

# 基于 Spark 的 pagerank 算法实现与原理解析

---

刘佳玮，计算机科学与技术学院，20031211496

<https://github.com/muyuuuu/Spark-learn>

2021 年 3 月 24 日



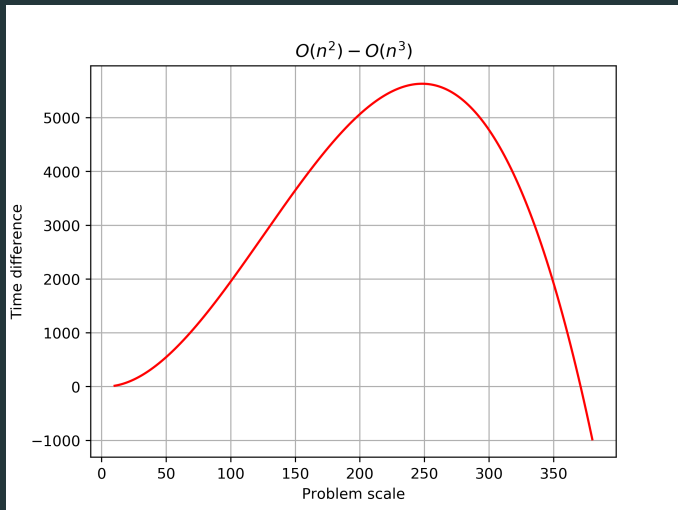
## 好的算法与好的设备

---

# 算法时间复杂度

同一个问题，一个时间复杂度  $O(n^2)$  与  $O(n^3)$  的时间对比，代码开放于：

[https://github.com/muyuuuu/Algorithm/tree/master/Insert\\_sort](https://github.com/muyuuuu/Algorithm/tree/master/Insert_sort)

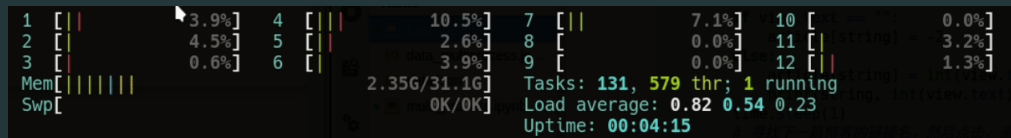


## 发挥设备优势

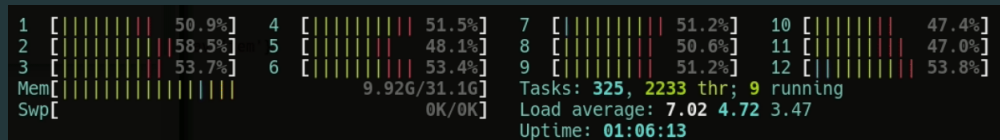
---

## 发挥设备优势

一个耗时 1153 秒的单进程任务：



使用多进程改进，相同任务耗时 105 秒，且多核利用率较为均衡：



任务必须可以并行化。代码地址：

<https://muyuuuu.github.io/2020/03/18/multi-process/>

# 硬件与软件依赖

---

# 硬件与软件

## 硬件部分

**CPU** 2 个 Intel(R) Xeon(R) Gold 5115 CPU 2.40GHz, 10 核心 20 线程

**GPU** 4 路 Tesla P40, 每路显存容量 22GB

**内存** 128GB

**外存** 520TB 可用, 已用 15TB

## 软件部分

**系统** CentOS Linux release 7.3.1611 (执行), Arch 5.9.6(开发)

**python** 3.8.2, 开发语言

**pytorch** 1.6.0, 模型实现, 借助其提供的 API 实现并行

**ssh** OpenSSH\_8.3p1, OpenSSL 1.1.1h: 实现远程登录

**scp** 文件传输



## 命令行内执行

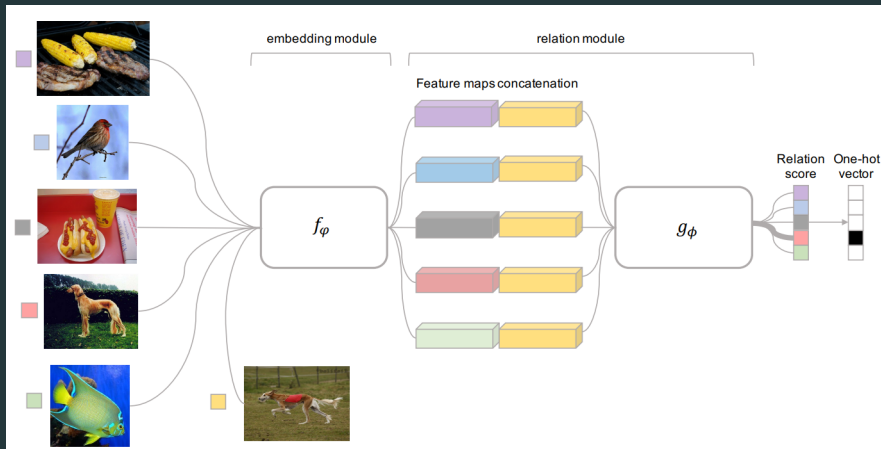
1. `mv`, `cd`, `ls`, `cp`, `cat` 等文件操作
2. `nohup python train.py > log` 挂起运行与重定向输出
3. `ps -f|grep python` 查看挂起程序是否执行

# 模型

---

# 模型结构

实现的模型为 Relation Network<sup>1</sup>。数据集为 miniImageNet<sup>2</sup>。



<sup>1</sup><https://ieeexplore.ieee.org/abstract/document/8778601>

<sup>2</sup><https://drive.google.com/file/d/0B3Irx3uQNoBMQ1F1NXJsZUdYWEE/view>

## 实验结果

---

# DataParallel 单机多卡

■ `nvidia-smi` 查看显卡利用率:

NVIDIA-SMI 396.26				Driver Version: 396.26			
GPU	Name	Persistence-M	Bus-Id	Disp.A	Volatile	Uncorr. ECC	
Fan	Temp	Perf	Pwr:Usage/Cap	Memory-Usage	GPU-Util	Compute M.	
0	Tesla P40	On	00000000:3B:00.0	Off		0	
N/A	55C	P0	157W / 250W	22825MiB / 22919MiB	99%	Default	
1	Tesla P40	On	00000000:86:00.0	Off		0	
N/A	58C	P0	185W / 250W	22015MiB / 22919MiB	91%	Default	
2	Tesla P40	On	00000000:AF:00.0	Off		0	
N/A	37C	P0	138W / 250W	11736MiB / 22919MiB	51%	Default	
3	Tesla P40	On	00000000:D8:00.0	Off		0	
N/A	33C	P0	151W / 250W	6451MiB / 22919MiB	28%	Default	

程序执行时间:  $T_1 = 137172$  秒, 约 2286 分钟, 约 1.59 天。

# DistributedDataParallel 单机多卡

■ nvidia-smi 查看显卡利用率:

```
文件(F) 编辑(E) 视图(V) 书签(B) 设置(S) 帮助(H)
Wed Nov 18 21:56:15 2020
```

NVIDIA-SMI 396.26				Driver Version: 396.26			
GPU	Name	Persistence-M	Bus-Id	Disp.A	Volatile Uncorr. ECC		
Fan	Temp	Perf	Pwr:Usage/Cap	Memory-Usage	GPU-Util	Compute M.	
0	Tesla P40	On	00000000:3B:00.0	Off	0		
N/A	63C	P0	146W / 250W	5891MiB / 22919MiB	99%	Default	
1	Tesla P40	On	00000000:86:00.0	Off	0		
N/A	57C	P0	94W / 250W	22839MiB / 22919MiB	98%	Default	
2	Tesla P40	On	00000000:AF:00.0	Off	0		
N/A	52C	P0	53W / 250W	22449MiB / 22919MiB	95%	Default	
3	Tesla P40	On	00000000:D8:00.0	Off	0		
N/A	55C	P0	96W / 250W	22239MiB / 22919MiB	89%	Default	

程序执行时间:  $T_2 = 89856$  秒, 约 1498 分钟, 约 1.03 天。

■ 加速比:  $\frac{T_1}{T_2} = 1.53$

■ 准确率对比: Dataparallel: 0.566, DDP: 0.582。

■ 代码开放于: <https://github.com/muyuuuu/Algorithm/tree/master/meta-learning/Metric-based/Relation-Netowrk>

**感谢聆听！**