

SACG: Simulated Annealing Combined with Greedy Algorithm

An Approximation Solution to Resource Scheduling Problem in Hadoop

Zhenran Xiao, Zichun Ye, Mabiao Long

Student ID:{520030910281, 520030910302, 520030910308}

{xiaozhenran, alchemist, albertlong007} @ sjtu.edu.cn

Department of Computer Science and Engineering,
Shanghai Jiao Tong University, Shanghai, China

摘要 The course project focuses on resource scheduling problem in Hadoop. This document first introduce the background and two versions of resource scheduling on single or multiple hosts. Then it designs an algorithm based on simulated annealing algorithm and greedy algorithm to give an approximate solution to the problems of two versions. It also does some experiments on parameters in the algorithm and analyses the complexity and other quality of the algorithm.

Keywords: Distributed Computing System, Resource Scheduling, Hadoop, Simulated Annealing Algorithm, Greedy Algorithm, NP-complete problem

1 Background and Motivation

Hadoop is an open-source software framework for storing data and running applications on clusters of commodity hardware. In Hadoop, the data of a running job can be divided into several data blocks, stored on different hosts. Each host has one or multiple CPU cores to process data blocks. In this project, our goal is to achieve effective parallel computing. That is to say, given a set of jobs with massive data, we need to design a resource scheduling algorithm to minimize the overall executing time of all jobs.

2 Problem Formulation and Explanation

In this section, we will formulate the problem and explain some definitions.

2.1 A Simplified Version with Single Host

First, let us consider a simple case with a single host storing all data blocks. To simplify the problem, we give the following assumptions and specifications.

1. There are n jobs that need to be processed by a single host, which has m CPU cores of the same computing capability. Let J be the set of jobs, and C be the set of cores for the host, where $J = \{job_0, job_1, \dots, job_{n-1}\}$, and $C = \{c_0, c_1, \dots, c_{m-1}\}$ (The labels of variables start from 0 because we use C/C++ source codes in the following tasks).

2. We treat data block as the smallest indivisible unit in our project, while a job can be divided into multiple data blocks with different sizes for storage and computing. Assume job_i is split into n_i data blocks, denoted by $B^i = \{b_0^i, b_1^i, \dots, b_{n_i-1}^i\}$. For block b_k^i of job_i , define its size as $size(b_k^i)$.
3. Assume job_i is assigned to e_i cores for processing, and naturally $e_i \leq n_i$. That is to say, one core can process multiple data blocks of one job sequentially. Let $B_j^i \subseteq B^i$ denote the set of data blocks of job_i allocated to core c_j , and $B_j^i \cap B_{j'}^i = \emptyset$ if $j \neq j'$ (they should be disjointed).
4. For job_i , the processing speed of its data blocks is s_i when job_i is assigned to a single core. However, when multiple cores process job_i in parallel, the computing speed of each core all decays because of some complicated interactions. We formulate such speed decay effect caused by multi-core computation as a coefficient function $g(\cdot)$ with respect to core number e_i , as described in Equation (1):

$$g(e_i) = 1.00 - \alpha \cdot (e_i - 1), \quad \text{for } 1 \leq e_i \leq 10, \quad (1)$$

where α is a decay factor satisfying $0 < \alpha < 1$, and usually the number of cores for processing a single job is no more than 10. Then, the speed of each core can be rewritten as $s_i \cdot g(e_i)$ for job_i respectively. (Note that although the speed of each core decays, the overall processing time using e_i cores in parallel should be faster than that of using just one core. Otherwise we do not need to implement parallel computing. Thus the setting of α should guarantee this principle.)

Correspondingly, the processing time tp_j^i of core c_j for job_i can be expressed as Equation (2):

$$tp_j^i = \frac{\sum_{b_k^i \in B_j^i} size(b_k^i)}{s_i \cdot g(e_i)}. \quad (2)$$

5. For consistency issues, if we assign a job to multiple cores, all cores must start processing data blocks at the same time. If one or several cores are occupied by other affairs, then all other cores should wait for a synchronous start, and keep idle. It means that the processing of job_i for every core should all start at time t_i , whereas their processing duration might be different. Let tf_j^i be the finishing time of core c_j for job_i , which is calculated by Equation (3):

$$tf_j^i = t_i + tp_j^i. \quad (3)$$

However, the occupied cores of job_i are released synchronously when the computing process of the last data block is finished. Thus the finishing time $tf(job_i)$ of job_i is given as Equation (4):

$$tf(job_i) = \max_{c_j \in C} tf_j^i, \quad \text{for } c_j \in C. \quad (4)$$

2.2 A Comprehensive Version among Multiple Hosts

In this section, we consider a more complex situation, where we need to schedule resources among multiple hosts. Now the data transmission process between pairwise hosts should be taken into consideration. The data blocks of jobs could be initially stored on different hosts, but one data block can only be initially stored on one specified host. If data block b_k^i and its assigned computing core c_j are not on the same host, b_k^i will be transmitted to the host containing c_j (Here we assume

that the bandwidth between hosts is sufficient for data transmission). The transmission process will influence the finishing time of jobs, further affecting the resource scheduling process.

Besides the descriptions and specifications of Section 2.1, here are more notations and explanations.

1. Assume we have q hosts $H = \{h_0, h_1, \dots, h_{q-1}\}$, and host h_l has m_l cores (may have different number of cores). Let C_l be the set of cores on host h_l , $C_l = \{c_0^l, c_1^l, \dots, c_{m_l-1}^l\}$. Easy to see, $\sum_{l=0}^{q-1} m_l = m$.
2. If core c_j^l on host h_l computes a data block b_k^i of job_i which is initially stored on another host h_l' , then b_k^i needs to be transmitted from h_l' to h_l at a transmission speed s_t (this speed is fixed in our system). An example is shown in Figure 1, where hosts h_0 and h_1 both have two cores, and many jobs need to be processed. Core c_0^1 on host h_1 is assigned to compute the data block b_2^1 of job_1 , which is initially stored on host h_0 . In this case, b_2^1 needs to be transmitted from h_0 to h_1 at a transmission speed s_t first, and then be computed by c_0^1 . Whenever b_2^1 starts transmission, other cores can work in parallel to process job_1 .

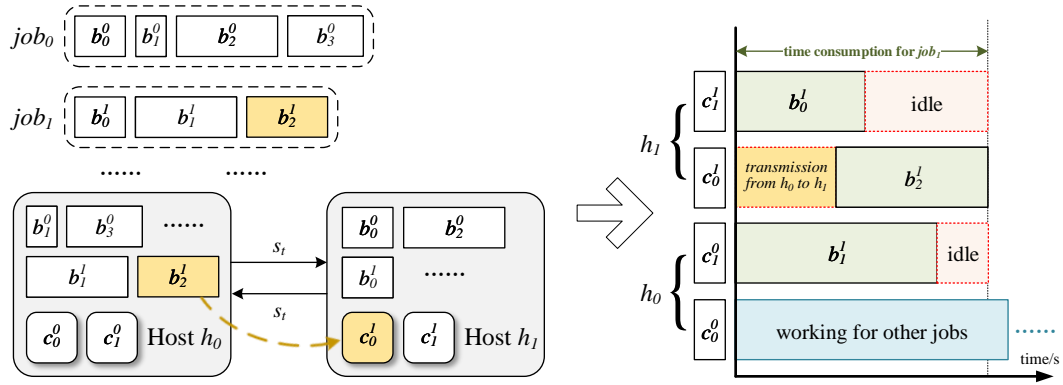


图 1. An example of data transmission between 2 hosts

3. Any core cannot call for data transmission when it is calculating other data blocks. Likewise, a core cannot start computing any data block until this block is ready, i.e. initially on the same host or transmitted from a remote host to the local host. For example, the core c_0^1 on host h_1 must wait for the data transmission of block b_2^1 from host h_0 to h_1 , and then start computation. What is more, the transmission time of b_2^1 from h_0 to h_1 affects the finishing time of job_0 , further affecting the finishing time of the whole system.

For core c_j^l on host h_l , let $\tilde{B}_{l_j}^i$ be the set of data blocks of job_i allocated to c_j^l but not initially stored on host h_l . All the data blocks in $\tilde{B}_{l_j}^i$ need to be transmitted to host h_l before computing. Let $B_{l_j}^i$ be the set of data blocks of job_i allocated to core c_j^l . Then, the processing time $tp_{l_j}^i$ of core c_j^l for job_i can be reformulated as Equation (5):

$$tp_{l_j}^i = \frac{\sum_{b_k^i \in \tilde{B}_{l_j}^i} size(b_k^i)}{s_t} + \frac{\sum_{b_k^i \in B_{l_j}^i} size(b_k^i)}{s_i \cdot g(e_i)}. \quad (5)$$

4. If the processing of job_i starts at time t_i , then the finishing time of core c_j^l for job_i is

$$tf_{lj}^i = t_i + tp_{lj}^i.$$

Then the finishing time $tf(job_i)$ of job_i is formulated as:

$$tf(job_i) = \max_{c_j^l} tf_{lj}^i, \text{ for } c_j^l \in C.$$

3 Problem Solving

Based on the descriptions in Section 2, we want to design a resource scheduling algorithm to minimize the overall finishing time of all jobs, whose objective function is shown as:

$$\min \max_{job_i} tf(job_i), \text{ for } job_i \in J.$$

And we hope the algorithm works for both versions of the problem.

3.1 SA: Simulated Annealing Algorithm

Simulated annealing algorithm was proposed by N. Tropolis scholar in 1953. It comes from the solid annealing principle, benefiting from the research results of statistical mechanics of materials. The algorithm is based on probability, mainly used to solve optimal solution problems, such as function extremum problems. In statistical mechanics of materials, different structures of particles in materials correspond to different energy levels of particles. At high temperature, particles have high energy and can move and rearrange freely. Low temperature particles have low energy. If the particles are cooled slowly from high temperature, namely annealed, they can reach thermal equilibrium at each temperature. When the system is completely cooled, a low-energy crystal is eventually formed.

The Principle of SA Algorithm Suppose the energy of the material in $state_i$ is $E(i)$, then the law of the material entering $state_j$ from $state_i$ at temperature T is as follows:

1. If $E(j) \leq E(i)$, we accept the $state_j$, namely the material enters $state_j$ from $state_i$.
2. If $E(j) > E(i)$, it is probable to accept the $state_j$. The probability is

$$e^{\frac{E(i) - E(j)}{KT}}$$

K is Boltzmann's constant, and we usually make K equal to 1 in fact. T is the temperature of the material.

When the temperature is very low, we think it is highly possible that the material is in the minimum energy state.

The Flow of SA Algorithm The flow of the algorithm is following:

1. Generating the initial state $state_0$. Set the initial temperature T_0 and the final temperature T_e . Set the rate of temperature fall β . Set the current number of iterations i and the maximum number of iterations M .
2. In each iteration, $T_i = \beta \cdot T_{i-1}$. Move the particles to generating a new state $state_i$ by some scheme.
3. Calculate the change of the energy by an objective function:

$$\delta E = E(state_i) - E(state_{i-1})$$

If $\delta E \leq 0$, accept the $state_i$. If not, accept the $state_i$ by probability $e^{\frac{-\delta E}{T_i}}$.

4. If the number of iterations arrives M or the temperature is lower than T_e , stop the iteration and obtain the final state as the result.

3.2 The Overview of Our Algorithm

In the resource scheduling problem, we can regard every feasible solution as a state, the finishing time as the energy, the blocks as the particles. We generate a feasible solution as the initial state at the beginning and set the initial temperature, final temperature and maximum number of iterations for the problem. Next, in each iteration we move the blocks by some schemes to generate a new solution as a new state. However, the moving scheme needs to be limited by the problem formulation in Section 2. According to the new state, we can calculate its energy and decide whether to accept it. When the iteration stops, we can get a final state as the approximation solution — we think it is highly possible that this solution is good enough.

Considering the following restrictions:

1. If we assign a job to multiple cores, all cores must start processing data blocks at the same time. If one or several cores are occupied by other affairs, then all other cores should wait for a synchronous start, and keep idle.
2. The occupied cores of a job are released synchronously when the computing process of the last data block is finished.

We can easily know that, for one job, if we have decided which cores will be allocated to it, it will be like a cube existing on the simulated image of the solution no matter where it is on the image. The width of it is the number of the cores allocated to the job. The length is the longest processing time on these cores. (Maybe the width of the cube is not connected together, but it doesn't matter too much.)

Therefore, we can use Greedy Algorithm for a job to make the length of its cube shorter, namely make the processing time for one job shorter. In both versions, Greedy Algorithm can give a good enough solution for one job.

When the cubes are fixed, we can just put them in the execution queues by some order, like the game *Tetris*.

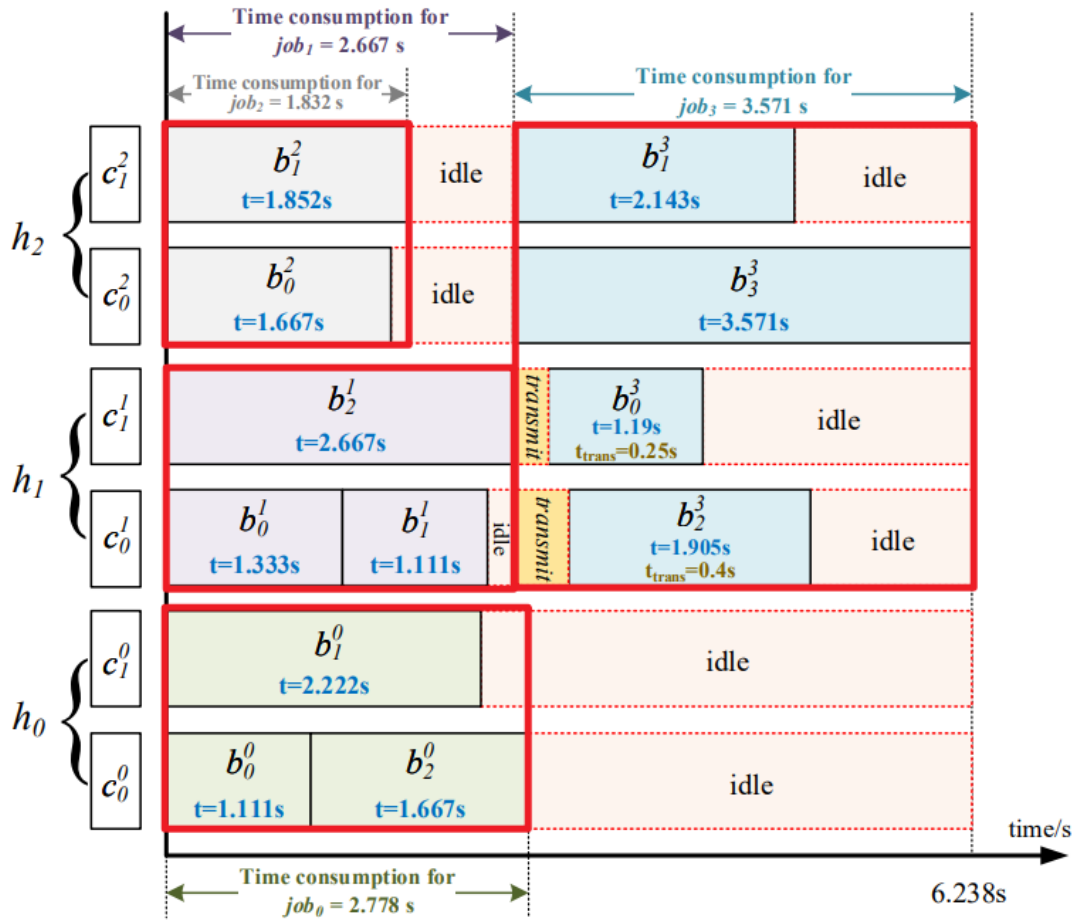


图 2. A simulated image of a solution

So the overview of our algorithm is:

1. Regard each job as a cube, generate the initial solution namely initial state in the Simulated Annealing Algorithm.
2. Move these cubes and change their width to generate a new state.
3. Considering the fixed width of cubes, use Greedy Algorithm to make their length shorter.
4. Simulate the calculating process to obtain the final finishing time, namely the energy of the state.
5. Perform the annealing, decide whether to accept the new state.
6. If the temperature and the number of iterations don't reach the standard, go to 2.
7. Otherwise, end the iteration, the final state is the final solution.

Here is the pseudo code:

Algorithm 1: *schedule()*

Input: all the information about hosts, cores, jobs, blocks

Output: the final state

```

1  init();
2   $E \leftarrow 0$ ;
3   $tempE \leftarrow 0$ ;
4   $T \leftarrow 1000$ ;
5   $T_e \leftarrow 0.0001$ ;
6   $\beta \leftarrow 0.999$ ;
7   $\delta E \leftarrow 0$ ;
8   $M \leftarrow 100000$ ;
9  repeat
10   move();
11    $E \leftarrow finishTime()$ ;
12    $tempE \leftarrow temp\_finishTime()$ ;
13    $\delta E \leftarrow tempE - E$ ;
14   if  $\delta E < 0$  then
15     | Accept the new state;
16   end
17   else
18     | There is a probability of  $\exp(\frac{-\delta E}{T})$  to accept the new state;
19   end
20    $T \leftarrow T * \beta$ ;
21   if  $T < T_e$  then
22     | break;
23   end
24    $M \leftarrow M - 1$ ;
25 until  $M = 0$ ;

```

3.3 The Data Structure Mainly Used in the Program

We mainly use these data structure in the program:

```

1  vector<int> hostCore;           // The number of cores for each host
2  vector<int> jobBlock;          // The number of blocks for each job
3  vector<double> Sc;             // Speed of calculation for each job
4  vector<vector<double>> dataSize; // Job-> block number-> block size
5  vector<vector<int>> location;    // Job-> block number-> block location (host
                                 // number)
6
7  vector<double> jobFinishTime;  // The finish time of each job
8

```

```

9      vector<vector<tuple<int, int, int>>> runLoc;
10     // 1. Block perspective: job number->block number->(hostID, coreID,rank), rank=1
        means that block is the first task running on that core of that host
11
12     vector<vector<vector<tuple<int, int, double, double>>>> hostCoreTask;
13     // 2. Core perspective: host->core->task-> <job,block,startRunningTime,
        endRunningTime>
14
15     vector<vector<double>> hostCoreFinishTime; // host->core->finishTime
16
17     /* 3. This representation of a solution is used for SAA */
18     vector<int>jobCore;           // The number of cores allocated to each job.
19     vector<int>jobOrder;          //jobOrder[1] stores the first job
20     vector<vector<tuple<int, int, vector<int> > > > total_sol;
21     //the information of the number of cores are stored in jobCore
22     //job-> core ->(hostID, coreID, vector of blockID)
23     vector<vector<int> >blockLocation;
24     //job-> blockID-> block destination (hostID)
25
26     //modify the solution temporarily
27     vector<int>temp_jobCore;
28     vector<int>temp_jobOrder;
29     vector<vector<tuple<int, int, vector<int> > > > temp_total_sol;
30     vector<vector<int> >temp_blockLocation;

```

We will use the third representation of solutions during the annealing. It means different *jobOrder*, *jobCore* and *total_sol* together represent different states. When we generate a new state, we can use it to obtain the finishing time, namely the energy of the state by simulating the calculating process. When we get the final state, we can transform this representation to the block perspective and core perspective solution.

3.4 The Scheme to Generate the Initial State

We already know for one job, if we have decided which cores will be allocated to it, it will be like a cube existing on the simulated image of the solution no matter where it is on the image. So we only need to generate the *jobOrder* and the cores which are allocated to each job. We use random method to generate the two information above. The flow of the scheme is as followed:

1. Initialize the *jobOrder* to 1, 2, 3, ..., *n*.
2. For each job, randomly decide the number of cores allocated to it, but restrict that the cores can't more than its blocks.
3. For each job, randomly decide which cores will be allocated to it. The probability of each core that will be chosen is equal.

4. For all the cores allocated to one job, put its blocks into the queues of cores, one for each by order, until the blocks are all put in.

The step 4 above is meant to make the energy of the initial state higher, to enlarge the effect of moving at the beginning.

Here is the pseudo code:

Algorithm 2: *init()*

Input: all the information about hosts, cores, jobs, blocks

Output: the initial state

```

1  jobOrder  $\leftarrow (1, 2, \dots, n)$  ;
2  for i  $\leftarrow 0$  to n-1 do
3      jobCore[i]  $\leftarrow \text{rand}() \% \min(10, \text{jobBlock}[i], m) + 1$  ;
4      for j  $\leftarrow 0$  to jobCore[i] - 1 do
5          pick a core that hasn't been allocated to jobi;
6          total_sol[i][j]  $\leftarrow (\text{host\_id}, \text{core\_id}, \text{process queue for this job})$ ;
7      end
8      for k  $\leftarrow 0$  to jobBlock[i] - 1 do
9          total_sol[i][k%jobCore[i]].queue.push(k);
10         blockLocation[i][k]  $\leftarrow \text{total\_sol}[i][k \% \text{jobCore}[i]].\text{host\_id}$ ;
11     end
12 end
13 temp_jobOrder  $\leftarrow \text{jobOrder}$  ;
14 temp_jobCore  $\leftarrow \text{jobCore}$  ;
15 temp_total_sol  $\leftarrow \text{total\_sol}$  ;
16 temp_blockLocation  $\leftarrow \text{blockLocation}$  ;

```

3.5 The Scheme to Move Blocks to Generate New State

The same reason, we only need to change the *job_order* and the cores allocated. We have two strategys to change them:

1. Randomly choose two adjacent jobs, change their order. This will happen at a probability of P_{co} . The probability of each job being chosen is equal.
2. Randomly choose one job, randomly decide the number of cores allocated to it, but restrict that the cores can't be more than its blocks and ten, then randomly decide which cores will be allocated to it. This will happen at a probability of P_{rc} . The probability of each core being chosen is equal.

Here is the pseudo code:

Algorithm 3: *move()*

Input: the last state

Output: the new state

```

1  $p \leftarrow rand()\%2;$ 
2 if  $n \geq 2$  and  $p \neq 0$  then
3    $i \leftarrow rand()\%(n - 1) + 1;$ 
4    $swap(temp\_jobOrder[i], temp\_jobOrder[i + 1]);$ 
5 end
6  $p \leftarrow rand()\%2;$ 
7 if  $p \neq 0$  then
8    $i \leftarrow rand()\%n;$ 
9    $temp\_jobCore[i] \leftarrow rand()\%min(10, jobBlock[i], m) + 1;$ 
10  for  $j \leftarrow 0$  to  $temp\_jobCore[i] - 1$  do
11    pick a core that hasn't been allocated to  $job_i$ ;
12     $temp\_total\_sol[i][j] \leftarrow (host\_id, core\_id, \text{process queue for this job});$ 
13  end
14   $greedy\_allocation(job_i);$ 
15 end

```

3.6 Make the length of Cubes Shorter by Greedy Algorithm

For one job, when the cores allocated to it have been decided, namely the width of the cube has been fixed, we use Greedy Algorithm to make its length shorter:

1. Sort all blocks of the job in descending order of data size.
2. From the biggest block, compare which core to place it on can minimize the current processing time. Put it on that core.
3. Decide all blocks position by 2.
4. During the process, we can update the *blockLocation*.

Here is the pseudo code:

Algorithm 4: *greedy_allocation*(job_i)

Input: job_i

Output: the schedule scheme for this job

```

1 sort all blocks of  $job_i$  in descending order of data size, use
    $jobBlock\_sorted.push((block\_id, datasize));$ 
2  $coreLength \leftarrow (0, 0, \dots, 0);$ 
3 for  $k \leftarrow 0$  to  $jobBlock[i]$  do
4    $b \leftarrow jobBlock\_sorted[k].block\_id;$ 
5    $min \leftarrow MAX\_LENGTH;$ 
6   for  $j \leftarrow 0$  to  $temp\_jobCore[i]$  do
7     if  $temp\_total\_sol[i][j].host\_id = location[i][b]$  then
8        $ratio \leftarrow 1$ 
9     end
10    else
11       $ratio \leftarrow 1 + Sc[i] * g(jobCore[i])/St;$ 
12    end
13    if  $coreLength[j] + dataSize[i][b] * rate < min$  then
14       $min \leftarrow coreLength[j] + dataSize[i][k] * ratio;$ 
15       $index \leftarrow j;$ 
16    end
17  end
18   $coreLength[index] \leftarrow min;$ 
19   $temp\_total\_sol[i][index].queue.push(b);$ 
20   $temp\_blockLocaton[i][b] = temp\_total\_sol[i][index].host\_id;$ 
21 end

```

3.7 The Scheme to Calculate the Energy

We will simulate the calculating process to obtain the final time, namely the energy. The simulating scheme is following:

1. Initialize a data structure tf . It is a two-dimensional array, meaning $host- > core- > releaseTime$.
2. According to $jobOrder$ and $total_sol$, put the cube of each job into the execution queue one by one, like playing *Tetris*.
3. For every next job, first find the core whose end time is the latest among all the cores allocated to this job and set its end time as the start time of this job.
4. Then add this job's cube length, namely the total processing time of the job, to the start time, obtain finish time of the job, update tf for all the cores allocated to this job.
5. Go to 3.
6. If all the jobs are scheduled, end the cycle.

After simulating, we obtain the final finish time as the energy of the state.

Here is the pseudo code:

Algorithm 5: *finishTime()*

Input: *jobOrder, jobCore, total_sol*

Output: *finishtime*

```

1 for  $i \leftarrow 1$  to  $n$  do
2    $job_i \leftarrow jobOrder[i]$ ;
3   for  $j \leftarrow 0$  to  $jobCore[job_i] - 1$  do
4      $host\_id \leftarrow total\_sol[job_i][j].host\_id$ ;
5      $core\_id \leftarrow total\_sol[job_i][j].core\_id$ ;
6      $block\_id \leftarrow total\_sol[job_i][j].queue$ ;
7      $tp_j \leftarrow 0$ ;
8     for  $k = 0$  to  $block\_id.size() - 1$  do
9       if  $location[job_i][block\_id[k]] = blockLocation[job_i][block\_id[k]]$  then
10         $tp_j \leftarrow tp_j + data\_Size[job_i][block\_id[k]] / (Sc[job_i] \cdot g(jobCore[job_i]))$ ;
11      end
12      else
13         $tp_j \leftarrow tp_j + data\_Size[job_i][block\_id[k]] / St +$ 
14           $data\_Size[job_i][block\_id[k]] / (Sc[job_i] \cdot g(jobCore[job_i]))$ ;
15      end
16    end
17     $tp_i[j] \leftarrow tp_j$ ;
18    if  $t_i < tf[host\_id][core\_id]$  then
19       $t_i \leftarrow tf[host\_id][core\_id]$ ;
20    end
21  end
22   $tf\_job_i \leftarrow t_i + \max(tp_i[])$ ;
23  update  $tf[host\_id][core\_id]$  for all cores allocated to  $job_i$ ;
24  if  $finishtime < tf\_job_i$  then
25     $finishtime \leftarrow tf\_job_i$ ;
26  end
27 return  $finishtime$ ;

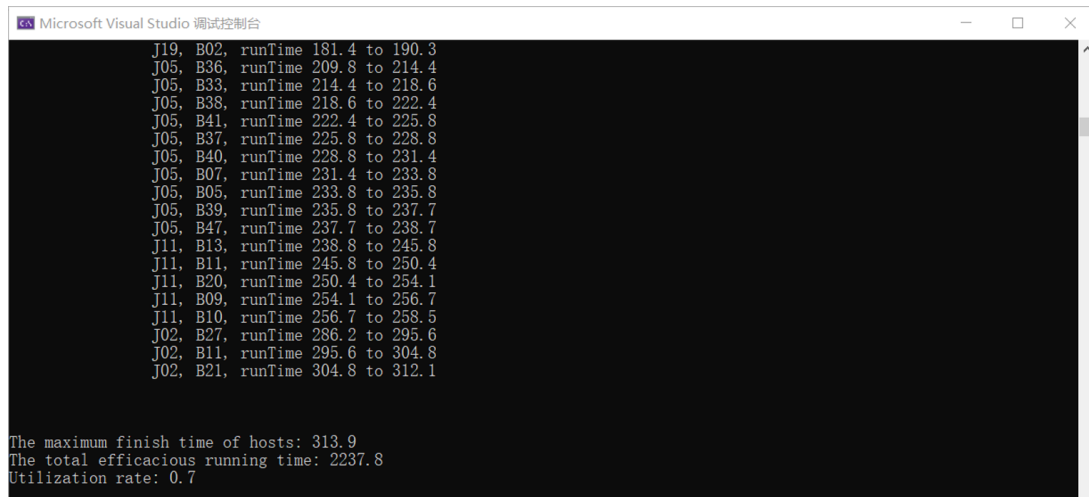
```

4 Results and Parameter Experiments

4.1 Final Results

We first give a sample result of our algorithm:

For *Task1*:



```

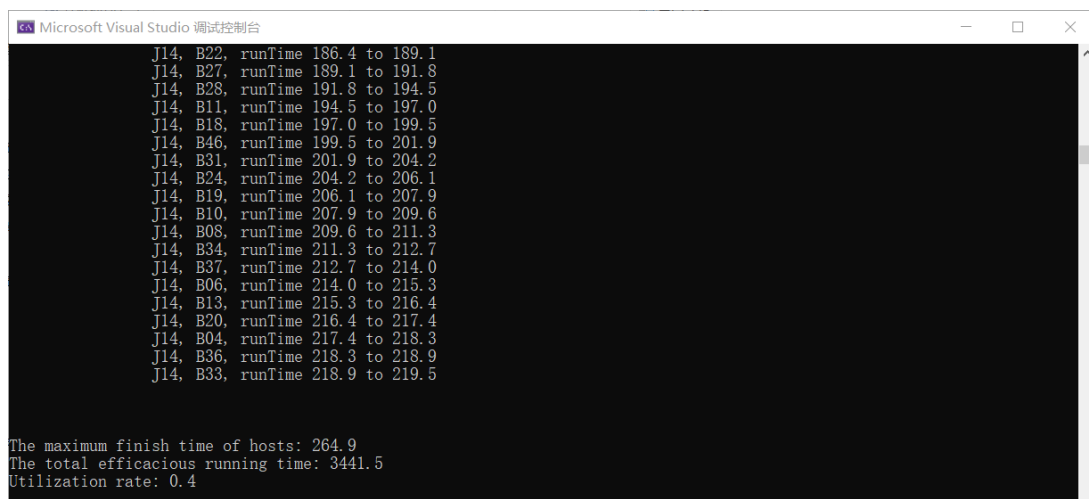
J19, B02, runTime 181.4 to 190.3
J05, B36, runTime 209.8 to 214.4
J05, B33, runTime 214.4 to 218.6
J05, B38, runTime 218.6 to 222.4
J05, B41, runTime 222.4 to 225.8
J05, B37, runTime 225.8 to 228.8
J05, B40, runTime 228.8 to 231.4
J05, B07, runTime 231.4 to 233.8
J05, B05, runTime 233.8 to 235.8
J05, B39, runTime 235.8 to 237.7
J05, B47, runTime 237.7 to 238.7
J11, B13, runTime 238.8 to 245.8
J11, B11, runTime 245.8 to 250.4
J11, B20, runTime 250.4 to 254.1
J11, B09, runTime 254.1 to 256.7
J11, B10, runTime 256.7 to 258.5
J02, B27, runTime 286.2 to 295.6
J02, B11, runTime 295.6 to 304.8
J02, B21, runTime 304.8 to 312.1

The maximum finish time of hosts: 313.9
The total efficacious running time: 2237.8
Utilization rate: 0.7

```

图 3. The sample result of Task1

For *Task2*:



```

J14, B22, runTime 186.4 to 189.1
J14, B27, runTime 189.1 to 191.8
J14, B28, runTime 191.8 to 194.5
J14, B11, runTime 194.5 to 197.0
J14, B18, runTime 197.0 to 199.5
J14, B46, runTime 199.5 to 201.9
J14, B31, runTime 201.9 to 204.2
J14, B24, runTime 204.2 to 206.1
J14, B19, runTime 206.1 to 207.9
J14, B10, runTime 207.9 to 209.6
J14, B08, runTime 209.6 to 211.3
J14, B34, runTime 211.3 to 212.7
J14, B37, runTime 212.7 to 214.0
J14, B06, runTime 214.0 to 215.3
J14, B13, runTime 215.3 to 216.4
J14, B20, runTime 216.4 to 217.4
J14, B04, runTime 217.4 to 218.3
J14, B36, runTime 218.3 to 218.9
J14, B33, runTime 218.9 to 219.5

The maximum finish time of hosts: 264.9
The total efficacious running time: 3441.5
Utilization rate: 0.4

```

图 4. The sample result of Task2

The complete output can be checked in `task1_case1_01.txt` and `task2_case1_01.txt`.

We can visualize the annealing process:

For *Task1*:

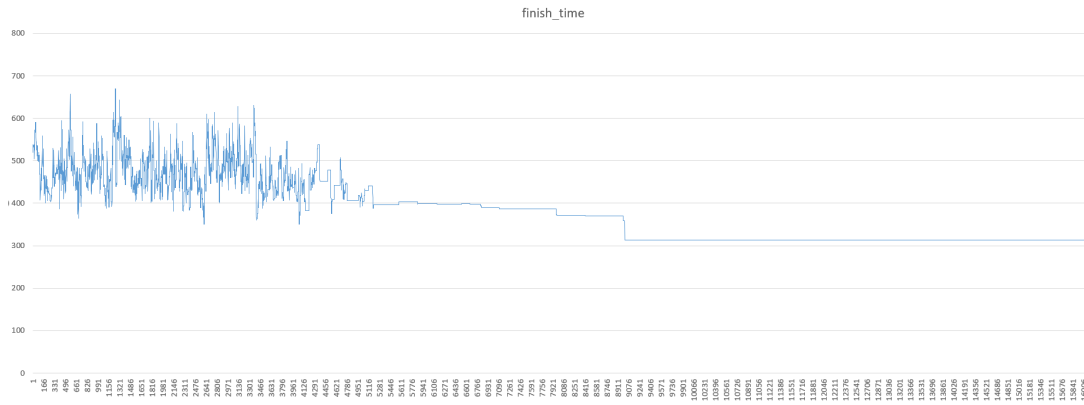


图 5. The annealing process of Task1

For *Task2*:

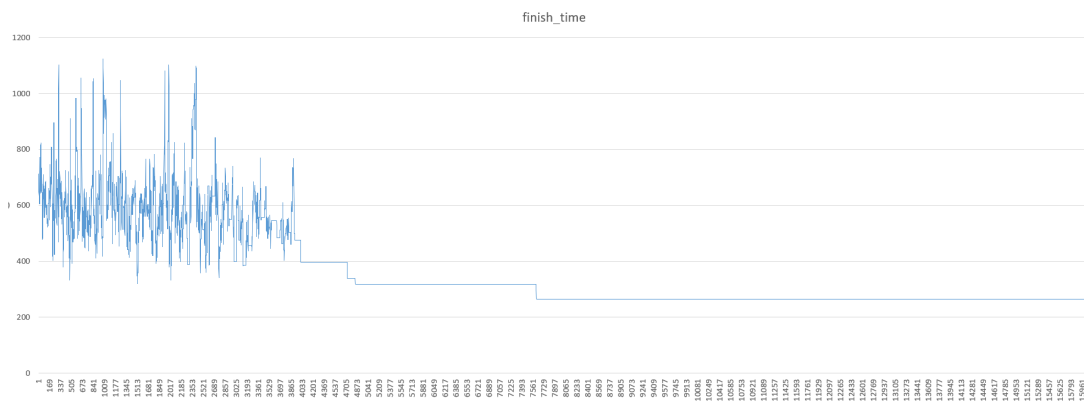


图 6. The annealing process of Task2

4.2 Parameter Experiments

In this section, we will use different parameters to test the simulated annealing algorithm and compare the different results.

For *Task1*:

表 1. The average result of Task1 for different parameters

	$P_{co} = 0.25$	$P_{co} = 0.5$	$P_{co} = 0.75$
$P_{rc} = 0.25$	361.10	356.70	353.16
$P_{rc} = 0.5$	360.54	358.16	349.98
$P_{rc} = 0.75$	344.74	354.98	354.90

For *Task2*:

表 2. The average result of Task2 for different parameters

	$P_{co} = 0.25$	$P_{co} = 0.5$	$P_{co} = 0.75$
$P_{rc} = 0.25$	352.64	350.56	354.18
$P_{rc} = 0.5$	332.76	315.92	331.48
$P_{rc} = 0.75$	316.10	324.28	323.02

However, the parameters will be affected by the test case itself. The results we get may not be universal.

5 Algorithm Analysis

In this section, we will analyze the time complexity of our algorithm. And we will justify whether this problem is NP-complete.

5.1 Time Complexity

1. The times of annealing is

$$\min(M, \log_{\beta} \frac{T_e}{T_0})$$

2. In *move()*, the time complexity of changing order is $O(1)$. As for choosing cores, the max number of cores allocated to one job is 10. Assume that we need X times to pick out 10 different cores for one job, X_i is how many times it needs to pick a new core when it has already picked $i - 1$ cores. X_i follows the a geometric distribution.

$$E(X_i) = \frac{1}{p_i} = \frac{m}{m - i + 1}$$

p_i means the probability of choosing a new core when it has already picked $i - 1$ cores.

Because $X = \sum_{i=1}^{10} X_i$,

$$E(X) = \sum_{i=1}^{10} \frac{m}{m-i+1} \rightarrow 10$$

3. The time complexity of *greedy_allocation*(job_i) is $O(jobBlock[i] \cdot jobCore[i])$. Because $jobCore[i] \leq jobBlock[i]$, the time complexity is actually $O(jobBlock[i]^2)$. Therefore, the time complexity of *move*() is $O(10 \cdot \max(jobBlock[])^2) \sim O(\max(jobBlock[])^2)$.
4. *finishTime*() essentially iterates through all the blocks and then finds the max finish time on cores. So the time complexity is $O(\sum_{i=0}^{n-1} jobBlock[i] + m)$.
5. The other functions just have the same time complexity of the functions above or smaller.

So the time complexity of the whole algorithm is

$$O(\min(M, \log_{\beta} \frac{T_e}{T_0}) \cdot [\max(jobBlock[])^2 + \sum_{i=0}^{n-1} jobBlock[i] + m])$$

5.2 Explanation to Greedy Thoughts

First of all, we would like to explain how we figure out the greedy algorithm used for shortening the length of cubes of jobs.

We do the following qualitative analysis:

1. The length of a job's cube must be larger than the max size of its blocks.
2. Every time we put a new block into the cores' queue, the increase of the length is not bigger than the size of this block.
3. We want the idle part to be as small as possible. And the max idle part is the gap between the shortest queue and the longest queue.
4. If adding a block makes the longest queue changes, we want the increase of the length of the cube is the smallest.

Based on them, we designed the *greedy_allocation*(). It obtain a good enough solution for each job in the simplified version. Because the processing time for each block is proportional to its size. But it isn't optimal. We will explain it in Section 5.3.

In the comprehensive version, the transmission time will affect the actual processing time of a job. Because some blocks have no need to transmit, while the others need. The actual processing time is still proportional to its size, but the ratio is different for each block. So the Greedy Algorithm can't work better than it does in the simplified version. However, because the transmission time of each job is one tenth of the processing time or less, the greedy algorithm can still give a good solution.

5.3 NP-complete Proof

NP If we know the minimum finish time and the list of places of all blocks, we use function *finishtime*() to certificate, which leads to the polynomial time complexity.

Partition Problem The partition problem, or number partitioning, is the task of deciding whether a given multiset S of positive integers can be partitioned into two subsets S_1 and S_2 such that the sum of the numbers in S_1 equals the sum of the numbers in S_2 . The partition problem is NP-complete[4].

NP-complete of Simplest Case We consider the simplest case that there is only one job allocated onto 2 cores.

Consider the partition problem about the set of blocks' working time A . If we find a solution of the problem that there is no idle time ($finishtime = \frac{1}{2} \sum_{a \in A} a$), we get a solution of partition problem about A , and if there is idle time existing ($finishtime \neq \frac{1}{2} \sum_{a \in A} a$), we know that there is no solution of this partition problem.

The reduction takes $O(n)$ time complexity to sum all elements in A and $O(1)$ time of the "black box" algorithm, so it is a polynomial reduction, which means the simplest case is NP-complete.

NP-complete of General Case It is obvious that the general case, with more jobs and more blocks, is harder than the simplest case. The discussion below is based on the assumption that we have an algorithm $f(m, n)$ to solve m job case on n cores, where m and n are arbitrary. Single-host or multi-host doesn't matter, because we only need to add a special parameter on the size of blocks.

Considering the simplest case, we can set $m = 1$ and $n = 2$, then $f(1, 2)$ is an algorithm to solve it, which leads to a polynomial reduction. Since the simplest case is NP-complete, the general case is NP-complete as well.

5.4 Experiment on Greedy Quality

In Section 5.3, we prove that minimizing the length of one job cube with fixed core number is an NP-complete problem. If $P \neq NP$, we can't find a optimal solution in polynomial time. The greedy algorithm we used for the simplified version is not optimal. (For the comprehensive version it is obviously not the optimal.) Therefore, in the simplified version, we compare the solution which our greedy algorithm gives to the theoretical optimal solution, in order to test whether it gives a good enough solution. The theoretical optimal cube length for job_i is

$$\frac{\sum_{k=0}^{jobBlock[i]-1} datasize[i][k]}{jobCore[i]}$$

It is easy to know that the actual optimal cube length must not be shorter than the theoretical one.

```

Microsoft Visual Studio 调试控制台

The maximum finish time of hosts: 340.8
The total efficacious running time: 1884.8
Utilization rate: 0.6

For job0, the cube length greedy algorithm gives is 1765.0. The theoretical shortest length is 1765.0.
The ratio is 1.0.
For job1, the cube length greedy algorithm gives is 945.0. The theoretical shortest length is 705.0.
The ratio is 1.3.
For job2, the cube length greedy algorithm gives is 1662.0. The theoretical shortest length is 1246.5.
The ratio is 1.3.
For job3, the cube length greedy algorithm gives is 1679.0. The theoretical shortest length is 1004.8.
The ratio is 1.7.
For job4, the cube length greedy algorithm gives is 1066.0. The theoretical shortest length is 636.2.
The ratio is 1.7.
For job5, the cube length greedy algorithm gives is 2914.0. The theoretical shortest length is 1943.0.
The ratio is 1.5.
For job6, the cube length greedy algorithm gives is 574.0. The theoretical shortest length is 338.6.
The ratio is 1.7.
For job7, the cube length greedy algorithm gives is 2044.0. The theoretical shortest length is 1200.0.
The ratio is 1.7.
For job8, the cube length greedy algorithm gives is 2334.0. The theoretical shortest length is 1402.2.
The ratio is 1.7.
For job9, the cube length greedy algorithm gives is 1225.0. The theoretical shortest length is 818.0.
The ratio is 1.5.
For job10, the cube length greedy algorithm gives is 2838.0. The theoretical shortest length is 2838.0.
The ratio is 1.0.
For job11, the cube length greedy algorithm gives is 1591.0. The theoretical shortest length is 1067.3.
The ratio is 1.5.
For job12, the cube length greedy algorithm gives is 1280.0. The theoretical shortest length is 741.7.
The ratio is 1.7.
For job13, the cube length greedy algorithm gives is 2417.0. The theoretical shortest length is 1511.0.
The ratio is 1.6.
For job14, the cube length greedy algorithm gives is 677.0. The theoretical shortest length is 397.4.
The ratio is 1.7.
For job15, the cube length greedy algorithm gives is 881.0. The theoretical shortest length is 554.2.
The ratio is 1.6.
For job16, the cube length greedy algorithm gives is 2181.0. The theoretical shortest length is 1272.7.
The ratio is 1.7.
For job17, the cube length greedy algorithm gives is 783.0. The theoretical shortest length is 518.0.
The ratio is 1.5.
For job18, the cube length greedy algorithm gives is 1829.0. The theoretical shortest length is 1049.9.
The ratio is 1.7.
For job19, the cube length greedy algorithm gives is 476.0. The theoretical shortest length is 359.5.
The ratio is 1.3.

```

图7. Comparison between greedy solution and theoretical optimal solution

The complete output can be checked in `task1_case1_01 (2).txt`.

We can see the cube length greedy algorithm gives is not more than 2 times the theoretical optimal one. For the actual optimal solution this ratio may be smaller. And sometimes it gives a result as good as the theoretical optimal one. Therefore, we think our greedy algorithm gives a good enough solution.

6 Experiment Based on A Completely Random Move

In Section 5, we explain and test the quality of the greedy algorithm we used. And we can know for multiple host version greedy algorithm is not optimal. So we design another move method which not only moves the job cubes, but also moves the blocks in each job cube. Namely, we will use a completely random strategy to generate new states and check if it can give us a better result. The new move strategy is following:

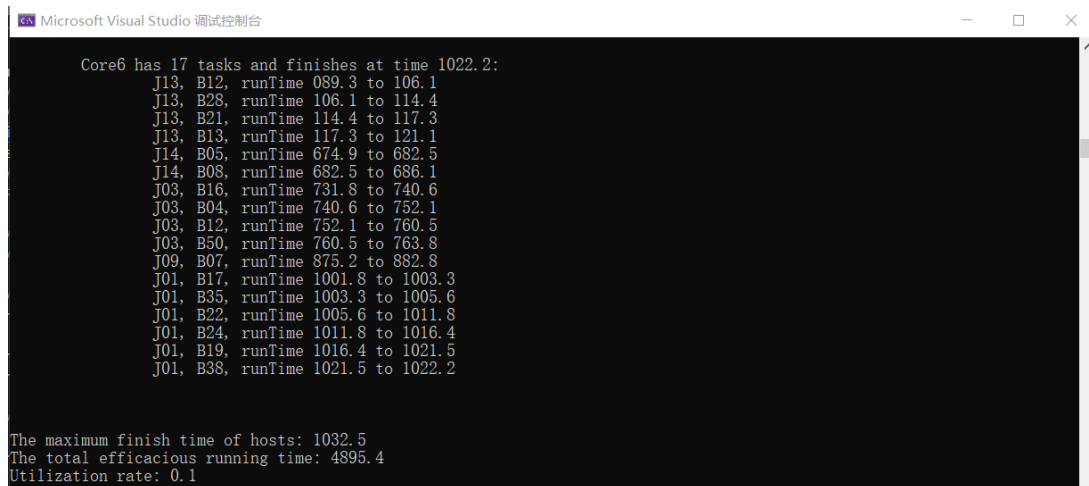
1. Randomly choose two adjacent jobs, change their order. This will happen at a probability of P_{co} . The probability of each job that will be chosen is equal.
2. Randomly choose one job, randomly choose a core that has been allocated to it, change it to another core that hasn't been allocated to it. This will happen at a probability of P_{rc} . The probability of each core that will be chosen is equal.

3. Randomly choose one job job_i , randomly choose a core $core_j$ that has been allocated to job_i , randomly choose a block B_k that on this core. Then randomly generate a core index. The core index will be in the range of $[0, temp_jobCore[i]]$, but it can't be the $core_j$. This will happen at a probability of P_{mb} .

If the core index is equal to $temp_jobCore[i]$, it means we will give a new core to job_i . We randomly choose this new core from the cores which haven't been allocated to job_i . If there is only one block on $core_j$, we will directly change $core_j$ to this new core. If there are more than one block on $core_j$, we add this new core to job_i and move B_k from $core_j$ to this new core.

If the core index is smaller than $temp_jobCore[i]$, it means we will move B_k from $core_j$ to another core of job_i . If there is only one block on $core_j$, we will delete $core_j$ from job_i after moving.

However, in experiment, though we amplify the probability of each move strategy and the iteration times, we can't get a better result than that the greedy algorithm gives.



```

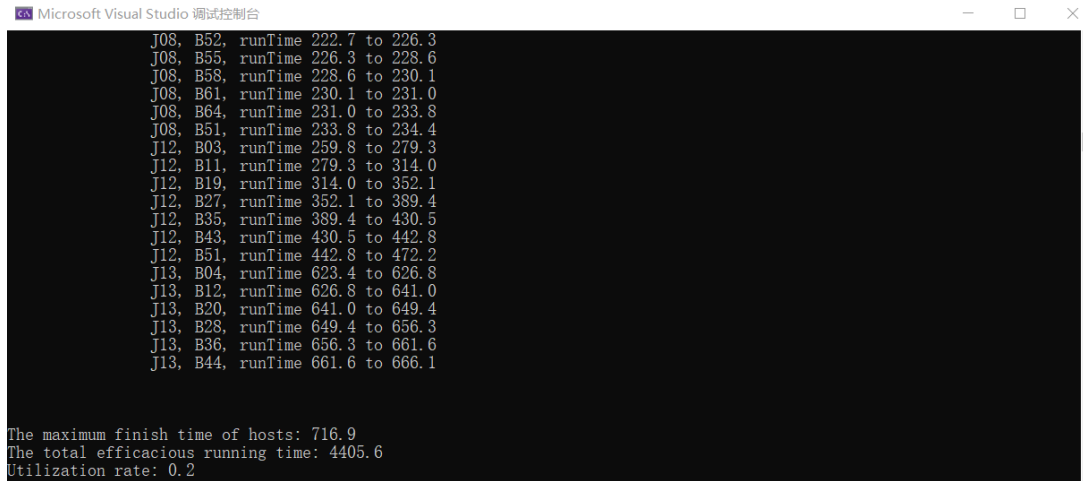
Core6 has 17 tasks and finishes at time 1022.2:
J13, B12, runTime 089.3 to 106.1
J13, B28, runTime 106.1 to 114.4
J13, B21, runTime 114.4 to 117.3
J13, B13, runTime 117.3 to 121.1
J14, B05, runTime 674.9 to 682.5
J14, B08, runTime 682.5 to 686.1
J03, B16, runTime 731.8 to 740.6
J03, B04, runTime 740.6 to 752.1
J03, B12, runTime 752.1 to 760.5
J03, B50, runTime 760.5 to 763.8
J09, B07, runTime 875.2 to 882.8
J01, B17, runTime 1001.8 to 1003.3
J01, B35, runTime 1003.3 to 1005.6
J01, B22, runTime 1005.6 to 1011.8
J01, B24, runTime 1011.8 to 1016.4
J01, B19, runTime 1016.4 to 1021.5
J01, B38, runTime 1021.5 to 1022.2

The maximum finish time of hosts: 1032.5
The total efficacious running time: 4895.4
Utilization rate: 0.1

```

图8. The result of Task2 given by completely random move strategy

This is because each movement of this strategy is too small. If we change the Simulated Annealing Algorithm, let it only accepts the state with lower energy, the result will be improved a little. But it is still much worse than the greedy.



```

Microsoft Visual Studio 调试控制台
J08, B52, runTime 222.7 to 226.3
J08, B55, runTime 226.3 to 228.6
J08, B58, runTime 228.6 to 230.1
J08, B61, runTime 230.1 to 231.0
J08, B64, runTime 231.0 to 233.8
J08, B51, runTime 233.8 to 234.4
J12, B03, runTime 259.8 to 279.3
J12, B11, runTime 279.3 to 314.0
J12, B19, runTime 314.0 to 352.1
J12, B27, runTime 352.1 to 389.4
J12, B35, runTime 389.4 to 430.5
J12, B43, runTime 430.5 to 442.8
J12, B51, runTime 442.8 to 472.2
J13, B04, runTime 623.4 to 626.8
J13, B12, runTime 626.8 to 641.0
J13, B20, runTime 641.0 to 649.4
J13, B28, runTime 649.4 to 656.3
J13, B36, runTime 656.3 to 661.6
J13, B44, runTime 661.6 to 666.1

The maximum finish time of hosts: 716.9
The total efficacious running time: 4405.6
Utilization rate: 0.2

```

图9. The result of Task2 given by completely random move strategy if SAA only accepts states with lower energy

If we apply this changed SAA on the greedy version, sometimes it can give a better solution than the average result of the original SAA. But it can hardly give a solution better than the best result the original SAA ever gives. The reason is that it may cause the state to stay at the locally optimal solution. When the solution is good enough, it is hard for the changed SAA to move anymore.

7 Summary

In this article, we focus on designing an algorithm for solving resource scheduling problems in Hadoop. The algorithm we design is based on simulated annealing algorithm and greedy algorithm to give an approximate solution to the problems of both simplified and comprehensive versions. We also do some experiments on parameters in the algorithm but the results we get are for the special test cases and not universal. During analyzing the complexity and other quality of our algorithm, we found that the problem is NP-complete. Though we hope our greedy algorithm is optimal for one job in the simplified version, it is an NP-complete problem, we can't find its optimal solution in polynomial time. So we compare the greedy solution to the theoretical optimal solution to prove the greedy one is good enough. Finally, we use a completely random move strategy to try to find a better solution for the comprehensive version. But the result is that it can't work better than the greedy.

In general, the SACG is our feasible but imperfect answer to resource scheduling problem in Hadoop.

Acknowledgements

This problem is motivated from a real-world cloud service corporation. It is formulated by Prof. Xiaofeng Gao (gao-xf@cs.sjtu.edu.cn) from Department of Computer Science and Engineering at Shanghai Jiao Tong University. Yue Ma scribed and modified the project. Yangyang Bao provided the example in Task 1. Zhen Sun provided the example in Task 2. Wanghua Shi provided the

test data and source code. Jiale Zhang, Tianyao Shi helped on proofreading and provided many suggestions.

Thank assistants Fang Hongjie, Lv Jialin and Shi Wanghua. They gave so many instructive suggestions to our group.

Thank professor Gao Xiaofeng. The problem she assigned helps us understand the actual case in real world better and improves our abilities.

Every group member has learned a lot from this project.

Zhenran Xiao: Figuring out this problem is more and more interesting during the whole process. At the beginning, I just feel anxious and nervous. I was worried whether we can give a result. But as the project progressed, I had fun, especially when I drew the picture of the annealing process. I feel the charm of the Simulated Annealing Algorithm.

When I gave out the greedy algorithm used for one job on single host(I need to state I think out this by myself instead of searching it out on the Internet), I actually don't know how to prove whether it is optimal. Finally, we found out that the problem is an NP-complete problem. We can't find a optimal solution in polynomial time.

Really thankful to my groupmates, they did much difficult work in the project and answered many questions I asked. Trying to express everyone's awesome ideas clearly in the article is a challenge as well. I really learned a lot.

Zichun Ye: As a member of the group, my job is implementing the main part of the code, debugging and providing the test answer. During this process, I faced many difficulties including setting the variables, finding errors through hundreds of lines and reading others' code. I took notes of all these variables to make coding easier. In the end, the program runs smoothly without annoying bugs.

I learned that before writing the code, I shall have clearheaded logic, otherwise it's destructive to modify the code. Writing notes in the code is also essential because your teammates also need to understand what you write.

Thank Xiao Zhenran and Long Mabiao. They supported me a lot and did an excellent job in this project.

Mabiao Long: Since we haven't learned about the NP problem at first, we learned from other similar problems on the Internet and finally chose the SA algorithm, instead of dynamic programming, to deal with the case that a large amount of data may appear. I believe that learning from others will be a shortcut while facing engineering problems, though it may be a hard problem or even an unsolvable problem now, there may be some tips or thoughts that lead to the correct way.

I worked on two functions, *finishtime()* and *write_back()*. Based on a concise and effective data structure, it is easier to work out. I also worked on the proof of NP-complete, which deepened my understanding of computational complexity and taught me that some seemingly simple questions can be difficult.

Thanks for my teammates, Zhenran Xiao and Zichun Ye, they inspired me and corrected my mistakes, which taught me a lot during these weeks.

参考文献

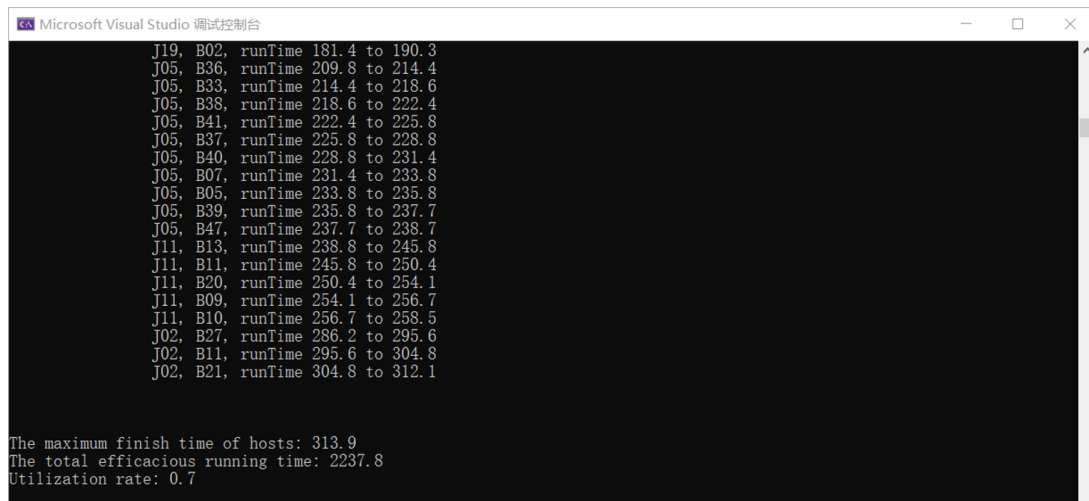
1. shiwanghua/SharedFiles, <https://github.com/shiwanghua/SharedFiles/tree/Project-CodeDemo>. Last accessed 15 May 2022
2. yunzhe99/Hadoop_Scheduling, https://github.com/yunzhe99/Hadoop_Scheduling. Last accessed 11 Dec 2021
3. 模拟退火 (simulated annealing) 算法详解, https://blog.csdn.net/weixin_60737527/article/details/123936955?utm_medium=distribute.pc_relevant.none-task-blog-2~default~baidujs_baidulandingword-default-0.pc_relevant_default&spm=1001.2101.3001.4242.1&utm_relevant_index=3. Last accessed 3 Apr 2022
4. Hayes, Brian (March–April 2002), "The Easiest Hard Problem" (PDF), American Scientist, Sigma Xi, The Scientific Research Society, vol. 90, no. 2, pp. 113–117, JSTOR 27857621

Appendix A

表 3. Symbols and Definitions

Symbols	Definitions
n	The number of jobs
m	The number of cores
q	The number of hosts
job_i, J	job_i is the i -th job. The job set is $J = \{job_0, \dots, job_{n-1}\}$.
h_l, H	h_l is the l -th host. The host set is $H = \{h_0, \dots, h_{q-1}\}$.
m_l	The number of cores on host h_l
c_j^l, C_l	c_j^l is the j -th core on host h_l . C_l is the set of cores on host h_l .
C	The set of cores. $C = \{c_0, \dots, c_{m-1}\}$ for single-host. $C = \cup_{l=0}^{q-1} C_l$ for multi-host.
b_k^i	The block of job_i whose id is k
B_j^i	The set of data blocks of job_i allocated to core c_j
B^i	The set of data blocks of job_i
B_{lj}^i	The set of data blocks of job_i allocated to core c_j^l
\tilde{B}_{lj}^i	The set of data blocks of job_i allocated to core c_j^l but not initially stored on h_l
$size(\cdot)$	The size function of data block
$g(\cdot)$	The computing decaying coefficient caused by multi-core effect
s_i	The computing speed of job_i by a single core
s_t	The transmission speed of data
e_i	The number of cores processing job_i
t_i	The time to start processing job_i
tp_j^i, tf_j^i	The processing time / finishing time of core c_j for job_i
tp_{lj}^i, tf_{lj}^i	The processing time / finishing time of core c_j^l for job_i
$tf(job_i)$	The finishing time of job_i
$state_i$	The i th state of the material or solution
$E(i)$	The energy of $state_i$
K	Boltzmann's constant
T	The temperature of current state
T_i	The temperature of $state_i$
T_e	The final temperature
β	The rate of temperature fall
M	The maximum number of iterations
δE	The change of energy between two states
P_{co}	The probability changing two job's order
P_{rc}	The probability re-choosing cores for one job
X	Times to pick out 10 different cores for one job
X_i	Times it needs to pick a new core when it has already picked $i - 1$ cores
p_i	The probability of choosing a new core when it has already picked $i - 1$ cores
P_{mb}	The probability move blocks for one job

Appendix B



```

Microsoft Visual Studio 调试控制台

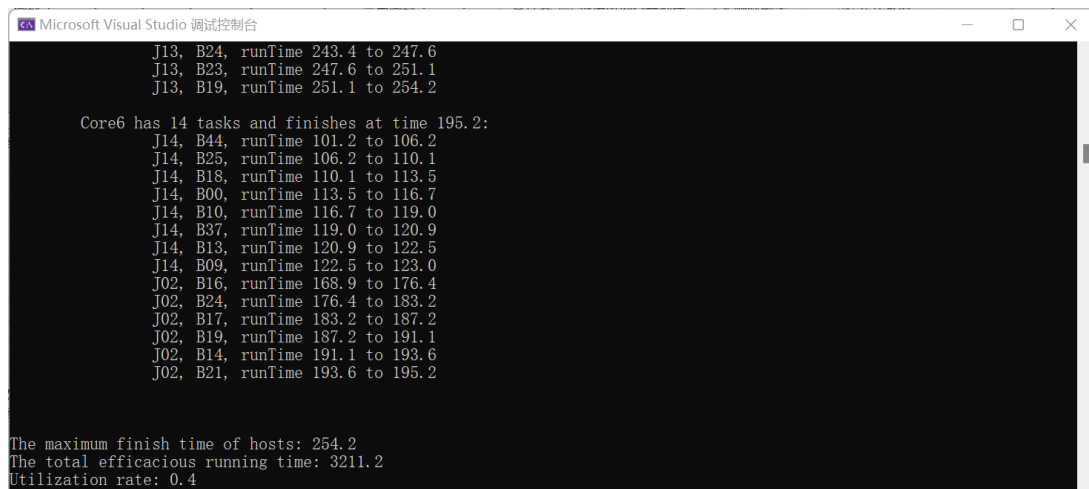
J19, B02, runTime 181.4 to 190.3
J05, B36, runTime 209.8 to 214.4
J05, B33, runTime 214.4 to 218.6
J05, B38, runTime 218.6 to 222.4
J05, B41, runTime 222.4 to 225.8
J05, B37, runTime 225.8 to 228.8
J05, B40, runTime 228.8 to 231.4
J05, B07, runTime 231.4 to 233.8
J05, B05, runTime 233.8 to 235.8
J05, B39, runTime 235.8 to 237.7
J05, B47, runTime 237.7 to 238.7
J11, B13, runTime 238.8 to 245.8
J11, B11, runTime 245.8 to 250.4
J11, B20, runTime 250.4 to 254.1
J11, B09, runTime 254.1 to 256.7
J11, B10, runTime 256.7 to 258.5
J02, B27, runTime 286.2 to 295.6
J02, B11, runTime 295.6 to 304.8
J02, B21, runTime 304.8 to 312.1

The maximum finish time of hosts: 313.9
The total efficacious running time: 2237.8
Utilization rate: 0.7

```

图 10. The best result of Task1 we have ever gotten

The complete output can be checked in `task1_case1_01.txt`.



```

Microsoft Visual Studio 调试控制台

J13, B24, runTime 243.4 to 247.6
J13, B23, runTime 247.6 to 251.1
J13, B19, runTime 251.1 to 254.2

Core6 has 14 tasks and finishes at time 195.2:
J14, B44, runTime 101.2 to 106.2
J14, B25, runTime 106.2 to 110.1
J14, B18, runTime 110.1 to 113.5
J14, B00, runTime 113.5 to 116.7
J14, B10, runTime 116.7 to 119.0
J14, B37, runTime 119.0 to 120.9
J14, B13, runTime 120.9 to 122.5
J14, B09, runTime 122.5 to 123.0
J02, B16, runTime 168.9 to 176.4
J02, B24, runTime 176.4 to 183.2
J02, B17, runTime 183.2 to 187.2
J02, B19, runTime 187.2 to 191.1
J02, B14, runTime 191.1 to 193.6
J02, B21, runTime 193.6 to 195.2

The maximum finish time of hosts: 254.2
The total efficacious running time: 3211.2
Utilization rate: 0.4

```

图 11. The best result of Task2 we have ever gotten

The complete output can be checked in `task2_case1_01 (2).txt`.