudaMalloc((void**)&res,sizeof(int)*N); //分配 N*4 KB 显存用于盛放两个数组之和
14.     cudaMalloc((void**)&a,sizeof(int)*N); //分配 N*4 KB 显存用于盛放数组 a
15.     cudaMalloc((void**)&b,sizeof(int)*N); //分配 N*4 KB 显存用于盛放数组 b
16.     cudaMemcpy(a,arr1,sizeof(int)*N,cudaMemcpyHostToDevice); //向显存拷贝数据
17.     cudaMemcpy(b,arr2,sizeof(int)*N,cudaMemcpyHostToDevice); //向显存拷贝数据
18.     add<<<N,1>>>>(a,b,res); //调用核函数，定义 N 个 Block，每 Block 包含 1 个线程，并行执行该函数
19.     cudaMemcpy(sum,res,sizeof(int)*N,cudaMemcpyDeviceToHost); //把结果从 res 拷贝回到 sum
20.     cudaFree(a); //通知 GPU 释放显存
21.     cudaFree(b); //通知 GPU 释放显存
22.     for(int i=0;i<N;i++) {printf("%d\n",sum[i]);} //显示计算结果
23. }
```



逻辑视图



物理视图



R : Register File , 位于GPU内的SM内 , SRAM形态 , 访问延迟1周期 , 每个线程独享各自的一份。  
 S : Shared Memory , 位于GPU片上 , SRAM形态 , 访问延迟数个周期 , 多个Block均分容量 , 每个Block内所有线程共享。  
 L : Local Memory , 位于GPU片外板载显存中 , DDR SDRAM形态 , 访问周期约500周期 , 每个线程独享各自的一份。  
 G : Global Memory , 位于GPU片外板载显存中 , DDR SDRAM形态 , 访问周期约500周期 , 所有Grid中线程共享。

### 9.5.3.6 基于 CUDA 的 PhysX 库效果

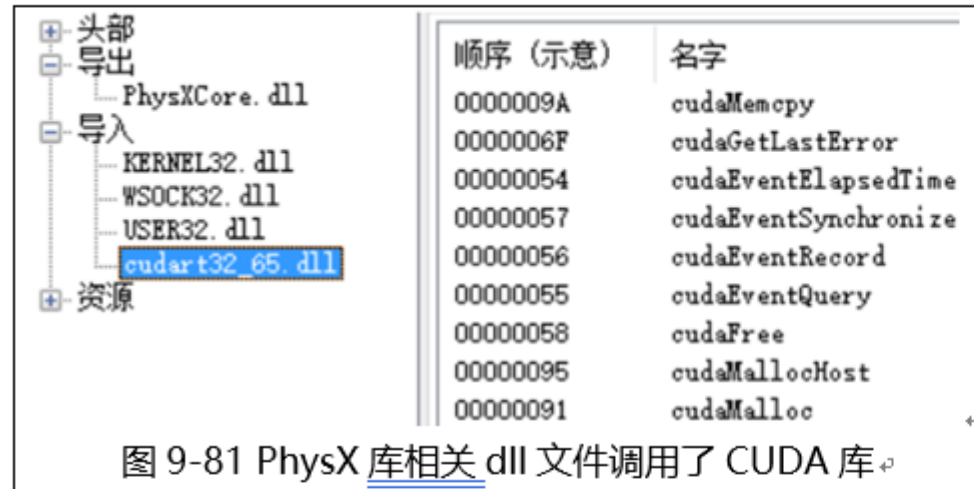


图 9-81 PhysX 库相关 dll 文件调用了 CUDA 库

基于 NVidia GPU 的 CUDA 加速计算越来越流行，在各行各业得到了广泛应用。Nvidia 自家搞的物体物理碰撞计算加速库 PhysX 本身就是基于 CUDA 实现的（如图 9-81 所示）。游戏程序调用 PhysX 库中的物理碰撞模拟算法，将模型的顶点坐标位置利用 CUDA 加速计算出来，返回到 Host 主存，然后游戏再走 D3D 渲染流程，将顶点新位置被计算好的模型渲染出来。物理碰撞特效非常适合使用并行计算来加速，比

如有一千颗冰晶，游戏中的角色朝某个方向释放后，每颗冰晶都有个初速度和模拟阻力（甚至可以引入当时的环境风速），然后针对每一颗冰晶计算其下一个时间单位（比如如果决定该物理特效动画流畅度为每秒 30 帧，那么每个时间单位就是 33ms）所在的位置，所使用的算法再简单不过，就是牛顿运动定律，根据合力求时间  $t$  后的位移量。如果前方有遮挡物，则根据位移量和遮挡物坐标算出发生碰撞的时间点和位置（这个过程被称为**碰撞检测**），然后在该时间点针对该粒子做反弹计算，利用的公式，我想任何一个高中毕业的朋友都做过小球以一定入射角撞墙反弹的题，不再赘述（其实冬瓜哥已经忘了，尴尬一下）。经过大量运算之后，这些冰晶颗粒就按照物理规律运动，产生震撼的物理模特

域 2》对 PhysX 达到了“滥用”的程度，对于当年的显卡当然会力不从心，但是时隔数年之后，可以特效全开 2K 分辨率下达到 70 帧以上，冬瓜哥甚是欢欣。写到这里，冬瓜哥不由自主的又打开了 PC 机载入《圣域 2》进去享受了一下里面的风景，以及令人心旷神怡的物理特效。如图 9-82 所示，在 GTX980 显卡强悍算力的支撑下，人物脚下的火圈在地上留下了按照真实物理相互作用的火星拖尾；火圈发出的热风将地上的石子、火星等向外吹；大量的（据冬瓜哥观察起码有数千个）魔法冰晶颗粒被爆发出来碰到场景中的物体上反弹并相互碰撞然后按照物理规律运动掉落并堆积在地上；风刮起地上数百片落叶并随风飘动下落；人物走路驱赶起地上的石子和落叶，所有石子/树叶都会有投影，所有的石子落叶都可以按照物理规律自旋转。而如果关闭 PhysX 特效，石子、落叶和冰晶都不会出现。《圣域 2》一度让冬瓜哥感觉没有物理特效的游戏都索然无味。

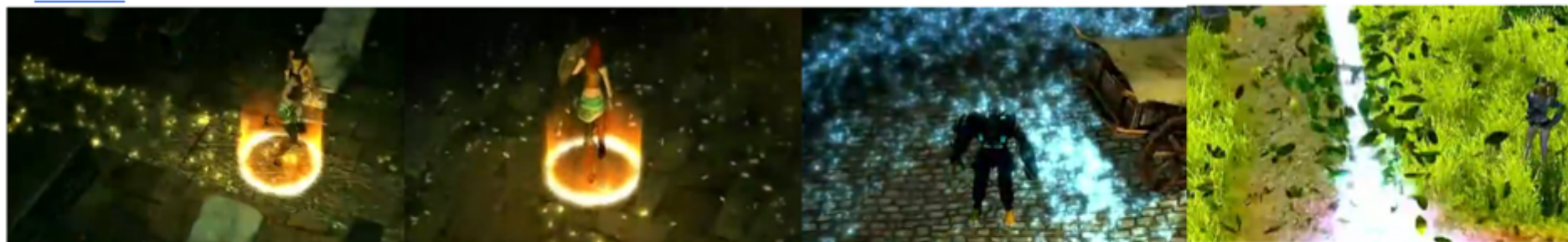


图 9-82 PhysX 物理加速库以及游戏《圣域 2》中的 PhysX 特效截图

值得一提的是，PhysX 可以被配置为使用 CPU+软件方式来计算，但是在冬瓜哥的 Intel 酷睿 i7 CPU 下，使用软计算运行圣域 2，当放出冰晶魔法时，已经不是每秒几帧的帧率了，而是几秒一帧，此时显卡基本处于闲置状态，CPU 使用率 100%，然而，CPU 就算忙的满头汗也赶不上显卡弹指一挥间。游戏画面成了名副其实的幻灯片速度。而且软算法考虑到 CPU 的算力，对物理效果做了大幅简化，无论是粒子数量还是排布的凌乱程度上，真实度差了一大截。



### 9.6.3 FPGA 是如何工作的

	a	b	c	out
a →	0	0	0	0
	0	0	1	1
b →	0	1	0	0
	0	1	1	1
c →	1	0	0	0
	1	0	1	1
	1	1	0	1
	1	1	1	1

(a AND b) OR c

图 9-93 用存储器模拟逻辑电路体现真值表

与 CPLD 采用与/或门阵列将多个乘积项相或的方式来实现逻辑的方式不同的是，FPGA 采用另一个奇妙的解决办法。我们在第 1 章中介绍过，任何（注意，是任何，没有之一也没有例外）逻辑电路都可以表达成真值表对外进行响应，输入值与输出值严格一一对应。可以这样认为：只要实现了这个真值表，就可以完成该电路的功能，而实现这个真值表并不一定非要用门电路，可以用存储器。比如，当输入值为 0000 时输出 10，输入为 0001 时输出 11，那么可以准备 16 行、每行 2bit 的存储器，然后用 4bit 输入值作为行号索引，将对应行

号的输出值保存在该行的 2bit 上，每次得到某个输入值，就到对应行号读出输出值输出，这也可以真实体现真值表中的逻辑，而不需要任何逻辑门电路。如图 9-93 所示为使用 4 字节的存储器实现了一个 (a AND b) OR c 的逻辑真值

表。当然，可以使用容量更大的存储器实现（或者说存储）更复杂的逻辑。用于存放真值表的存储器被称为 LUT (Lookup Table)。

利用 LUT 查表来“计算”数据时，其过程并不是“计算”，而是“查询”，也就是说，LUT 中保存的数据，其实都是预先被计算好的，相当于把一块逻辑电路的所有结果展开到了 LUT 中。如果使用逻辑电路来搭建这个逻辑，则只需要一个双输入与门加一个双输入或门，如果换算成开关的话（不考虑实际工程设计），只需要与门两个开关和或门两个开关共 4 个开关即可实现。而 LUT 实现相同逻辑所耗费的晶体管数量显然要大于逻辑门。



为了实现高速查表，LUT 一般使用 SRAM

## 9.6.4 FPGA 编程及应用形态

基于 PLD 器件进行开发逻辑，首先确定该逻辑的功能，然后将其翻译成 FPGA LUT 内的真值表。这个过程如果全靠人脑，将会非常复杂。为此，人们开发了**硬件描述语言（Hardware Description Language, HDL）**。开发者先把逻辑功能描述成 HDL，然后由 HDL 编译器自动分析并将这些逻辑编排落地到 FPGA 中。

比如“有个逻辑模块，它有 3 输入 8 输出，当输入为 xxx 时输出为 xxxxxxxx；当输入为 xxx 时输出为 xxxxxxxx；.....”，我们知道这是在描述一个 3-8 译码器。很显然，这种描述方式中并没有出现任何对逻辑门的描述，它只是在描述这个电路的行为是什么。这种 HDL 描述方式被称为**行为级描述**。

而如果是这样来描述：“有这么个器件，有 xxx 这些输入和输出信号，其内部有个驱动能力 x 延迟 x 的双输入与门连接了 x 和 x 及 x 信号，其内部有个驱动能力 x 延迟 x 的双输入异或门连接了 x 和 x 及 x 信号，.....”。其被称为**门级描述**。

在行为级和门级之间还存在一种叫做**寄存器传输级描述（Register Transfer Level, RTL）**的中间描述方式。行为级描述层次太高太抽象，代码非常精简，有种你咋不上天的感觉。而门级描述非常底层（其实还有更底层的开关级描述，因为门电路是由开关组成的），代码量非常大，写起来不方便。而 RTL 级描述则成为比较主流的描述方式，其对行为级描述的更加具体，但是却又避免谈及底层的电路结构，其只聚焦在数据在进入组合逻辑和时序逻辑电路时的更具体的行为。

上述描述方式都属于 HDL。目前有两种比较流行的 HDL，分别为 VHDL 和 Verilog HDL。如图 9-104~106 所示为一些简单的行为级和门级描述的例子。

```
1 module q_decode_38(data_in,data_out);
2
3   input[2:0] data_in;      //端口声明
4   output[7:0] data_out;
5   reg[7:0] data_out;
6
7   always@(data_in)
```

```
1 module ActiveStructure_mux2(datain_A, datain_B, sl, dataout);
2   input datain_A,
3   datain_B;
4   inputs;
5   output dataout;
6   reg dataout;
7
```

```
1 module TwoMux(datain_A,datain_B,sl,dataout);
2   input datain_A,
3   datain_B,
4   sl;
5   output dataout;
6
7   not U11(nsl,sl);
```

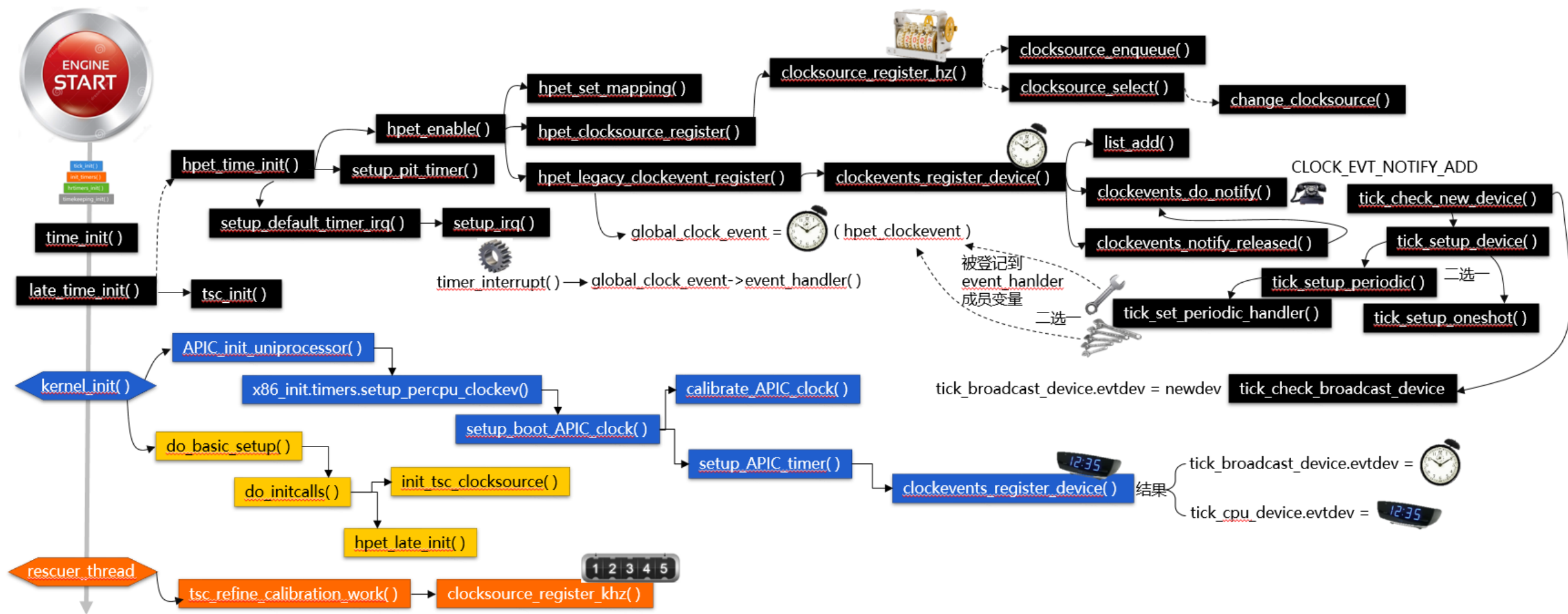
# 《大话计算机》第10章部分图示

20 19 18 17 16 15 14 13 12 11 10 9 8 7 6 5 4 3 2 1 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20

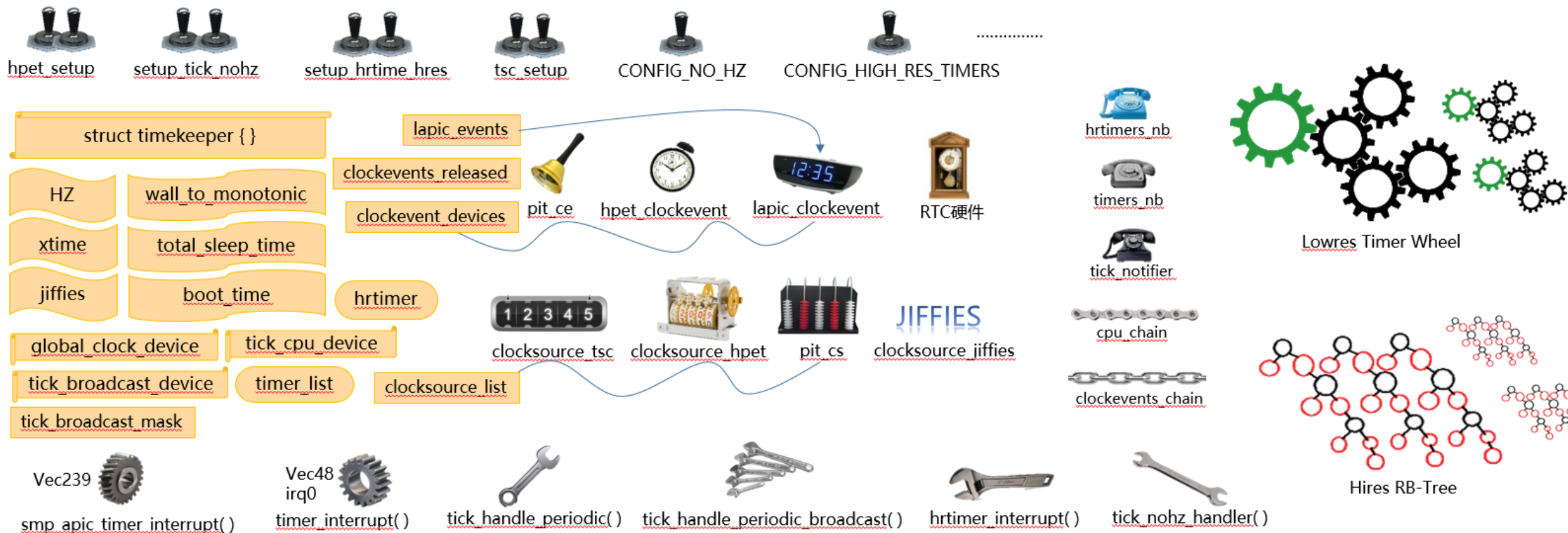




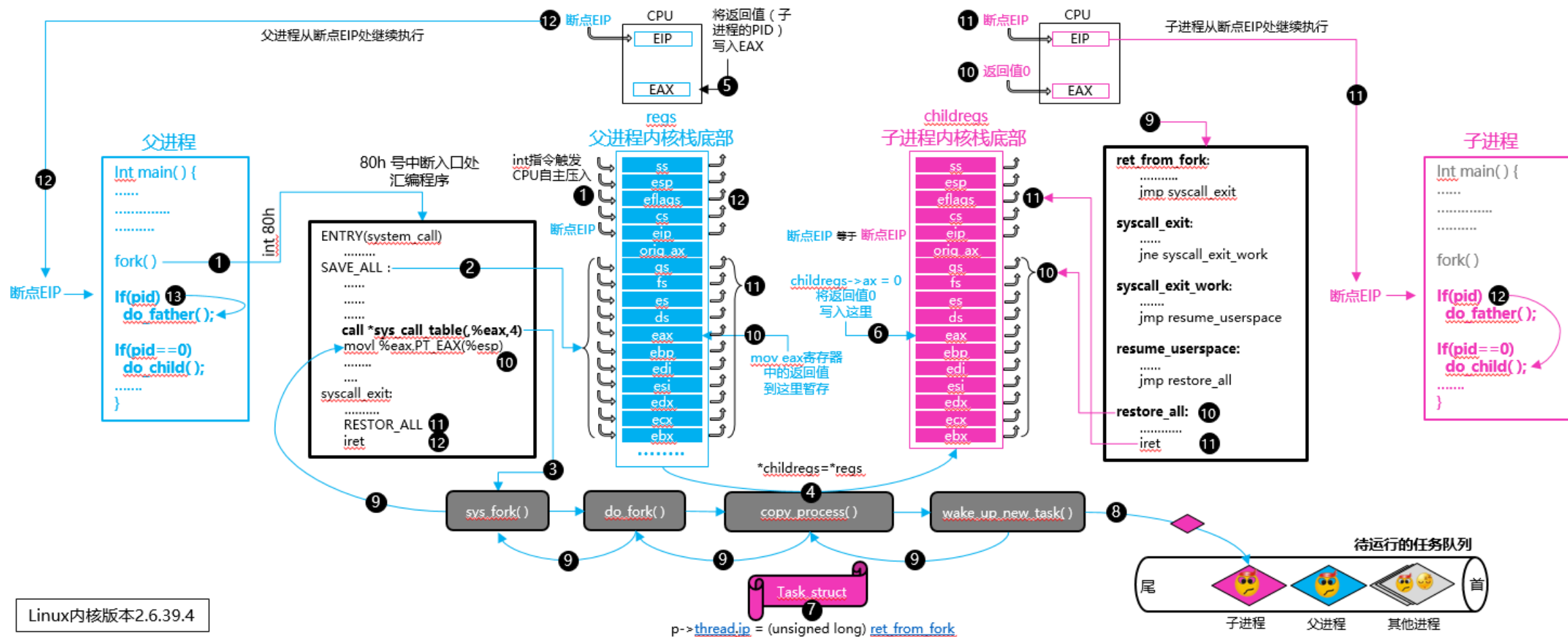
# 《大话计算机》第10章部分图示



# 《大话计算机》第10章部分图示

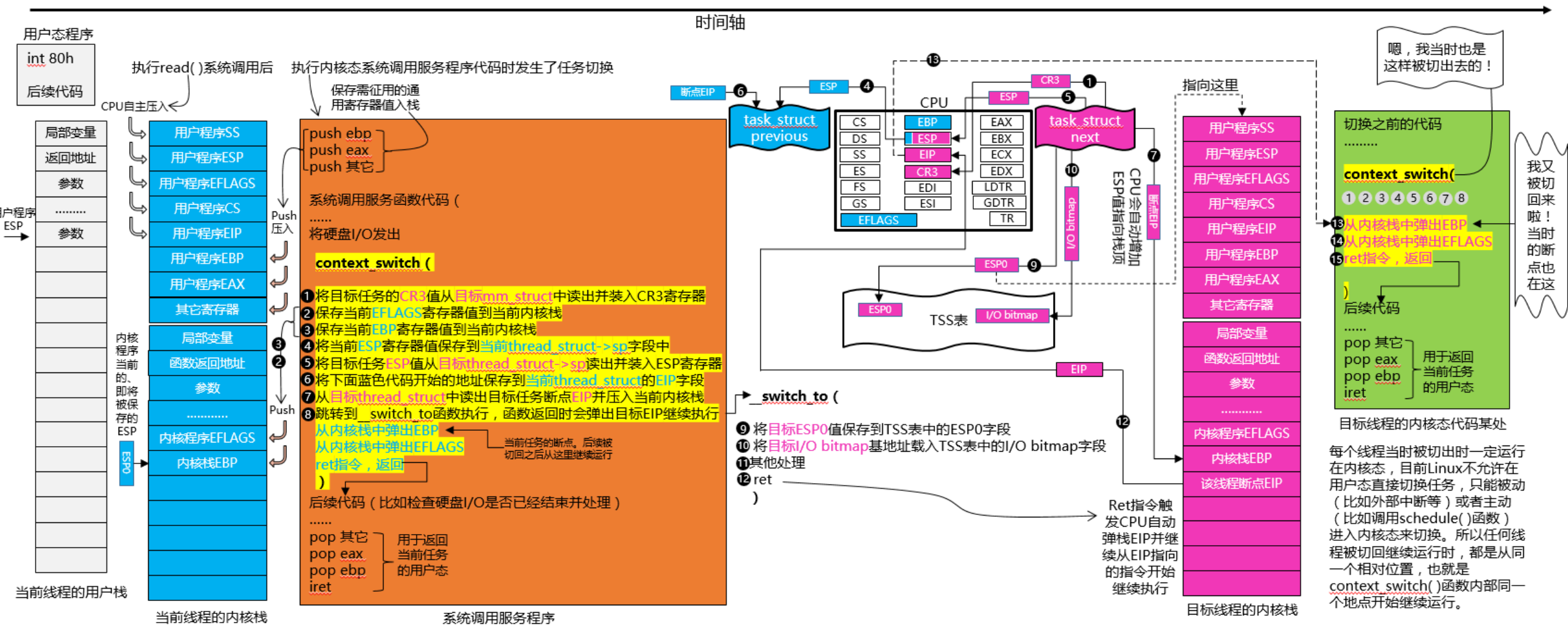


## 《大话计算机》第10章部分图示





# 《大话计算机》第10章部分图示

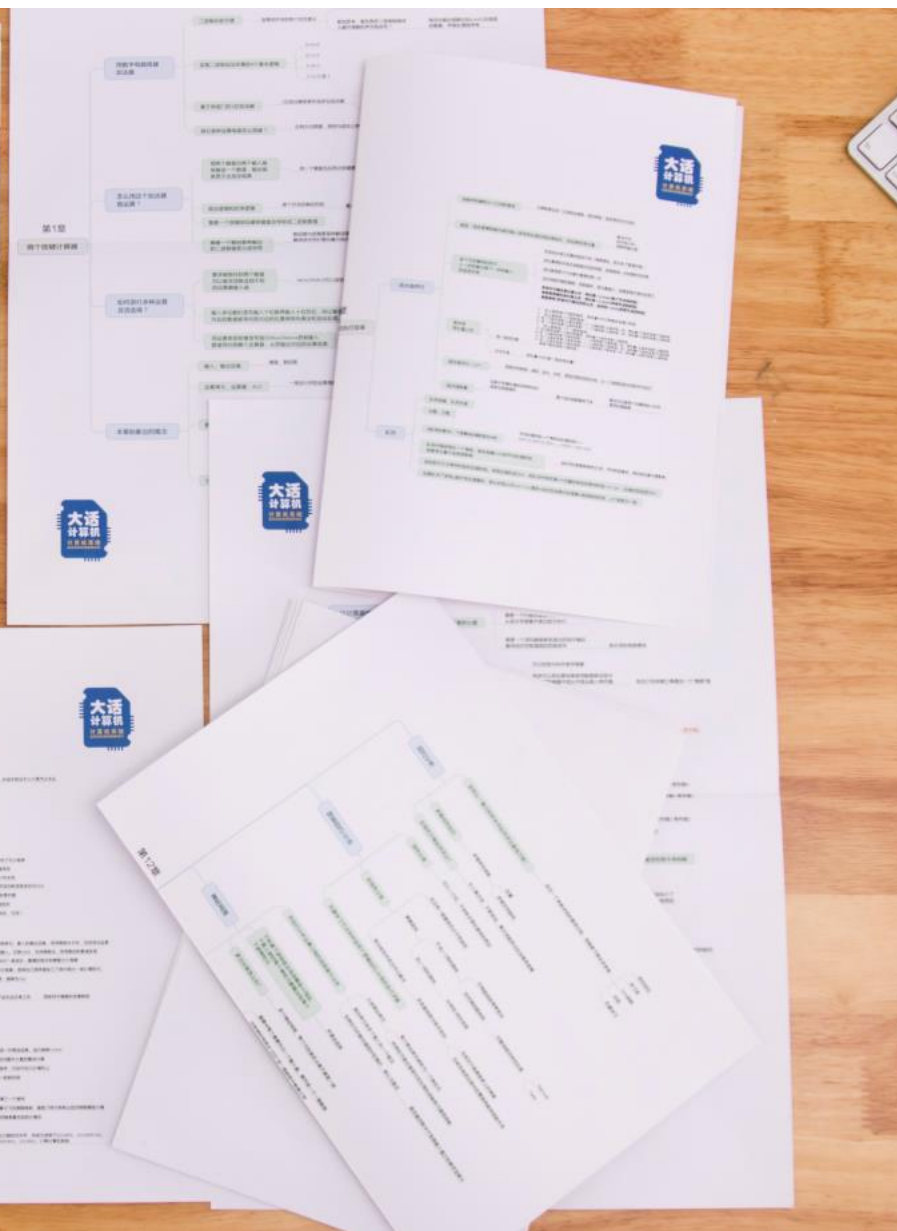












转交。Eager模式则是节点在收到消息之后就立即转发出去。同时节点会异步查询。直到网络上所有节点都收到该消息。一次之后回到该节点。这种方式浪费了很多带宽。如果网络带宽很宽。与广播无异。这两种方式都是没有过滤目录的。也就是说每个节点消息都必须所有人都知道。免不了。除非某个节点判断出最新数据就在自己。没必要往下转发了(比如Lazy模式)。第三种方法称为“Oracle”(先知)。之所以未知。当然是因为每个节点都实现了精确的过滤目录。可以点对点传递消息。

6.9.10.3 增设远程目录过滤片外广播

上文中介绍了单CPU芯片内部Inclusive的LLC缓存行或者Exclusive的所有缓存行中存放过滤bitmap目录。以实现片内广播过滤。该bitmap又被称为目录。由于其只能追踪本CPU片内的缓存状况。所以暂且称之为本地目录。同时也提到了。由于不知道其他CPU中是否缓存了该行。所以一些作成的消息不得不向其他所有核心广播。如果能够存在某种专门用于追踪其他CPU片内可能缓存有哪些缓存行的bitmap的话。这个bitmap就可暂且称之为远程目录。那样就可以避免广播而使用单播。节省网络流量。提升性能。下面我们使用一些流程图来帮助大家梳理清楚本地目录和远程目录在访问操作时的过滤作用全景。本系统采用Inclusive模式的缓存设计。直接在LLC缓存行中存储本地目录bitmap。

图6-167为一个在没有远程过滤目录情况下的CC案例示意图。C<sub>0</sub>核心发出一个Stor请求。C<sub>0</sub>的CA通过查询本地目录成功地发现C<sub>3</sub>核心的私有缓存中不会命中。那就不需要向其发送WrtIvid消息了。但是C<sub>0</sub>

的CA无法判断右侧的CPU内都有哪些缓存行。所以发送了个目标地址为广播地址的WrtIvid消息到该CPU控制器。右侧CPU的OP控制收到该广播。便会在该侧各独立的缓存行分别发给它的4个核心。后者各自查询本地目录(假设该访问的地址会被分配到C<sub>0</sub>核心的L3缓存分片)以判断自己的私有缓存里是否有该行。

如图6-168所示。C<sub>0</sub>和C<sub>3</sub>通过读取本地目录bitmap后发现自己对应的位为1。表明自己的私有缓存命中该行。于是去将对应行作废处理。然后返回IvidRspF消息。C<sub>0</sub>和C<sub>3</sub>通过bitmap发现自己核心里并没有缓存该地址。所以直接返回IvidRspN消息。同时。C<sub>0</sub>还从其他CPU收到了C<sub>3</sub>核心返回的IvidRspN消息。那么。右侧CPU收到了C<sub>0</sub>核心返回的IvidRspF消息。为对应的行已经被作废。不存在于该CPU内了。那么。该由哪个CA来更新目录呢?这是个问题。就像图中这样由这两个CA都去更新一遍。那会做无用功。降低性能。而且还会带来潜在的一致性问题。这个问题下文中再来解释。

可见。上述过程发送了太多广播。我们强烈期望下。本地CPU如何知道远程CPU到底缓存了哪些行呢?一开始。所有CPU的缓存都是空的。一旦任意一个CPU发生Load操作。那么必然会发出RdPrb消息。由于发送RdPrb消息的CPU的本地目录是空的。其必将该消息广播到全网。那么本地CPU就可以收到该广播。很自然地就知道了“哪个CPU读取了这一行并缓存了”。如果是Stor指令。一旦命中缓存。也会发起RdIvidPrb消息广播到全网。这样。其他CPU也就知道“这个CPU缓存了这一行”。只要经过是够长的时

没有远程目录时无法源过滤



图6-167 没有远程过滤目录时的CC示意图(3)







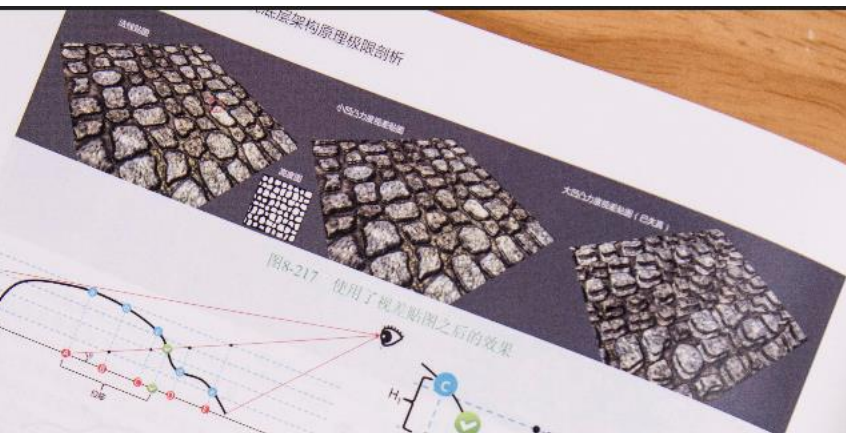


图8-217 使用了透视贴图之后的效果

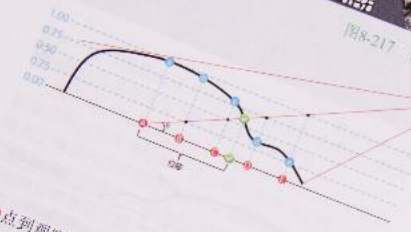


图8-218 迭代逼近法求交示意图

将A点到观察者的光线分别与每一层的等高线进行求交，由于等高线为直线，所以很容易求得A点与等高线交点的纹理坐标，然后从高度图中对应纹理坐标处读出这些纹理对应的假表面高度。显而易见，本例中的A和B，它们的表面高度值和层级高度值是相反的，也就是说，A点的表面高度值大于其对应的层级高度值，而B点的表面高度值小于其对应的层级高度值。交点一定位于被这样两个点夹住的表面区间内。所以，具体计算时并不需要将所有等高线对应的交点全求出，而是可以先求最高层级的点，考察其对应的表面高度值与层级高度值的关系，如果表面高度值小于层级高度值，则其继续与下一个层级高度求交并做相同比较，直到找到发生反转（表面高度大于层级高度值）的那个点，也就可以确定两个关键点了。或者采用二分法，先从中间层级高度开始求交计

算，如果发现表面高度大于层级高度，则从再高一层级的层级求交并计算、比较。然后，利用这两个点各自的表面高度与层级高度之间的差值，用两者的比例去插值A和B两点间的距离，最终得到I点的纹理坐标值。I点与理想值并不完全重合，只是更加接近了，如图8-218右侧所示。这种利用寻找两个相邻且表面高度比例反转的等高线交点，并在其之间水平距离上插值寻找光线与假表面交点的透视贴图技术，被称为**陡峭视差贴图 (Steep Parallax Mapping)**。如图8-219所示为SPM视差贴图效果，可以看到其凹凸得更加真实，失真度更小，而且其精度已经高到足以实现近处像素遮挡远处像素的真实效果了。看上去很难想象这个模型的表面原本是平坦的，不过，可以观察其平坦的边缘，表明其一定是障眼法，而没有使用置换贴图。

要提升陡峭视差贴图的近似度，就

层级切分得个优化是，是越发正对看以可以动态降低的视差遮蔽贴图 (Mapping)。其中Oc实现遮挡效果，其实时已经可以达到以个贴图已经可以达到在确定两个反转点之前的但是在这两个点之间寻找用了更精确的算法，篇幅家可以自行了解。

最后，如图8-220所示，直接把3D游戏的画面抬升了一幕无比。其中就使用了视差贴图。图技术出现的年代：凹凸贴图 (197 (1984年)、法线贴图 (1996年)、浮雕贴图 (2005年)。

### 8.2.7.4 物体投影 (Shadow)

我们前文中介绍的光照效果，都是致力型的表面凹凸细节明暗区域，虽然表面可能投不出影子，但是通过法线贴图处理依然可以

地安  
组成  
A分子

上的某个氨基酸残基上的电学基团产生电力，从而驱动蛋白质分子上的具有钾离子选择性的分子马达的ATP结合点（油门）暴露，ATP一拥而上将其旋转起来，从一侧吸引钾离子，转到另一面，释放钾离子。但是模拟量底层并不是无限连续的，其最小单元是一个正弦波形状。那么，一种可能的模型是：细胞一侧的钾离子浓度如果较高，那么该蛋白质分子在这一侧的某些残基上的能够与钾离子相互产生物理电作用或者化学力作用的基团每次受力都会对整个蛋白质分子深处的各个原子进行牵拉，当钾离子浓度较高的时候，这种被牵拉的频率就会越高，当牵拉频率达到某个值的时候，引起分子的构象产生一个比较大的跃变，将分子马达的ATP结合点暴露从而转运离子。当钾离子浓度不够的时候，阈值频率未达到，分子马达停止工作。随着钾离子被运输到另一侧，本侧浓度持续降低，最后形成一个闭环的负反馈控制系统。

蛋白质分子就是利用原子来搭建的精密机械，是个冲程，比如结合ATP分子时，构象改变，导致蛋白质转子旋转一定角度，ATP变为ADP后，转子再次旋转一定角度，ADP被释放后，转子再次旋转。这就像热机的吸气、压缩、爆炸、排气这四个冲程不断循环一样。而蛋白质分子马达纯粹依靠ATP分子形成的分子间作用力来拉动整个转子旋转。扫描二维码观看蛋白质分子马达的冲程动画。



一切全靠实验。讲，通电可以让转子转动，但是如果让转子转动，则可以反过来产生电流。对于蛋白质分子马达也是一样，ATP生成器也是一个分子马达，其在细胞膜内的区域是一个转子，转子表面包含多个质子结合点，可以利用高酸性环境驱动转子转动，转子转动导致胞外部分构象形变，从而将ADP+Pi合成为ATP，这相当于发电机。

对于图4左侧，诗意一些的表达则是：我要送你一棵树，象征着你我共同的目标。树根之下这是能量的源泉。质子驱动着马达把ADP和Pi强行结合成ATP。好吧，编

我们前文中介绍的光照效果，都是致力型的表面凹凸细节明暗区域，虽然表面可能投不出影子，但是通过法线贴图处理依然可以











## 大话 计算机

准备开脑洞了，  
请读者准备好耐心和勇气，  
与作者一同畅游  
缤纷计算机世界。

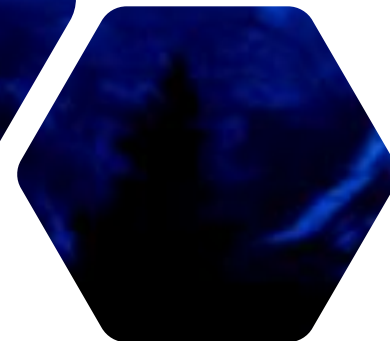
Enjoy!

清华大学出版社

清华大学出版社

## 大话 计算机

计算机系统  
原理架构原理及应用







还在看盗版？  
Out了！ 这书  
真的值得买  
纸张版！

绝对良心，这个作者绝对为了读者负责的角度去写作，就怕读者看着看着蒙了，解释的非常细致。书中内容真的是包罗万象，不仅仅是计算机。全篇透着作者对事物本质的追求和思索。

全彩印刷，作者用ppt做的图堪比autocad，世上绝无仅有的美图，将计算机运行原理动态展现出来，其他书中从没见过这种手法。盗版完全享受不到纸张版的质感。

这本书是一件**艺术品**，良心性价比。四五百块钱，拿出一两次玩网游吃大餐的钱就够了，可以说他价格高，但是千万别说贵。You can check out any time, but you will never leave.

三卷彩印精装书（正文共1506页）+11张海报+典藏礼盒。原价598元，折后价508.30元，券后**到手价458.30元**，这个品质这个价格，应该说很合适了。**预售入口**：长按下面二维码进入购买链接（请确保从下面入口进入购买京东自营，否则无法用代金券，切记！另外，一些京东第三方店铺会出现低价，不要相信！）。买买买买买买买买买买买买！

仅限预售期有效，预售期过后本价格失效





大话计算机



大话存储

作者的两个微信公众号，里面全是宝贝！