

并行分布计算

任务负载分配

算法

①

26-28

并行分布计算中随机任务负载的分配*

Distribution of Random Workloads in Parallel and Distributed Computing

陈华平 石琴[†] 黄刘生 陈国良

TP301.6

(国家高性能计算中心 合肥 230027) (合肥工业大学 合肥 230009)[‡]

Abstract As one of the hard problems in Parallel Distributed Computing(PDC), task distribution has great influence on the execution efficiency of PDC. In this paper, we first discuss the distribution model of random workloads, which is engaged to make trade-off between the balancing of the execution costs among processors and the reducing of communication and synchronization overheads. Finally we analyze the PDC meaning of the model and discuss the way of assigning these random workloads to processors.

Keywords Parallel and distributed computing, Task distribution, Model

为了提高并行分布计算的执行效率,必须尽量减少处理结点间的通讯和平衡处理结点间的计算成本,这两个目标看上去是冲突的^[1]。从减少通讯量角度出发,应把所有任务放在同一处理结点上,但很明显,这样会造成系统资源的极大浪费而且各处理结点计算成本极不平衡;从求得处理结点间的计算成本平衡及提高资源利用率出发,应把任务均匀分配到各处理结点上,但这样有可能大大增加处理结点间的通讯。因此,并行分布系统中任务分配的目标就是寻找一个任务分配方案,使得整个任务负载的执行时间最短。

目前分布系统下的大部分任务分配算法是基于网络流方法,其中最著名的是 H. Stone^[2]于1977年提出的由两个处理结点组成的并行分布系统上的最佳任务分配方法,以及 C. Lee^[3]等于1992年提出的基于二维线性网络连接并行分布系统上的最佳任务分配方法,其它绝大多数是启发式分配方法,其中最为典型的是 m-路网络流算法^[4]。本文主要讨论如何通过数学方法,来获得随机任务负载在分布系统下的最佳分配。

1 问题模型

设系统负载由相互作用的 n 个任务 {a₁, a₂, ..., a_n} 组成,并行分布计算系统由 m 个处理结点 {P₁, P₂, ..., P_m} 组成。我们用一组随机函数变量来描述并行分布系统特性及任务负载特性。假定任务是随机产生的,用随机变量 t_i 来表示任务 a_i 在结点 P_i 上的执行时间,并且这些随机变量的数学期望值为 t_i = E{t_i}。e_{jk} 是表示任务 a_j 与任务 a_k 之间是否存在数据交换的二进制随机

变量,如果 a_j 与 a_k 之间存在数据交换,那么 e_{jk} = 1, 否则 e_{jk} = 0。我们假定这随机变量的数学期望值为 e = E{e_{jk}}。

设 B_i 为分配给处理结点 P_i 的任务块子集, L(i) 为 B_i 中任务块的下标集,即 L(i) = {j: a_j ∈ B_i}。用 x_i 表示分配给结点 P_i 的任务块数目,即 x_i = |B_i| = |L(i)|, 考虑到实际情况下每个处理结点的计算能力会受到限制,并且不同处理结点的计算性能也有很大差异,我们假定每个结点 P_i 对其能运行的最大任务数有一上限 c_i。对于一个任务分配方案 B = {B₁, B₂, ..., B_m} 把 x = (x₁, x₂, ..., x_m) 称为负载分配向量,它满足下列约束条件: ∑ x_i = n, 0 ≤ x_i ≤ h_{i}, 其中 h_{i} = min{c_i, n}。我们所要解决的问题就是如何确定负载分配向量 x, 使得整个任务负载具有最短的执行时间。}}

一般说来,在整个任务负载的执行过程中,每个结点 P_i 所处的可能状态有三种:计算状态、通讯状态和同步等待状态,这三个状态的交替反复构成了每个处理结点的运行过程。相应地,每个结点 P_i 执行完成 B_i 中的 x_i 个任务所需的整个执行时间 T_i 也由三部分 T_{i1}, T_{i2}, T_{i3} 构成。

(1) 计算时间 T_{i1}。很明显, T_{i1} 为 B_i 中所有任务的计算时间和,即:

$$T_{i1} = \sum_{a_j \in B_i} t_j \quad (1)$$

(2) 通讯时间 T_{i2}。如果任务 a_j 和 a_k 被分配在不同结点上,即 a_j ∈ B_i, a_k ∈ B_l, 那么它们之间通过外部网络进行通讯,我们用随机变量 D_{jk} 来表示这些外部通讯

* 1 本文受国家863重点项目(863-306-ZD01-02-3)和安徽省自然基金的资助。

开销,它们的数学期望值为 $D_i = E\{D_{ik}\}$ 。再根据任务 a_i 和 a_k 发生数据变换的随机变量 e_{ik} , T_{12} 为:

$$T_{12} = \sum_{i \in L_{11}} \sum_{k \in L_{11}} e_{ik} D_{ik} \quad (2)$$

如果任务 a_i 和 a_k 被分配在相同结点上,那么不管它们之间是否存在数据交换,我们均可忽略其内部通讯延迟。

(3) 由于同步等待引起的延迟时间 T_{13} 。对于任务 a_i 和 a_k ($a_i \in B_i, a_k \in B_j$) 的外部数据交换,我们可通过二进制随机变量 q_{ik} 来表示 P_i 是否因该次通讯而引起同步开销,并且假定如果存在延迟,用随机变量 Q_{ik} 来表示这延迟开销。设这两组随机变量的期望值分别为 $E\{q_{ik}\} = q$ 和 $E\{Q_{ik}\} = Q$, 那么 T_{13} 为:

$$T_{13} = \sum_{i \in L_{11}} \sum_{k \in L_{11}} e_{ik} q_{ik} Q_{ik} \quad (3)$$

在上面(1)(2)(3)式中, $t_i, D_{ik}, e_{ik}, q_{ik}$ 和 Q_{ik} 这些随机变量均是相互独立的。另外,计算、通讯和同步等待这三种状态间的切换会带来一些额外开销,并且任务间的切换也会带来一些开销,我们假定已把这些开销合并到 t_i, D_{ik}, Q_{ik} 中去了。

根据以上讨论可知,处理结点 P_i 执行完成任务子集 B_i 所需的整个时间 T_i 为:

$$T_i = T_{11} + T_{12} + T_{13} = \sum_{j \in L_{11}} t_j + \sum_{j \in L_{11}} \sum_{k \in L_{11}} e_{jk} D_{jk} + \sum_{j \in L_{11}} \sum_{k \in L_{11}} e_{jk} q_{jk} Q_{jk} \quad (4)$$

对等式两边同取数学期望值:

$$T_i = E(T_i) = \sum_{j \in L_{11}} E(t_j) + \sum_{j \in L_{11}} \sum_{k \in L_{11}} E(e_{jk} D_{jk}) + \sum_{j \in L_{11}} \sum_{k \in L_{11}} E(e_{jk} q_{jk} Q_{jk})$$

进一步化简可得:

$$T_i = t_i \sum_{j \in L_{11}} 1 + eD_i \sum_{j \in L_{11}} \sum_{k \in L_{11}} 1 + eqQ \sum_{j \in L_{11}} \sum_{k \in L_{11}} 1$$

根据 $|L(i)| = x_i$, 可把 T_i 定义成 x_i 的函数,上式化为:

$$T_i(x_i) = t_i x_i + eD_i x_i (n - x_i) + eqQ x_i (n - x_i)$$

进一步化简,上式可化为:

$$T_i(x_i) = [-eD_i - eqQ]x_i^2 - [-t_i - enD_i - eqQn]x_i \quad (5)$$

$$\text{令 } a_i = -eD_i - eqQ \quad (6)$$

$$b_i = -t_i - enD_i - eqQn \quad (7)$$

那么(5)式可化为:

$$T_i(x_i) = a_i x_i^2 - b_i x_i \quad (8)$$

假定每个处理结点同时开始执行,那么最后执行完成的处理结点所需的执行时间即为整个系统任务负载的执行时间 $T(x)$, 也就是说,

$$T(x) = \text{MAX}_i \{a_i x_i^2 - b_i x_i\} = \text{MAX}_i \{T_i(x_i)\} \quad (9)$$

最佳任务分配问题就是如何确定一个负载分配向量 $x^* = (x_1^*, x_2^*, \dots, x_n^*)$, 使得 $T(x)$ 具有最小值, 设最小值为 T^* , 那么

$$\begin{aligned} T^* &= T(x^*) = \text{MIN}_x T(x) = \text{MIN}_x \text{MAX}_i \{T_i(x_i)\} \\ &= \text{MIN}_x \text{MAX}_i \{a_i x_i^2 - b_i x_i\} \end{aligned} \quad (10)$$

其中 x_i 为整数, 并且 $\sum x_i = n, 0 \leq x_i \leq h_i \leq n$ 。

2 问题模型的并行分布计算意义

下面对函数 $T_i(x_i) = a_i x_i^2 - b_i x_i$ ($x_i \in [0, h_i]$) 进行分析, 其中设 $y_i = T_i(h_i) = a_i h_i^2 - b_i h_i$, 对 $T_i(x_i)$ 分别求一阶导数和二阶导数可得: $T_i'(x_i) = 2a_i x_i - b_i, T_i''(x_i) = 2a_i$ 。很明显, 当 $x_i < b_i/2a_i$ 时, $T_i'(x_i) < 0, T_i(x_i)$ 是单调上升的, $x_i > b_i/2a_i$ 时, $T_i'(x_i) > 0, T_i(x_i)$ 是单调下降的。所以必须根据 h_i 的大小分别讨论。

(I) ($h_i < b_i/2a_i$): 这时 $x_i \leq h_i < b_i/2a_i, T_i'(x_i) > 0, T_i(x_i)$ 在 $[0, h_i]$ 上恒单调上升。

(II) ($h_i > b_i/2a_i$): 这时当 $x_i \in [0, b_i/2a_i]$ 时, $T_i'(x_i) < 0, T_i(x_i)$ 是单调上升的, 当 $x_i \in (b_i/2a_i, h_i]$ 时, $T_i'(x_i) > 0, T_i(x_i)$ 是单调下降的, 所以这时 $T_i(x_i)$ 在 $[0, h_i]$ 上是非单调的。

$T_i(x_i) = a_i x_i^2 - b_i x_i = 0$ 有两个根, 即 $x_{i1} = 0, x_{i2} = b_i/a_i$ 。虽然情形(I)中 h_i 可能大于 $b_i/2a_i$, 但 $h_i < b_i/a_i$, 这是因为根据(6)(7)式, $b_i/a_i = t_i/(eD_i + eqQ) + n > n > h_i$ 。

从并行分布计算意义上我们也可以这样理解, 如果给 P_i 上多分配些任务, 那么由此可以减少任务间的外部通讯及同步开销, 从减少整个通讯开销和同步开销的角度出发, 这样是有利的, 但由于新增任务本身计算时间的加入, 使得 $T_i(x_i)$ 一开始还是随任务数目的增加而增长, 但 $T_i'(x_i) = 2a_i < 0$, 也就是说 $T_i(x_i)$ 的增加速度逐步减小, 即 $T_i(x_i)$ 随 x_i 的增加而呈亚线性增长。如果继续不断增加 P_i 上的任务数目, 那么由此减少的外部通讯及同步开销所获得的收益, 逐渐逼近或超过由此带来的新增计算时间, 任务数 $x_i = b_i/2a_i$ 就是它们之间的平衡点。如果 $h_i < b_i/2a_i$, 说明由于处理结点 P_i 受到计算能力的限制, 其上所分配的任务数目到达不了这个平衡点; 如果 $h_i > b_i/2a_i$, 那么 P_i 上的任务数目超过 $b_i/2a_i$ 时, 由于外部通讯及同步开销的大大减小, 使得 $T_i(x_i)$ 反而随任务数的增加而减小。

3 负载分配向量的确定

下面我们先按 x_i 是连续变量这个假定求解, 然后在此基础上求整数解, 根据以上讨论, 我们可把所有的 $T_i(x_i)$ ($1 \leq i \leq m$) 组成的函数集合 R 分为两部分, 一部分为单调递增函数, 记为 R_1 , 对应于上一节的第(I)种情形; 另一部分为非单调递增函数, 记为 R_2 ($R_2 = R - R_1$), 对应于上一节的第(II)种情形。如果 R_2 非空, 那么设 u 是 R_2 中使 $T_i(h_i)$ 取得最小值的下标, 即定义: $y_u = T_u(h_u) = \text{MIN}_{T_i(x_i) \in R_2} \{T_i(h_i)\} = \text{MIN}_{T_i(x_i) \in R_2} \{y_i\}$ (11)

如果 R_2 中取得最小值 y_u 的函数多于一个, 那么任取其中一个即可。如果 $R_2 - \{T_i(x_i)\}$ 非空, 那么设 v

是 $R_1 = \{T_1(x_i)\}$ 中使 $T_1(h_1)$ 取得最小值的下标,即定义:

$$y = T_1(h) = \min_{T_1(x_i) \in R_2 - \{T_1(x_0)\}} \{T_1(h_i)\} = \min_{T_1(x_i) \in R_2 - \{T_1(x_0)\}} \{y_i\} \quad (12)$$

同样,如果 $R_2 = \{T_2(x_i)\}$ 中取得最小值 y_0 的函数多于一个,那么任取其中一个即可。这样, y_0 和 y_0' 是非单调递增函数集 R_2 中最小的两个值,并且 $y_0 \leq y_0'$ 。

给 R 中的每个 $T_i(x_i)$ 定义如下反函数:

$$x_i = G_i(y) = \begin{cases} T_i^{-1}(y) & 0 \leq y \leq y_i \\ h_i & y > y_i \end{cases} \quad (13)$$

根据上一节对二次曲线 $T_i(x_i)$ 的系数 a_i 和 b_i 的分析可知, R_1 中的 $T_i(x_i)$ 在区间 $y \in [0, y_i]$ 上是一一对应的,所以它们都存在严格定义的反函数 $T_i^{-1}(y) = (b_i + \sqrt{b_i^2 + 4a_i y}) / 2a_i$ 。对于 R_2 中的 $T_i(x_i)$,在点 $y = y_i$ 上, $T_i^{-1}(y_i)$ 存在两个值,一个是 $x_{i1} = h_i$,另一个是 $a_i x_{i2} - b_i x_{i2} = y_i$ 的小根,由于 $a_i < 0, b_i < 0$,所以小根值为 $x_{i2} = (b_i + \sqrt{b_i^2 + 4a_i y_i}) / 2a_i$ 。如果我们规定 x_{i2} 作为 $T_i^{-1}(y_i)$ 的唯一值,那么这时, $T_i^{-1}(y)$ 在 $y \in [0, y_i]$ 上也是严格定义的。

很明显,对于 R_1 中的函数,上面定义的 $T_i^{-1}(y)$ (即 $G_i(y)$) 是连续的,而对于 R_2 中的函数 $G_i(y)$ 在 $y = y_i$ 处有一跳跃,即: $G_i(y_i) = (b_i + \sqrt{b_i^2 + 4a_i y_i}) / 2a_i, G_i(y_i^+) = h_i$ 。

我们定义 $S(y) = \sum G_i(y) (y \geq 0)$, 由于 $G_i(y)$ 是 y 的非减函数,所以 $S(y)$ 也是 y 的非减函数,对于 R_2 中函数 $T_i(x_i)$ 的反函数 $G_i(y)$, $S(y)$ 在其对应每个 y_i 处出现一个跳跃性的不连续,它们对应了 $G_i(y)$ 函数的不连续,另外,在两个不连续点的区间内函数 $S(y)$ 是连续且单调递增的,具有严格定义的反函数 $S^{-1}(\cdot)$,其中 $S^{-1}(\cdot)$ 的自变量为构成系统负载的任务总数。

如果 R_2 至少包含两个函数,那么 $S(y_0^+) = \sum G_i(y_0^+) = G_1(y_0^+) + G_2(y_0^+) + \sum_{i=3, \dots, m} G_i(y_0^+)$, 由于 $y_0^+ \geq y_0 \geq y_0'$, 根据(13)的定义, $G_1(y_0^+) = h_1$, 从而上式可化为: $S(y_0^+) = h_1 + h_2 + \sum_{i=3, \dots, m} G_i(y_0^+)$ 。由于 $h_2 > (b_2 / 2a_2), h_3 > (b_3 / 2a_3)$, 所以 $S(y_0^+) > (b_2 / 2a_2) + (b_3 / 2a_3)$, 根据上一节说明, $b_2/a_2 > n$, 因此 $S(y_0^+) > n$ 。这样,在 R_2 非空的一般情况下, $S(y)$ 最多只可能包含 y_0 和 y_0' 处两个断点,函数 $S(y)$ 在区间 $[0, y_0]$ 及 (y_0^+, y_0') 上连续且单调递增,相应地 $S^{-1}(n)$ 在 $[0, S(y_0)]$ 及 $(S(y_0^+), S(y_0'))$ 上连续且单调递增,并且 $S^{-1}(n)$ 与任务总数 n 构成了——对应关系,而在 $(S(y_0), S(y_0^+)]$ 及 $(S(y_0'), S(y_0'))$ 这两个区间上, $S^{-1}(n)$ 与任务总数 n 不满足——对应条件。

如果点 $(n, S^{-1}(n))$, 落在 $(S(y_0), S(y_0^+)]$ 或者 $(S(y_0'), S(y_0'))$ 非——对应区间上时,可分别令 $x_i^* = h_i$ 或者 $x_i^* = h_i$, 对于其它的 x_i^* , 直接选择满足 $(T_i(x_i^*) \leq y_0)$ 或者 $(T_i(x_i^*) \leq y_0')$, 并且 $\sum x_i^* = n$ 的整数解即可,这时 $T(x)$ 取得最小值,其值为 $T^* = y_0$ 或 $T^* = y_0'$ 。如果点 $(n, S^{-1}(n))$ 落在 $[0, S(y_0)]$ 或 $(S(y_0^+), S(y_0'))$ ——对应区间上,或者 R_2 为空(这时 $S(y)$ 均由 R_1 中的单调函数构成), $S(y)$ 和 $S^{-1}(n)$ 在任何点上均是一一对应并且单调递增的,这时 $T^* = S^{-1}(n)$ 是其最小值,令 $x_i^* = G_i(T^*) = G_i(S^{-1}(n)) (1 \leq i \leq m)$, 这时 $T(x)$ 取得最小值,其最小值为 $(T^*) = S^{-1}(n)$, 但 x_i^* 并不一定是整数,可对非整数解 $x^* = (x_1^*, x_2^*, \dots, x_m^*)$ 进一步调整,使其每个元素 x_i^* 均为整数。

设调整后的整数解为 $r = \{r_1, r_2, \dots, r_m\}$ 。一开始 r_i 取 x_i^* 的整数部分,因此 $r_i \leq x_i^*$, 结果 $\sum r_i \leq \sum x_i^* = n$, 并且 $\sum r_i$ 一般都要比 n 小,除非初解 $x_i^* (1 \leq i \leq m)$ 刚好均为整数,然后每次选择 r 中的一个元素 r_i 进行加1操作,选择的元素满足条件 $T_i(r_i + 1) \leq \forall_i (i \neq z) (T_i(r_i + 1))$, 经过 $n - \sum_{i=1, \dots, m} r_i$ 次加1操作后, $\sum_{i=1, \dots, m} r_i = n$, 下面是具体的调整算法。

```

procedure integer_adjust(x*, r)
begin
  for (1 ≤ i ≤ m) do r_i = round(x_i*); endfor
  k = n - ∑_{i=1}^m r_i;
  for (1 ≤ j ≤ k) do
    选择满足 T_i(r_i + 1) ≤ ∑_{i ≠ z} (T_i(r_i + 1)) 的 r 向量下标 z
    r_z = r_z + 1;
  endfor
end
    
```

参考文献

- 1 Hesham El-Rewini, et al. Task scheduling, PTR Prentice Hall, Englewood Cliffs, New Jersey, 1994, 07532
- 2 陈华平, 计永昶, 陈国良. 分布式动态负载平衡调度的一个通用模型. 软件学报, 1998, 9(1): 25~29
- 3 Feitelson D G, Rudolph L. Toward Convergence in job schedulers for parallel super-computers. In: IPPS '96 Workshop on Job Scheduling Strategies for Parallel Processing. Hawaii, 1996. 1~9
- 4 Stone H. Multiprocessor scheduling with the aid of network flow algorithms. IEEE Trans. Software Eng., 1977
- 5 Lee C, et al. Optimal task assignment in linear array networks. IEEE Trans. on Computers
- 6 Lo V. Heuristic algorithms for task assignment in distributed systems. IEEE Trans. on Computers, 1988(11)