

基于阿姆达尔定律和兰特法则 计算多核架构的加速比

李文石, 姚宗宝

(苏州大学电子信息学院微电子学系, 江苏苏州 215006)

摘 要: 在评价多核 CPU 加速比已知模型的基础上, 基于第一性计算原理融合理解阿姆达尔定律和兰特法则, 提出描述多核 CPU 加速比的一个新模型. 研究方法是从传统的阿姆达尔定律切入, 论述的逻辑顺序分别基于约束固定任务, 固定时间, 存储器和互连复杂性; 兼顾了举例论述同构多核的 NoC 带宽性质和最大温度特性. 计算表明: 基于固定时间模型与存储器模型预测多核的加速能力, 容易得到估计结果的乐观上限; 我们提出的基于兰特法则的模型计算结果, 在并行比例较大时稍小于但接近前述模型估计值, 而比固定任务模型的保守结果要好; NoC 带宽和最大温度的结果提示, 多(同构)核 CPU 期盼相对高的并行度架构.

关键词: 多核处理器; 阿姆达尔定律; 加速比; 兰特法则; 第一性原理; 带宽; 温度

中图分类号: TP301.6 **文献标识码:** A **文章编号:** 0372-2112 (2012) 02-0230-05

电子学报 URL: <http://www.ejournal.org.cn>

DOI: 10.3969/j.issn.0372-2112.2012.02.004

Multicore Architecture Speedup Computation Based on Amdahl's Law and Rent's Rule

LI Wen-shi, YAO Zong-bao

(School of Electronics and Information Engineering, Soochow University, Suzhou, Jiangsu 215006, China)

Abstract: We made Amdahl's Law and Rent's Rule understood and proposed one novel model to describe the multicore CPU speedup based on evaluating the known speedup models. Our methodology starts with traditional Amdahl's Law, in the steps of fixed-size, fixed-time, memory-bounded, and interconnection-bounded. Also the NoC-bandwidth and max-temperature of multicore features are discussed. The results showed, fixed-time model and memory-bounded model grasp speedups' up-limits, our model touches the middle gap between formers and Amdahl's Law, at high parallel ratio, and NoC bandwidth and Max temperature show that homogeneous multicore architectures long for higher ratio of parallelization.

Key words: multicore CPU; Amdahl's law; speedup; Rent's rule; first principle; bandwidth; temperature

1 引言

针对散热、速度、存储和带宽等性能的新挑战, 半导体学术界和工业界在 20 世纪中期开始探索多核 CPU 解决方案的原理可行性与技术可行性.

1996 年, Stanford 大学首先提出片上多处理器和首个多核结构原型^[1]; 2001 年, IBM 推出第一个商用多核处理器 POWER4^[1]; 2005 年, Intel 和 AMD 开始挺进多核处理器市场, 先后推出 2、4、8 与 16 核 CPU. 多核处理器架构成为了拉动摩尔定律发展新机制^[2].

2007 年, MIT 教授兼 Tiler 公司创始人 A. Agarwal 与 EEMBC 创始人兼总裁 Markus Levy 发表专论预测: 十年内将出现千核处理器^[3]; 2010 年 12 月 29 日国外媒体报

道, 英美科学家基于 FPGA 成功开发出一款千核处理器, 可在提高计算速度的同时降低能耗^[4].

评价多核架构性能的参数是加速比(speedup), 定义为多核串行所需时间除以并行所需时间. 例如, 文献[5]研究认为: 串行 PCA 用时 95.08 单位, 而并行只需 5.86 单位, 则其加速比就是 16.22. 文献[6]重视研究多核加速比极限问题, 为本工作提供了研究范式.

首先评述阿姆达尔定律; 接着基于三维图解, 梳理已知的加速比极限研究结果(三个已知模型), 重点结合兰特法则(Rent's Rule)描述的互连约束, 基于第一性原理重新计算了加速比极限(构造一个新模型), 比较讨论了 4 个模型的结果, 例析了芯片温度与核数的关系.

2 多核处理器的加速比研究演进

根据模型提出的时间顺序(1967, 1988, 1990), 分别讨论多核处理器的三种加速比, 构造三维图解, 以便突出其间的可比性^[6].

2.1 固定任务(fixed-size)模型

1967年, IBM大型机之父阿姆达尔(Gene M. Amdahl)博士图解了并行计算系统设计的关键.

阿姆达尔定律指出:系统某一部件由于采用某种更快的执行方式后, 整个系统功效的提高与这种执行方式使用频率占总执行时间的比例有关. 由并行方法所能获得的加速比为:

$$Speedup_{Amdahl} = \frac{1}{1 - f + \frac{f}{m}} \quad (1)$$

其中, f 为问题中可被并行处理部分的比例, $(1 - f)$ 是串行的比例, m 为并行处理器的数量, $Speedup$ 为并行时相比于串行时的加速比.

如下择要推导阿姆达尔定律, 旨在还原与加深理解大型机的并行架构发明理念.

设若多核芯片最多内置 n 个基本核(BCE, Base core equivalent), 运用多个 BCE 资源可组成一个更高性能的内核: 令单个 BCE 的性能(理解为运算速度)为 1, 设用 r 个 BCE 内核所创建结构的串行性能为 $perf(r)$ ^[7].

根据加速比的原始定义, 有

$$Speedup = \frac{\text{加速后的性能}}{\text{原性能}} = \frac{T_{ori}}{T_{enh}} \quad (2)$$

设问题的工作任务为 w , 那么单个 BCE 执行时间为 $T_{ori} = w / perf(1) = w$, 而 n 个 BCE 基本核执行时间则为

$$T_{enh} = \text{串时} + \text{并时} = \frac{(1 - f)w}{perf(r)} + \frac{fw}{\frac{n}{r} \cdot perf(r)} \quad (3)$$

其中, $n/r = m$ 是核数, 且每个核的串行性能为 $perf(r)$, 将 T_{ori} 和 T_{enh} 代入式(2), 则加速比实为

$$\begin{aligned} Speedup &= \frac{w}{\frac{(1 - f)w}{perf(r)} + \frac{fw}{m \cdot perf(r)}} \\ &= \frac{1}{\frac{1 - f}{perf(r)} + \frac{f}{m \cdot perf(r)}} \end{aligned} \quad (4)$$

对于给定的多核设计, r 个 BCE 的 $perf(r)$ 是常数, 设为 c , 则式(4)化简为

$$Speedup = \frac{c}{1 - f + \frac{f}{m}} \quad (5)$$

阿姆达尔定律的解析式(1)就是 $c = 1$ 的式(5). 讨论: 阿姆达尔定律默认工作任务(workload)是固定的, 这

种加速模式强调解决给定任务的耗时减少了, 所以阿姆达尔定律也叫固定任务(fixed-size)模型.

三维图解多核架构的固定任务模型, 参见图 1.

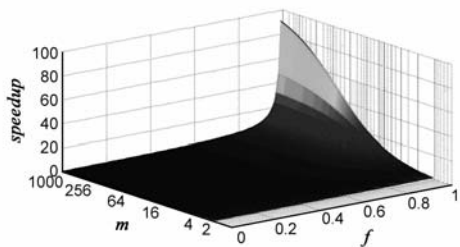


图1 多核架构的固定任务模型

述评: 阿姆达尔定律表明, 在问题的可并行部分不大时, 增加处理机的数量并不能显著地加快解题速度. 这曾让计算机界产生过悲观情绪, 有专家认为搞多处理器的机器没什么前途^[6].

分析可知提速的机会还是存在的. 虽然, 阿姆达尔定律基于约定: (1) 固定工作任务, (2) 固定的并行化比例; 但是, 其忽略的重要事实是: 在实际应用中, 当工作任务更大时, 通常情况下, 该任务也有更大的可能被分为可并行化的小任务(或者说处理多个相互独立的任务), 此也就意味着 f 更大(更接近 1), 可能得到更大的加速比.

2.2 固定时间(fixed-time)模型

直到 1988 年, J. L. Gustafson 提出了一个固定时间模型(fixed-time model), 也即 Gustafson 定律, 专家们重拾对大规模并行计算的信心^[8].

同样假设原始工作任务为 w , 比例扩增的工作任务(scaled workload)为 w' , 分别是串行条件下和并行条件下 m 个核在同样时间里完成的工作任务, 故有 $w' = (1 - f)w + fmw$. 基于类似速度关系的比例计算, 因此

$$Speedup = \frac{\text{串行解决 } w' \text{ 的时间}}{\text{并行解决 } w' \text{ 的时间}} = 1 - f + fm \quad (6)$$

此公式称为 Gustafson 定律.

图 2 展示了多核架构固定时间加速的性能改善的总趋势.

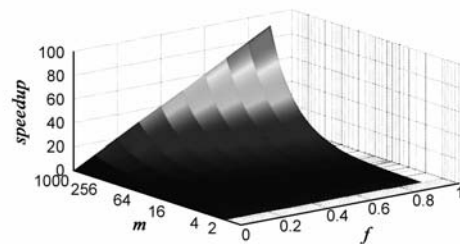


图2 多核架构的固定时间模型

Gustafson 定律表明: (1) 如果工作任务被扩大以保持固定执行时间, 则固定时间内的加速比 $speedup$ 是有关 m 的线性关系式, (2) 由于 $speedup$ 可以随系统规模(m)线性增长(更接近 1 的 f), 建立一个大规模的并行

系统将有裨益。

论述至此,研究多核加速比的提速问题,尚未考虑基本核中存储器的约束或曰贡献。

2.3 存储器约束(memory-bounded)模型

在 1990 年, X. H. Sun 和 L. Ni 提出了存储器约束模型(memory-bounded model)^[1,2], 即 Sun and Ni 定律。

设 w_e 为存储器空间限制下的规模扩大的工作任务, 则加速比形如^[7]

$$\text{Speedup} = \frac{\text{串行解决 } w_e \text{ 的时间}}{\text{并行解决 } w_e \text{ 的时间}} \quad (7)$$

设 $y = g(x)$ 是反映存储器容量(随核数)增长 m 倍时并行工作任务增长因数的方程式。

再设存储器节点容量是 M , 则有原始工作任务 $w = g(M)$, 而规模扩大的工作任务 $w_e = g(m * M)$, 我们有 $w_e = g(m * g^{-1}(w))$ 。

基于类似于处理式(3)的方法, 解算式(7), 有

$$\text{Speedup} = \frac{(1-f)w + f \cdot g(m \cdot g^{-1}(w))}{(1-f)w + \frac{f \cdot g(m \cdot g^{-1}(w))}{m}} \quad (7')$$

化简繁式(7'): 对于任何一个激励方程 $g(x) = a * x^b$ 和任意有理数 a 和 b , 我们有

$g(mx) = a(mx)^b = g(x)m^b = g'(m)g(x)$, 其中显然 $g'(m) = m^b$

$$\text{Speedup} = \frac{(1-f)w + f \cdot g'(m)}{(1-f)w + \frac{f \cdot g'(m)}{m}} \quad (7'')$$

图 3 对存储器约束模式的加速比进行了三维图解。

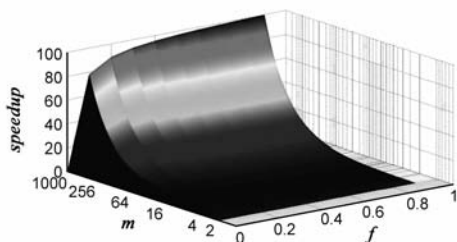


图3 多核架构的存储器约束模型

基于图 1 做整体概貌理解, 比较图 2 和图 3 显知: 并行比例 f 较小时, 加速比提升可观(数据参见表 1)。原因是仔细分析约束条件(①任务; ②时间; ③存储器), 逐渐精细建模得到加速比的阶段优化结果。

然而, 当我们重新审视约束条件(④互连复杂性; ⑤带宽), 可能对多核 CPU 的加速比优化问题, 产生全新的认识和理解。

3 基于兰特法则建立与修正加速比模型

1960 年 IMB 公司的工程师兰特(E. F. Rent)发现了基于同质模块构造计算机系统的互连复杂性规律, 这就是兰特法则(Rent's Rule), 描述为式(8),

$$T = kG^\beta \quad (8)$$

其中, T 为终端(引脚)数, G 为芯片上的模块(同构核)数, k 为平均每个模块上的终端数, β 是与芯片上并行比例有关的参数^[8]。

首先借鉴第一性原理推导兰特法则, 基于此构造加速比的新模型; 接着, 结合带宽的兰特法则, 建立加速比与带宽的关系模型; 最后, 利用表格对比展示各个模型的加速比的可比较性质, 也例析了芯片温度与同质核数的关系。

3.1 互连约束模型(应用 Rent's Rule 描述加速比)

借鉴文献[9]的研究思路, 设一个处理器有 m 个核时加速比为 S 。若核数有 Δm 的微小变化, 在没有其他信息改变的情况下, 我们只能认为加速比的相对变化与核数的相对变化, 存在比例关系 β , 其是与处理器并行比例有关的参数。

近似写作

$$(dS)/S = \beta(dm)/m \quad (9)$$

积分得:

$$S = km^\beta \quad (10)$$

k 为与平均每个核的加速相关的一个常数, 类似式(3)中的 r , 因此 $k = r = 1$ 。

由于 β 与处理器并行比例有关, 可以用 αf (α 为常数)代替。

最终

$$\text{Speedup} = m^{\alpha f} \quad (11)$$

令式(11)和式(6)在 $f = 0.99$ 及 $m = 1024$ 处相等, 可计算出 $\alpha \approx 1.0086$, 因此令 $\alpha = 1$ 对互连约束模型做三维图解示于图 4。

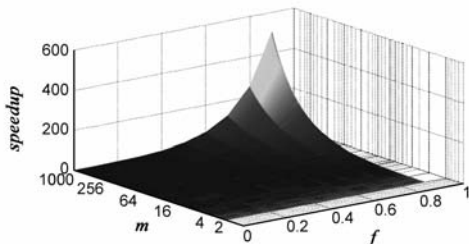


图4 基于兰特法则描述加速比

分析可知: 相对于图 1 与图 2, 我们的这种加速比估算值, 介于固定工作任务模型和固定时间模型之间。

3.2 带宽约束模型(建立加速比与带宽的关系)

由文献[10]知道: 研究 NoC(片上网络)技术得到带宽的新表述

$$B = kG^\beta \quad (12)$$

在多核处理器中, $k = r$, 可以设为 1, G 为芯片上的模块数, 相当于多核中的 m , β 可以用 αf 代替。则式(12)可以转化为

$$B = m^{\alpha f} \quad (13)$$

将式(13)代入阿姆达尔定律中,仍设 $\alpha = 1$,可得

$$Speedup = \frac{1}{1 - f + \frac{f}{B^{1/f}}} \tag{14}$$

利用 MATLAB 绘制三维图,得到图 5。

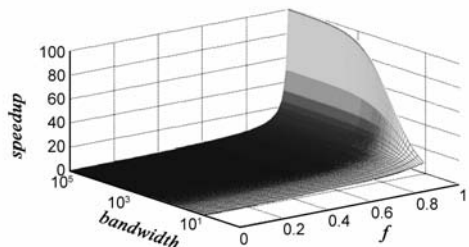


图5 加速比与带宽的关系

分析可知:该加速比的提升,强烈伴随着并行比例趋于 1,同时也给出宽带特性。

3.3 4 种加速比模型的结果比较

针对前面分析与提出的四种加速比模型,在核数 $m = 1024$ 条件下,分别计算加速比值,列入表 1。

表 1 4 种加速比模式的结果比较

模式 f	固定任务	固定时间	存储器约束	兰特法则
0.2	1.25	205.60	910.33	4.00
0.4	1.67	410.20	963.80	16.00
0.6	2.50	614.80	1003.12	64.00
0.8	4.98	819.40	1016.07	256.00
0.9	9.91	921.70	1020.46	512.00
0.99	91.18	1013.77	1023.68	955.42

小结:基于固定时间约束模型与存储器约束模型预测多核的加速能力,容易得到估计结果的乐观上限;而我们提出的基于兰特法则的模型计算结果,在并行比例较大时稍小于但接近上述两种模型估计值,确又比固定任务模型的保守结果要好。

3.4 核数与温度的关系

特别地,我们也关心多核优化的功耗限制问题^[11]。温度单位为摄氏度(℃), TDP 表征处理器可产生的最大功耗;由文献[12]芯片峰值温度与内部核数(m)和总散热设计功耗(TDP)的关系式描述如下

$$T_{\max} = TDP \cdot \left(R_{\text{conv}} + \frac{t_{\text{si}} - t_{\text{iso}(m)}}{kA} + \frac{t_{\text{iso}(m)}}{k} \frac{1}{A(1 - Ca(m))} \left| \frac{1}{1 + jw_s \tau_s} \right| \right)$$

其中的参数要点包括: R_{conv} 是热阻, A 是芯片面积。计算得到结果列入表 2。

表 2 芯片峰值温度值与核数 m 和散热设计功耗 TDP 的关系

TDP m	200W	250W
1	44.0℃	55.0℃
10	40.5	51.0
100	34.5	43.0
1000	32.5℃	41.0℃

小结:表 2 提示,处理器内的(同构)核数越多,峰值温度越低。

4 结论

阿姆达尔定律(1967 年)首次概括了加速比跟同质核的数目与并行度的深刻关系,基于固定任务约束。

Gustafson 定律(1988 年)针对工作任务很大时,设法提高并行比例,得到较大的加速比,基于固定时间约束。

Sun and Ni 定律(1990 年)的核心建模思想是加入了存储器约束,使加速比预测更接近实际峰值。

我们融合阿姆达尔定律和兰特法则思想,提出了一种新的表征多核处理器加速比的方法,经过验证,并行度 f 分别为(0.4,0.8,0.99)时,阿姆达尔定律计算出的加速比分别(1.67,4.98,91.18),而新表征方法计算值分别为(16,256,955.42),与固定时间加速模型比较,新方法在大 f 下比较接近估算加速比。

关于同构多核的 NoC 带宽性质和最大温度特性,本工作也给予了数据比较结果,结论是同质多核 CPU 的内部技术驱动,期盼相对高的并行度架构。

多核技术的未来趋向异构多核^[13~15],因此本工作权作多核建模优化计算的入门研究(针对同构多核)。

参考文献

[1] 黄国睿,张平,魏广博.多核处理器的关键技术及其发展趋势[J].计算机工程与设计,2009,30(10):2414-2418.

[2] 回首 05 多核之路:AMD 英特尔 Sun 的技术攻坚战[EB/OL] <http://news.pconline.com.cn/hy/0512/742256.html>, 2005-12.

[3] A Agarwal,M Levy. The kill rule for multicore[A]. Proc of IEEE Design Automation Conference[C]. San Diego, 2007. 750-753.

[4] 科学家开发千核处理器运算速度提升 20 倍 [EB/OL] <http://www.it.com.cn/news/cyxw/gjjy/2010123016/952063.html>, 2010-12-30.

[5] G Seshadri,R Jain,A Mittal. Parallelization of principal component analysis[A]. IEEE Advance Computing Conference[C]. Patiala, 2010. 44-49.

[6] Xian-he Sun,Yong Chen. Reevaluating Amdahl's law in the multicore era[J]. Journal of Parallel and Distributed Computing, 2010, 70(2): 183-188.

[7] Mark D Hill,Michael R Marty. Amdahl's law in the multicore era[J]. Computer, 2008, 41(7): 33-38.

[8] J L Gustafson. Reevaluating Amdahl's law[J]. Communications of ACM, 1988, 31(5): 532-533.

[9] P Christie,D Stroobandt. The interpretation and application of Rent's rule[J]. IEEE Trans on VLSI Systems, Special Issue on System-Level Interconnect Prediction, 2000, 8(6): 639-648.

- [10] Daniel Greenfield, Arnab Banerjee, Jeong-Gun Lee, Simon Moore. Implication of Rent's rule for NOC design and its fault-tolerance [J]. Networks-on-Chip, 2007, 7(9): 283 – 294.
- [11] 郝松, 都志辉, 王曼, 刘志强. 多核处理器降低功耗技术综述[J]. 计算机科学, 2007, 34(11): 259 – 263.
- [12] Wei Huang, Mircea R Stan, Karthik Sankaranarayanan, Robert J Ribando, Kevin Skadron. Many-core design from a thermal perspective [A]. IEEE Design Automation Conference [C]. Anaheim, 2008. 746 – 749.
- [13] 邓让钰, 陈海燕, 窦强, 等. 一种异构多核处理器的并行流存储结构[J]. 电子学报, 2009, 37(2): 312 – 317.
DENG Rang-yu, et al. Aparallel stream memory architecure for heterogeneous nlult-core processorz [J]. Acta Electronica

Sinica, 2009, 37(2): 312 – 317. (in Chinese)

- [14] 张饶, 武晓岛, 谢学军. 透过专利看微处理器的技术发展(四)——中国专利中的多核技术演进分析[J]. 中国集成电路, 2009, (4): 83 – 89.
- [15] 陈国良, 吴俊敏, 章锋, 章隆兵. 并行计算机体系结构[M]. 北京: 高等教育出版社, 2002.

作者简介

李文石 男, 1963 年生于黑龙江省哈尔滨市, 工学博士, 苏州大学电子信息学院微电子学系教授, 研究方向为模式分析与微电子系统设计. E-mail: lwshi@suda.edu.cn

姚宗宝 男, 1988 年出生于山东省临沂市, 苏州大学电子信息学院微电子学系硕士研究生, 研究方向为集成电路设计与测试. E-mail: yaozb1988@yahoo.com.cn