

# 对新时期编著高水平教材的若干思考

刘宇航 李浩铨

中国科学院计算技术研究所

**编者按：**教材是统筹科技、人才、教育的重要抓手，本文对计算机体系结构经典教科书中的两道例题进行勘误，提出应以学习、批判、建设的眼光阅读和编著高质量教科书。

## 引言

教科书在学生的学习中具有重要作用。教材建设意义大，难度高。如何撰写世界一流教材，是一个亟待解决但又容易被忽视的问题。在 20 世纪 50 年代，各条战线都急需人才，但毛主席指示“宁可把别的摊子缩小点，必须抽调大批干部编写教材”。2024 年 9 月 1 日，《求是》杂志上发表了一篇文章，再次强调了教材建设的重要性。

在计算机体系结构领域，有一本享誉世界的教科书《计算机体系结构：量化研究方法》(*Computer Architecture: A Quantitative Approach*)<sup>[1]</sup>，这本书有 30 多年的历史，由 2017 年 ACM 图灵奖得主约翰·亨尼西(John L. Hennessy)和大卫·帕特森(David A. Patterson)合著，目前已出版到第 6 版<sup>1</sup>。这本书被誉为计算机体系结构领域的“圣经”，是计算机体系结构方向学生的必读教材。

自 2006 年以来，笔者先后阅读过该教材第 3~6 版的英文版及对应的中文版。在阅读、研究该教材

的过程中，笔者时常思考：在新时期，我们如何编写出高质量的教材？是直接引进国外教材，还是出版国产教材？如果选择后者，是依靠少数名家还是集体创作？等等。

接下来，我们对这本享誉世界的计算机体系结构经典教科书最新版本中的两道例题进行勘误，然后进行关于高质量教材建设的若干延伸思考。

## 对其中一道例题及其解答的分析

我们首先对英文第 6 版第 96、97 页的一道例题及其解答进行分析。为便于读者结合中英文进行准确理解，这里给出该例题及解答的英文原文，如图 1 所示。

**例题：**利用（英文版教材）图 2.8 和附录 B 图 B.8 中的数据，判断 32 KiB 四路组相联 L1 高速缓存的存储访问时间是否快于 32 KiB 两路组相联 L1 高速缓存。假设 L2 高速缓存的缺失代价是较快的 L1 高速缓存的访问时间的 15 倍。忽略 L2 后续存

<sup>1</sup> 第 1、2、3、4、5、6 版分别于 1990 年、1996 年、2002 年、2007 年、2011 年、2017 年出版。

**Example** Using the data in Figure B.8 in Appendix B and Figure 2.8, determine whether a 32 KiB four-way set associative L1 cache has a faster memory access time than a 32 KiB two-way set associative L1 cache. Assume the miss penalty to L2 is 15 times the access time for the faster L1 cache. Ignore misses beyond L2. Which has the faster average memory access time?

**Answer** Let the access time for the two-way set associative cache be 1. Then, for the two-way cache,

$$\begin{aligned} \text{Average memory access time}_{2\text{-way}} &= \text{Hit time} + \text{Miss rate} \times \text{Miss penalty} \\ &= 1 + 0.038 \times 15 = 1.38 \end{aligned}$$

For the four-way cache, the access time is 1.4 times longer. The elapsed time of the miss penalty is  $15/1.4 = 10.1$ . Assume 10 for simplicity:

$$\begin{aligned} \text{Average memory access time}_{4\text{-way}} &= \text{Hit time}_{2\text{-way}} \times 1.4 + \text{Miss rate} \times \text{Miss penalty} \\ &= 1.4 + 0.037 \times 10 = 1.77 \end{aligned}$$

Clearly, the higher associativity looks like a bad trade-off; however, because cache access in modern processors is often pipelined, the exact impact on the clock cycle time is difficult to assess.

(a) 例题及解答的英文原文

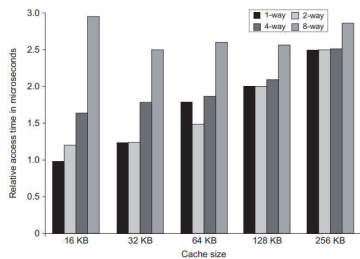


Figure 2.8 Relative access times generally increase as cache size and associativity are increased. These data come from the CACTI model 6.5 by Tarjan et al. (2005).

(b) 例题对应的图

Cache size (KiB)	Degree associative	Total miss rate	Miss rate components (relative percent) (sum = 100% of total miss rate)					
			Compulsory	Capacity	Conflict			
32	1-way	0.042	0.0001	0.2%	0.037	89%	0.005	11%
32	2-way	0.038	0.0001	0.2%	0.037	99%	0.000	0%
32	4-way	0.037	0.0001	0.2%	0.037	100%	0.000	0%
32	8-way	0.037	0.0001	0.2%	0.037	100%	0.000	0%

(c) 例题对应的表

图1 第一道例题、解答和对应图表的英文原文

储层次的缺失。哪一种高速缓存的平均存储访问时间较短？

### 原解答

设两路组相联高速缓存的访问时间为1。则对于两路高速缓存：

$$\begin{aligned} \text{平均存储访问时间}_{\text{两路}} &= \text{命中时间} + \text{缺失率} \times \text{缺失代价} \\ &= 1 + 0.038 \times 15 = 1.57 \end{aligned}$$

对于四路高速缓存，访问时间是两路高速缓存的1.4倍。缺失代价占用的时间为  $15/1.4 \approx 10.7$ 。为简单起见，设其为10：

$$\begin{aligned} \text{平均存储访问时间}_{\text{四路}} &= \text{命中时间}_{\text{两路}} \times 1.4 + \text{缺失率} \times \text{缺失代价} \\ &= 1.4 + 0.037 \times 10 = 1.77 \end{aligned}$$

显然，采用较高的相联度看起来是一种糟糕的权衡选择；不过，由于现代处理器中的高速缓存访

问通常都实现了流水化，很难评估对时钟周期时间的具体影响。

### 分析

这道例题立意很好，旨在考查对命中延迟与命中率之间的权衡的理解，但是教材提供的答案是不完全正确的，原因在于单位混用。正确的做法是，要么一直使用时钟周期（clock cycle）作为单位，要么一直使用秒（second）作为单位，而不能混合使用两者。

秒是一个绝对的时间单位，时钟周期是一个相对的时间单位。四路组相联的L1高速缓存的命中时间是两路组相联的L1高速缓存的1.4倍。时钟周期长度一般等于L1高速缓存的命中时间。

显然，使用秒作为单位更为方便。虽然使用时

钟周期作为单位也是可行的，但不同的微架构可能有不同的时钟周期，这会使性能比较变得更加繁琐。但无论如何，不能混用这两种时间单位。

具体而言，单位混用的错误出现在计算四路高速缓存的存储器平均访问时间中。原解答中，计算四路高速缓存的缺失代价使用的式子是  $15/1.4 \approx 10.7$ ，这样的计算方式说明此处使用的时间单位是“时钟周期”。根据题中假设，两个不同关联度的 L1 高速缓存所对应的 L2 高速缓存是相同的，所以两个不同关联度的 L1 高速缓存的缺失代价若用“秒”作为时间单位也应该是相同的。对“15 倍”除以 1.4 进行转化，意在把两路 L1 高速缓存的时钟周期转化为四路 L1 高速缓存的时钟周期。这样转化的结果是，原解答中的四路 L1 高速缓存命中时间以两路 L1 高速缓存对应的时钟周期为单位。这样在同一个公式中，使用两种不同的时钟周期作为时间单位，导致错误。

下面我们给出正确解法，分别使用绝对时间单位和相对时间单位（即不同情形下的时钟周期）作为计时单位。

## 正确解法1：使用绝对时间单位

第一种解法使用绝对时间单位（如秒），本解法使用单位时间“unit”。

设两路组相联 L1 高速缓存的访问时间为 1 unit，故对于该高速缓存：

$$\begin{aligned} \text{平均存储访问时间}_{\text{两路}} &= \text{命中时间}_{\text{两路}} + \text{缺失率}_{\text{两路}} \times \text{缺失代价}_{\text{两路}} \\ &= 1 + 0.038 \times 15 = 1.57 \text{ units} \end{aligned}$$

四路组相联 L1 高速缓存的访问时间是两路组相联 L1 高速缓存的 1.4 倍，即  $1 \times 1.4 = 1.4 \text{ units}$ 。对于同一个 L2 高速缓存，两个 L1 高速缓存的缺失代价相同，皆为 15 倍的两路组相联 L1 高速缓存的访问时间，即  $1 \times 15 = 15 \text{ units}$ 。故：

$$\begin{aligned} \text{平均存储访问时间}_{\text{四路}} &= \text{命中时间}_{\text{四路}} + \text{缺失率}_{\text{四路}} \times \text{缺失代价}_{\text{四路}} \\ &= \text{命中时间}_{\text{两路}} \times 1.4 + \text{缺失率}_{\text{四路}} \times \text{缺失代价}_{\text{四路}} = \end{aligned}$$

$$1 \times 1.4 + 0.037 \times 15 = 1.955 \text{ units}$$

## 正确解法2：使用相对时间单位

第二种解法使用时钟周期，该值对同一 CPU 而言是固定的<sup>2</sup>，但是对不同 CPU 可能不相同。假设两个 L1 高速缓存对应 CPU 的时钟周期分别等于它们的命中时间，故对于两路组相联 L1 高速缓存：

$$\begin{aligned} \text{平均存储访问时间}_{\text{两路}} &= \text{命中时间}_{\text{两路}} + \text{缺失率}_{\text{两路}} \times \text{缺失代价}_{\text{两路}} \\ &= 1 + 0.038 \times 15 = 1.57 \text{ 时钟周期}_{\text{两路}} \end{aligned}$$

无论是四路组相联 L1 高速缓存，还是两路组相联 L1 高速缓存，命中时间都是其对应的 CPU 的一个时钟周期，缺失代价取决于后备缓存的结构，因此也是相同的。题干假定访问 L2 高速缓存的缺失代价是较快的 L1 高速缓存的访问时间的 15 倍，这里，较快的 L1 高速缓存就是两路组相联 L1 高速缓存。本解法使用时钟周期作为单位，所以需要将此对两路组相联 L1 高速缓存而言的时钟周期转换为对四路组相联 L1 高速缓存而言的时钟周期。由于四路 L1 高速缓存的访问时间是两路组相联 L1 高速缓存的 1.4 倍，所以四路 L1 高速缓存的时钟周期是两路 L1 高速缓存的 1.4 倍，故：

$$\begin{aligned} \text{缺失代价}_{\text{四路}} &= \text{缺失代价}_{\text{两路}} = 1 \text{ 时钟周期}_{\text{两路}} \times 15 \\ &= 1 \text{ 时钟周期}_{\text{四路}} / 1.4 \times 15 = 10.7 \text{ 时钟周期}_{\text{四路}} \end{aligned}$$

$$\begin{aligned} \text{所以，平均存储访问时间}_{\text{四路}} &= \text{命中时间}_{\text{四路}} + \text{缺失率}_{\text{四路}} \\ &\times \text{缺失代价}_{\text{四路}} = 1 + 0.037 \times 10.7 = 1.3959 \text{ 时钟周期}_{\text{四路}} \end{aligned}$$

可以看到，以时钟周期为单位时，四路组相联 L1 高速缓存的平均存储访问时间小于两路组相联 L1 高速缓存，但是这并不意味着四路组相联 L1 高速缓存的访存速度更快，因为两者的时钟周期不相同。为了比较哪一种结构的访存速度更快，应该使用统一的时间单位，可以将两者的时钟周期转化为统一的单位再进行比较。最终可以发现，在题干的假设下，两路组相联 L1 高速缓存的访存速度快于四路组相联 L1 高速缓存。

<sup>2</sup> 当 CPU 不使用动态电压频率调节（Dynamic Voltage and Frequency Scaling, DVFS）时，CPU 的主频一般是固定的，时钟周期也因此是固定的。

## 对其中另一道例题及其解答的分析

我们再对英文第6版第103页的一道例题及其解答进行分析。为便于读者结合中英文进行准确理解,这里给出该例题及解答的英文原文,如图2所示。

**例题:** 假定主存储器的访问时间为36 ns, 存储器系统的持续传输速率为16 GiB/s。设块大小为64字节。如果在给定请求流的情况下能够保持峰值带宽, 而且访问从来不会冲突, 则需要支持的最大未决<sup>3</sup>缺失数目为多少? 如果一次访问与其之前的4次访问之一发生冲突的概率为50%, 并且每次冲突的访问都要等待更早的访问完成才能继续进行, 请估计最大未决访问数目。为简单起见, 忽略缺失之间的时间。

### 原解答

在第一种情况下, 假定我们可以保持峰值带宽, 存储器系统支持每秒 $(16 \times 10^9)/64=2.5$ 亿次访问。由于每次访问耗时36 ns, 因此可以支持 $2.5 \times 10^8 \times 36 \times 10^{-9}=9$ 次访问。如果发生冲突的概率大于0, 我们就需要更多的未决访问, 因为存储系统无法处理发生冲突的访问; 存储器系统需要更

多的独立的访问。为了简单估计这一数目, 假定有一半存储访问不能立即发送到存储器。这就意味着必须支持两倍的未决访问, 即18次。

### 分析

尽管上述解决方案的核心思想是正确的, 但存在一些瑕疵甚至错误, 在定量分析方面也有待提高。

瑕疵在于传输速率单位换算。题目中给出“存储器系统的持续传输速率为16 GiB/s”, 使用的单位是二进制速率单位, 而非十进制的GB/s, 但是原解答直接将其视为十进制单位来计算存储器系统每秒支持的访问数, 带来了误差。因为:

$$1 \text{ GiB/s} = (2^{10})^3 \text{ B/s} = 1073741824 \text{ B/s}$$

该值与 $10^9$ 相差了约7.4%, 不能随意忽略, 所以正确的解法应该使用正确的GiB/s换算方式。为了解答的方便, 我们建议将题干中的GiB/s改为GB/s。

除了瑕疵之外, 更重要的是, 原解法的定量分析过于粗略, 没有恰当地使用题目条件, 没有充分贯彻该书原本所倡导的量化研究的理念。题目条件指出“一次访问与其之前的4次访问之一发生冲突的概率为50%”, 原解答将其粗略理解为“一半存储访问不能立即发送到存储器”, 偏差较大。

**Example** Assume a main memory access time of 36 ns and a memory system capable of a sustained transfer rate of 16 GiB/s. If the block size is 64 bytes, what is the maximum number of outstanding misses we need to support assuming that we can maintain the peak bandwidth given the request stream and that accesses never conflict. If the probability of a reference colliding with one of the previous four is 50%, and we assume that the access has to wait until the earlier access completes, estimate the number of maximum outstanding references. For simplicity, ignore the time between misses.

**Answer** In the first case, assuming that we can maintain the peak bandwidth, the memory system can support  $(16 \times 10^9)/64 = 250$  million references per second. Because each reference takes 36 ns, we can support  $250 \times 10^6 \times 36 \times 10^{-9} = 9$  references. If the probability of a collision is greater than 0, then we need more outstanding references, because we cannot start work on those colliding references; the memory system needs more independent references, not fewer! To approximate, we can simply assume that half the memory references do not have to be issued to the memory. This means that we must support twice as many outstanding references, or 18.

### 正确解答

在第一种情况下, 假设可以保持带宽峰值, 存储系统每秒能够支持 $(16 \times 10^9)/64=2.5$ 亿次访问。若所有在1秒内以峰值带宽传输的数据对应的存储访问都串行执行, 则访问时间为 $2.5 \times 10^8 \times 36 \times 10^{-9}=9$ 秒。若要在1秒内发出2.68亿次访问, 存储系统的

图2 第二道例题的英文原文

<sup>3</sup> “outstanding”译为“未决的”, “outstanding misses”指的是那些已经被发起但还没有被解决或者完成的高速缓存缺失请求。



## 关于教材建设的思考

编写高水平教材，要脚踏实地、心无旁骛。在多样化媒体竞相抢占有限注意力的时代，沉下心来读书的人少了。教材是对一门学科的系统梳理，涉及较长的历史和较多的方向，往往篇幅较大，例如《计算机体系结构：量化研究方法（第6版）》英文版有1527页，《人工智能：一种现代方法（第4版）》英文版有2145页。一般而言，书越厚，越难写，也越难读。高质量的教材，需要作者沉下心来去写，也需要读者沉下心来去读。在很多领域，国内教材数量并不少，但往往为了出版而出版，良莠不齐。也有很多学生不读教材，导致知识体系碎片化。

编著高水平教材，要坚持长期主义、凝聚群智。教材建设非一日之功，非一人之功。编著教材考验的不仅是编著者的写作能力、语言能力，还包括编著者的专业造诣、思维能力。《计算机体系结构：量化研究方法》已有34年的历史，历经6个版本，据说第7版也在准备之中。各个版本都有分布在世界各地的数百位学者参与修订。从这个角度看，可能不宜运动式、集中式地编写教材。教材应该在丰富的、深刻的、真实的科研和教学实践中自然产生，要经得起时间和实践的检验，历久弥新。

编著高水平教材，要注重借鉴国际一流教材，也要看到其不足。通过本文所述的两道例题，大家可以看到计算机体系结构领域的“圣经”也有瑕疵甚至错误。书中还有更多的瑕疵和错误，限于篇幅，不再赘述。对国外教材，要虚心学习，但不宜盲目迷信。教师和学生都要认真钻研，勤于思考，思辨概念，实操习题。

编著高水平教材，要参考人类历史上已完成的知识工程的成功经验。我国历史上的《吕氏春秋》《资治通鉴》等书都有为某一类读者作教科书的动机，其编著团队人员造诣深、数量多、用时长，对质量的要求达到近乎苛刻的程度。成语“一

字千金”<sup>6</sup>的寓意值得深思。诸如此类的历史经验值得我们参考。

编著高水平教材，要融合科研和教育的成果和力量。教材要成为联系科研和教育的纽带。高水平教材应是长期科研和教学的思考的结晶。在现实中，科研工作者和教育工作者往往各自为战，科研脱离教育，教育脱离科研。科研工作者要以“能否写入教科书，能否经得起历史检验”这个角度审视研究成果，教育工作者要以“教学材料能否反映研究前沿，能否促进学生研究能力的培养”这个角度审视教学成果。笔者所在课题组在这方面进行了初步尝试<sup>[3]</sup>，相关的实践值得继续深入推进和研究，并在此基础上不断改进，臻于至善。

## 结束语

教材对教育、科学研究具有重要作用，关乎一门学科的人才培养质量高不高，关乎一门学科的发展后劲足不足。为了将教育、人才、科技融为一体并相互促进，科研工作者、教师和学生需要主动认识和遵循教材建设的规律，以学习、批判、建设的眼光阅读和撰写高质量教科书。 ■



刘宇航

CCF 高级会员，CCCF 专栏编委。中国科学院计算技术研究所副研究员。主要研究方向为计算机体系结构、高性能计算、存储系统、智能并发系统。liuyuhang@ict.ac.cn



李浩铨

中国科学院计算技术研究所硕士研究生。主要研究方向为计算机系统结构。lihaoquan24s@ict.ac.cn

（本文责任编辑：翟季冬）

延伸阅读请登录 <http://dl.ccf.org.cn/cccf/list>

<sup>6</sup> 吕不韦乃使其客人著所闻，集论以为八览、六论、十二纪，二十余万言，以为备天地万物古今之事，号曰《吕氏春秋》。布咸阳市门，悬千金其上，延诸侯游士宾客，有能增损一字者，予千金。——《史记·吕不韦列传》